

NASA/TP-2018-219822



# Navigation Filter Best Practices

*Edited by*

*J. Russell Carpenter*

*Goddard Space Flight Center, Greenbelt, Maryland*

*Christopher N. D'Souza*

*Johnson Space Center, Houston, Texas*

---

April 2018

## NASA STI Program... in Profile

Since its founding, NASA has been dedicated to the advancement of aeronautics and space science. The NASA scientific and technical information (STI) program plays a key part in helping NASA maintain this important role.

The NASA STI Program operates under the auspices of the Agency Chief Information Officer. It collects, organizes, provides for archiving, and disseminates NASA's STI. The NASA STI Program provides access to the NASA Aeronautics and Space Database and its public interface, the NASA Technical Report Server, thus providing one of the largest collection of aeronautical and space science STI in the world. Results are published in both non-NASA channels and by NASA in the NASA STI Report Series, which includes the following report types:

- **TECHNICAL PUBLICATION.** Reports of completed research or a major significant phase of research that present the results of NASA programs and include extensive data or theoretical analysis. Includes compilations of significant scientific and technical data and information deemed to be of continuing reference value. NASA counterpart of peer-reviewed formal professional papers, but having less stringent limitations on manuscript length and extent of graphic presentations.
- **TECHNICAL MEMORANDUM.** Scientific and technical findings that are preliminary or of specialized interest, e.g., quick release reports, working papers, and bibliographies that contain minimal annotation. Does not contain extensive analysis.
- **CONTRACTOR REPORT.** Scientific and technical findings by NASA-sponsored contractors and grantees.

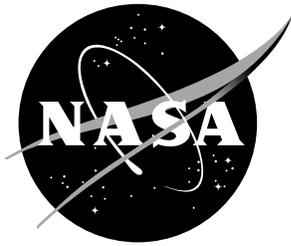
- **CONFERENCE PUBLICATION.** Collected papers from scientific and technical conferences, symposia, seminars, or other meetings sponsored or co-sponsored by NASA.
- **SPECIAL PUBLICATION.** Scientific, technical, or historical information from NASA programs, projects, and missions, often concerned with subjects having substantial public interest.
- **TECHNICAL TRANSLATION.** English-language translations of foreign scientific and technical material pertinent to NASA's mission.

Specialized services also include organizing and publishing research results, distributing specialized research announcements and feeds, providing information desk and personal search support, and enabling data exchange services.

For more information about the NASA STI Program, see the following:

- Access the NASA STI program home page at <http://www.sti.nasa.gov>
- E-mail your question to [help@sti.nasa.gov](mailto:help@sti.nasa.gov)
- Phone the NASA STI Information Desk at 757-864-9658
- Write to:  
NASA STI Information Desk  
Mail Stop 148  
NASA Langley Research Center  
Hampton, VA 23681-2199

NASA/TP-2018-219822



# Navigation Filter Best Practices

*Edited by*

*J. Russell Carpenter*

*Goddard Space Flight Center, Greenbelt, Maryland*

*Christopher N. D'Souza*

*Johnson Space Center, Houston, Texas*

National Aeronautics and  
Space Administration

NASA Engineering and Safety Center  
Hampton, Virginia 23681

---

April 2018

The use of trademarks or names of manufacturers in this report is for accurate reporting and does not constitute an official endorsement, either expressed or implied, of such products or manufacturers by the National Aeronautics and Space Administration.

Available from:

NASA STI Program / Mail Stop 148  
NASA Langley Research Center  
Hampton, VA 23681-2199  
Fax: 757-864-6500



NASA ENGINEERING AND SAFETY CENTER

# Navigation Filter Best Practices

First Edition

*J. Russell Carpenter and Christopher N. D'Souza, Editors*

with contributions from  
J. Russell Carpenter, Christopher N. D'Souza, F. Landis Markley, Renato Zanetti

April, 2018

This material is declared a work of the US Government and is not subject to copyright protection in the United States, but may be subject to US Government copyright in other countries.

If there be two subsequent events, the probability of the second  $b/N$  and the probability of both together  $P/N$ , and it being first discovered that the second event has also happened, from hence I guess that the first event has also happened, the probability I am right is  $P/b$ .

Thomas Bayes, c. 1760

The explicit calculation of the optimal estimate as a function of the observed variables is, in general, impossible.

Rudolph Kalman, 1960.

The use of Kalman Filtering techniques in the on-board navigation systems for the Apollo Command Module and the Apollo Lunar Excursion Module was an important factor in the overwhelming success of the Lunar Landing Program.

Peter Kachmar, 2002.

Dedicated to the memory of Gene Muller, Emil Schiesser, and Bill Lear.

# Contents

Foreword	iii
Editor's Preface	v
Notational Conventions	ix
Chapter 1. The Extended Kalman Filter	1
1.1. The Additive Extended Kalman Filter	1
1.2. The Multiplicative Extended Kalman Filter	7
Chapter 2. The Covariance Matrix	9
2.1. Metrics for Orbit State Covariances	9
2.2. Covariance Propagation	15
2.3. Covariance Measurement Update	25
2.4. Sigma-Point Methods	28
Chapter 3. Processing Measurements	29
3.1. Measurement Latency	29
3.2. Invariance to the Order of Measurement Processing	30
3.3. Processing Vector Measurements	31
3.4. Filter Cascades	32
3.5. Use of Data from Inertial Sensors	32
Chapter 4. Measurement Underweighting	33
4.1. Introduction	33
4.2. Nonlinear Effects and the Need for Underweighting	34
4.3. Underweighting Measurements	38
4.4. Pre-Flight Tuning Aids	39
Chapter 5. Bias Modeling	41
5.1. Zero-Input Bias State Models	41
5.2. Single-Input Bias State Models	43
5.3. Multi-Input Bias State Models	49
Chapter 6. State Representations	53
6.1. Selection of Solve-For State Variables for Estimation	53
6.2. Units and Precision	53
6.3. Coordinate and Time Systems	54
6.4. Orbit Parameterizations	54
6.5. Relative State Representations	55
6.6. Modeling Inertial Components	58

Chapter 7. Factorization Methods	63
7.1. Why Use the UDU Factorization?	63
7.2. Preliminaries	64
7.3. The Time Update of the Covariance	65
7.4. The Measurement Update	72
7.5. Consider Covariance and Its Implementation in the UDU Filter	74
7.6. Conclusions	76
Chapter 8. Attitude Estimation	77
8.1. Attitude Matrix Representation	77
8.2. Euler Axis/Angle Representation	79
8.3. Quaternion Representation	80
8.4. Rodrigues Parameter Representation	83
8.5. Modified Rodrigues Parameters	84
8.6. Rotation Vector Representation	85
8.7. Euler Angles	86
8.8. Additive EKF (AEKF)	87
8.9. Multiplicative EKF (MEKF)	88
Chapter 9. Usability Considerations	93
9.1. Editing	93
9.2. Reinitialization, Restarts, and Backup Ephemeris	93
9.3. Ground System Considerations	94
Chapter 10. Smoothing	95
Chapter 11. Advanced Estimation Algorithms	99
The Sigma-Point Estimator	99
Appendix A. Models and Realizations of Random Variables	105
Appendix B. The Mathematics Behind the UDU Factorization	107
B.1. The Partitioning into Two Subproblems	107
B.2. The Mathematics Behind the Second Subproblem	107
B.3. The Agee-Turner Rank-One Update	109
B.4. Decorrelating Measurements	112
B.5. The Carlson Rank-One Update	112
Appendix C. An Analysis of Dual Inertial-Absolute and Inertial-Relative Navigation Filters	117
C.1. Introduction	117
C.2. The Filter Dynamics	117
C.3. Incorporation of Measurements	120
C.4. Analysis of the Merits of the Inertial-Absolute and Inertial-Relative Filters	123
C.5. Conclusions	125
Bibliography	127

## Foreword

It certainly should not come as a surprise to the reader that navigation systems are at the heart of almost all of NASA's missions, either on our launch vehicles, on robotic science spacecraft, or on our crewed human exploration vehicles. Clearly navigation is absolutely fundamental to operating our space systems across the wide spectrum of mission regimes. Safe and reliably performing navigation systems are essential elements needed for routine low Earth orbiting science missions, for rendezvous and proximity operation missions or precision formation flying missions (where relative navigation is a necessity), for navigation through the outer planets of the solar system, and for accomplishing pinpoint landing on planets/small bodies, and many more mission types.

I believe the reader will find that the navigation filter best practices the team has collected, documented, and shared in this first edition book will be of practical value in your work designing, developing, and operating modern navigation systems for NASA's challenging future missions. I want to thank the entire team that has diligently worked to create this NASA Engineering and Safety Center (NESC) GN&C knowledge capture report. I especially want to acknowledge the dedication, care, and attention to detail as well as the energy that both Russell Carpenter and Chris D'Souza, the report editors, have invested in producing this significant product for the GN&C community of practice. It was Russell and Chris who had the inspiration to create this report and they have done a masterful job in not only directly technically contributing to the report but also coordinating its overall development. It should be mentioned that some high-level limited work was previously performed under NESC sponsorship to capture the lessons learned over the course of the several decades NASA has been navigating space vehicles. This report however fills a unique gap by providing extensive technical details and, perhaps more importantly, providing the underlying rationale for each of the navigation filter best practices presented here. Capturing these rationales has proven to be a greatly needed but very challenging task. I congratulate the team for taking this challenge on.

The creation, and the wide dissemination of this report, is absolutely consistent with the NESC's commitment to engineering excellence by capturing and passing along, to NASA's next generation of engineers, the lessons learned emerging from the collective professional experiences of NASA's navigation system subject matter experts. I believe this book will not only provide relevant tutorial-type guidance for early career GN&C engineers that have limited real-world on the job experience but it should also serve as a very useful memory aid for more experienced GN&C engineers, especially as a handy reference to employ for technical peer reviews of navigation systems under development.

As the NASA Technical Fellow for GN&C I urge the reader (especially the "navigators" among you obviously) to invest the time to digest and consider how the best practices provided in this report should influence your own work developing navigation systems for the Agency's future missions. The editors and I recognize this will be a living document

and we sincerely welcome your feedback on this first edition of the report, especially your constructive recommendations on ways to improve and/or augment this set of best practices.

Cornelius J. Dennehy  
NASA Technical Fellow for GN&C  
January 2018

## Editor’s Preface

As the era of commercial spaceflight begins, NASA must ensure that lessons the US has learned over the first 50 years of the Space Age will continue to positively influence the continuing exploration and development of space. Of the many successful strands of this legacy, onboard navigation stands out as an early triumph of technology whose continuous development and improvement remains as important to future exploration and commercial development as it was in the era of *Gemini* and *Apollo*. The key that opened the door to practical and reliable onboard navigation was the discovery and development of the extended Kalman filter (EKF) in the 1960s, a story that has been well-chronicled by Stanley Schmidt [65], and Kalman filtering has far outgrown NASA’s applications over the intervening decades. What are less well-documented are the accumulated art and lore, tips and tricks, and other institutional knowledge that NASA navigators have employed to design and operate EKFs in support of dozens of missions in the *Gemini/Apollo* era, well over one hundred Space Shuttle missions, and numerous robotic missions, without a failure ever attributed to an EKF. To document the best of these practices is the purpose that motivates the contributors of the present document.

Kernals of such best practices have appeared, scattered throughout the open technical literature, but such contributions are limited by organizational publication policies, and in some case by technology export considerations. Even within NASA, there has heretofore not been any attempt to codify this knowledge into a readily available design handbook that could continue to evolve along with the navigation community of practice. As a result, even with the Agency, it is possible for isolated practitioners “not to know any better:” to fail to appreciate the subtleties of successful and robust navigation filter design, and to lack an understanding of the motivations for, and the implied cost/benefit trade, of many of the tried and true approaches to filter design.

Some limited progress toward filling this void has been made at a summary level in reports and briefings prepared for the NASA Engineering and Safety Center (NESC) [13]. In particular, one of a series of recommendations in Reference 13 “...directed towards the development of future non-human rated [rendezvous] missions...” included as its fourteenth recommendation the admonishment to “[u]tilize best practices for rendezvous navigation filter design.” This recommendation listed eight such practices, as follows:

- a. **Maintain an accurate representation of the target-chaser relative state estimation errors, including an accurate variance-covariance matrix.** This allows the filter to compute an appropriate gain matrix. It also aids the filter in appropriately editing unsuitable measurements.
- b. **Provide a capability for measurement underweighting that adapts to the current uncertainty in the filters state estimation error, as required to be consistent with the suboptimality**

**of the navigation filters measurement update.** Effective means for accomplishing this have been found to include:

- i. Modified second-order Gaussian state update method [30];
- ii. Multiplicative adjustment of the mapping of the state error covariance matrix into the measurement subspace, which occurs within the computation of the residual covariance [78]; and
- iii. Schmidt-Kalman state update [3] that utilizes the covariance matrix of “consider” parameters (i.e., states that the filter does not update, but for which it maintains a covariance).

Multiplicative adjustment of the measurement noise covariance matrix within the computation of the residual covariance (the “bump up R” method [3]) has been found to be less effective, and is not recommended unless other methods are not feasible.

- c. **Estimate states that model biases in sensor measurements and account for unmodeled accelerations.** Gauss-Markov models for these biases have been found to be more effective than random constant or random walk models. Random constant models can become stale, and random walk models can overflow during long periods without measurement updates.
- d. **Provide commands that allow for selective processing of individual measurement types.** If the filter utilizes an automated residual edit process, then the recommended command capability should be able to override the residual edit test.
- e. **Maintain a backup ephemeris, unaltered by measurement updates since initialization, which can be used to restart the filter without uplink of a new state vector.**
- f. **Provide a capability for reinitializing the covariance matrix without altering the current state estimate.**
- g. **Ensure tuning parameters are uplinkable to the spacecraft, and capable of being introduced to the filter without loss of onboard navigation data.**
- h. **Provide flexibility to take advantage of sensors and sensor suites full capability over all operating ranges.**

A subsequent briefing given for an NESC webinar [7] listed these as well as the following “additional considerations:”

- State Representation
  - Translational states
    - \* Dual inertial
    - \* Inertial/relative
    - \* Relative-only
  - Attitude states, as required
    - \* 3-parameter vs. 4-parameter
    - \* Multiplicative vs. additive update
- Covariance Factorization (or not)
  - U-D
  - “Square Root” Methods
- Measurement correlation

- Non-simultaneous measurements
- Backward smoothing (for BETs)
- Error Budgets
- Sensitivity Analysis
- IMU/Accelerometer processing
- Observability

While these summary-level lists give the community a place to start, they are lacking in some respects. They lack sufficient rationale that would motivate a designer to adopt them. Even if so motivated, a designer needs much more detailed information concerning how to implement the recommendations.

The present work is an attempt to address these shortcomings. Each contributor has selected one aspect of navigation filter design, or several closely related ones, as the basis of a chapter. Each chapter clearly identifies best practices, where a consensus of the community of practice exists. While it is sometimes difficult to cast aside one's opinions and express such a consensus, each contributor has made a best effort in this regard. Where a diversity of opinion exists, the chapter will summarize the arguments for and against each approach. Also, if promising new developments are currently afoot, the chapter will assess their prospects.

While the contributors strive for consistency of convention and notation, each has his own preferences, and readers may need to accommodate subtle differences along these lines as they traverse the book. The first chapter, which summarizes the EKF, sets the stage, and should be briefly perused by even seasoned navigators in order to become familiar with the conventions adopted for this work. Subsequent chapters should stand on their own, and may be consulted in any order.

While this is a NASA document concerned with space navigation, it is likely that many of the principles would apply equally to the wider navigation community. That said, readers should keep in mind that hard-earned best practices of a particular discipline do not always carry over to others, even though they may be seemingly similar. To assume so is a classic example of the logical fallacy *argumentum ad verecundiam*, or the argument from [false] authority.

Finally, the contributors intend for this work to be a living document, which will continue to evolve with the state of the practice.



## Notational Conventions

$\mathbb{A}$	A set
$a, \alpha, A$	Scalars
$\mathbf{q}, \mathbf{M}$	An array of scalars, e.g. column, row, matrix
$\mathbf{x}$	A point in an abstract vector space
$\vec{\mathbf{r}}$	A physical vector, i.e. an arrow in 3-D space
$\mathcal{F}$	A coordinate frame
$\mathbf{M}^\top$	The transpose of the array $\mathbf{M}$
$\ \mathbf{x}\ $	The (2-)norm of the vector $\mathbf{x}$
$y$	A random variable
$\mathbf{z}$	A random vector
$p_x(x)$	The probability density function of the random variable $x$ evaluated at the realization $x$
$\Pr(y < Y)$	The probability that $y < Y$
$E[\mathbf{z}]$	The expectation of the random vector $\mathbf{z}$
$\exp(t)$	The exponential function of $t$
$e^t$	$\exp(t)$ written as Euler's number raised to the $t$ power
$dx$	Leibniz' (total) differential of $x$
$\frac{dy}{dx}$	(First) (total) derivative of $y$ with respect to $x$
$\frac{d^n y}{dx^n} = \frac{d^n}{dx^n} y$	$n$ th (total) derivative of $y$ with respect to $x$
$\frac{\mathcal{F} d^n \vec{\mathbf{r}}}{dt^n}$	$n$ th (total) derivative of $\vec{\mathbf{r}}$ with respect to $t$ in frame $\mathcal{F}$
$\frac{\partial \mathbf{M}}{\partial \mathbf{x}}$	(First) partial derivative of $\mathbf{M}$ with respect to $\mathbf{x}$
$\frac{\partial M}{\partial x} \Big _{x_o}$	(First) partial of $M$ with respect to $x$ , evaluated at $x_o$



## CHAPTER 1

# The Extended Kalman Filter

Contributed by J. Russell Carpenter

As described in the preface, use of the Extended Kalman Filter (EKF) for navigation has a long history of flight-proven success. The EKF thus forms the foundational best practice advocated by this work, and it forms the basis for many of the best practices later chapters describe. The purpose of the present chapter is not to derive the EKF and its relations, but rather to present them in a basic form, as a jumping off point for the rest of the material we shall present. As we shall show, while the EKF is a powerful and robust algorithm, it is based on a few *ad hoc* assumptions, which can lead to misuses and misunderstandings. Many of the best practices we shall describe are tricks of the trade that address such issues.

### 1.1. The Additive Extended Kalman Filter

The *additive* EKF is distinguished from the *multiplicative* EKF (MEKF) by the form of its measurement update. The additive EKF is the usual and original form of the EKF, and when we refer to the EKF without a modifier, one may assume we mean the additive form.

**1.1.1. The Dynamics Model** Suppose we have a list of  $n$  real quantities that we need to know in order to perform navigation, and we have a differential equation that tells us how these quantities evolve through time, such as

$$\dot{\mathbf{X}}(t) = \mathbf{f}(\mathbf{X}(t), t) \tag{1.1}$$

where  $\mathbf{X} \in \mathbb{R}^n$ , which we call the *state vector*, contains the quantities of interest; we shall call the  $\mathbf{f}(\mathbf{X}, t)$  the *dynamics function*. If we knew these quantities perfectly at any time, (1.1) would allow us to know them at any other time. For a variety of reasons, this is not the case however; both the initial conditions and the dynamics function are corrupted by uncertainty.

Suppose instead that the quantities of interest are realizations of a random process,  $\mathbf{X}(t)$ , whose distribution at some initial time  $t_o$  is known to us, and whose evolution (forward) in time follows the stochastic differential equation given by

$$d\mathbf{X}(t) = \mathbf{f}(\mathbf{X}(t), t) + \mathbf{B}(t) d\mathbf{w}(t) \tag{1.2}$$

where the presence of the *process noise*  $d\mathbf{w}(t)$  reflects uncertainty in the dynamics. To interpret (1.2), imagine  $d\mathbf{w}(t)$  as the limit of a discrete sequence of random increments, as the time between increments goes to zero. The result will be a continuous but non-differentiable process; hence the notation  $\dot{\mathbf{X}}(t)$  has ambiguous meaning. Henceforth, we shall define our notation such that when we write

$$\dot{\mathbf{X}}(t) = \mathbf{f}(\mathbf{X}(t), t) + \mathbf{B}(t)\mathbf{w}(t) \tag{1.3}$$

what we mean is really (1.2).

Finally, suppose that the initial distribution of  $\mathbf{X}(t)$  is Gaussian, with mean and covariance given by

$$\mathbb{E}[\mathbf{X}(t_o)] = \bar{\mathbf{X}}_o \quad \text{and} \quad \mathbb{E}[(\mathbf{X}(t_o) - \bar{\mathbf{X}}_o)(\mathbf{X}(t_o) - \bar{\mathbf{X}}_o)^\top] = \mathbf{P}_o \quad (1.4)$$

and suppose that infinitesimal increments of  $\mathbf{w}(t)$  are Gaussian, with

$$\mathbb{E}[\mathbf{w}(t)] = \mathbf{0} \quad \text{and} \quad \mathbb{E}[\mathbf{w}(t)\mathbf{w}^\top(\tau)] = \mathbf{Q}(t)\delta(t - \tau) \quad (1.5)$$

where  $\mathbf{Q}(t)$  is the power spectral density function of  $\mathbf{w}(t)$ , and  $\delta(t - \tau)$  is the Dirac delta function. We shall also assume that

$$\mathbb{E}[\mathbf{w}(t)(\mathbf{X}(t_o) - \bar{\mathbf{X}}_o)^\top] = \mathbf{0}, \quad \forall t \quad (1.6)$$

We shall take (1.3) – (1.6) to define the *dynamics model* for the additive EKF. Note that even though we have assumed  $\mathbf{X}(t_o)$  and  $\mathbf{w}(t)$  are Gaussian, we cannot assume that  $\mathbf{X}(t)$  remains Gaussian for  $t > t_o$ , because  $\mathbf{f}$  may be a nonlinear function.

**1.1.2. The Measurement Model** In an ideal world, we might have devices for measuring all of the state vector components directly; then state determination would be simply a matter of collecting enough such observations to reduce the state uncertainty to sufficient levels. Unfortunately, this is almost never the case. Instead, like Socrates’ prisoners, we can usually only perceive noisy projections of the state elements, at discrete times,  $t_i$ , in the form of *measurements*:

$$\mathbf{Y}(t_i) = \mathbf{h}(\mathbf{X}(t_i), t_i) + \mathbf{v}(t_i) \quad (1.7)$$

where  $\mathbf{h}$  is a surjection from  $\mathbb{R}^n$  to  $\mathbb{R}^m$ . We shall assume that  $\mathbf{v}(t_i)$ , which we call the *measurement noise*, is a Gaussian sequence, with mean and covariance given by

$$\mathbb{E}[\mathbf{v}(t_i)] = \mathbf{0} \quad \text{and} \quad \mathbb{E}[\mathbf{v}(t_i)\mathbf{v}(t_j)^\top] = \mathbf{R}(t_i)\delta_{ij} \quad (1.8)$$

where  $\delta_{ij}$  is the Kronecker delta function. We shall also assume that

$$\mathbb{E}[\mathbf{v}(t_i)(\mathbf{X}(t_o) - \bar{\mathbf{X}}_o)^\top] = \mathbf{0} \quad \text{and} \quad \mathbb{E}[\mathbf{w}(t)\mathbf{v}(t_i)^\top] = \mathbf{0}, \quad \forall i, t \quad (1.9)$$

We shall take (1.7) – (1.9) to define the *measurement model* for the additive EKF. Note that we cannot assume that  $\mathbf{Y}(t_i)$  is Gaussian, because  $\mathbf{h}$  may be a nonlinear function. For compactness of notation, we shall often suppress the time argument and write (1.7) as

$$\mathbf{Y}_i = \mathbf{h}_i(\mathbf{X}_i) + \mathbf{v}_i \quad (1.10)$$

**Bayes’ Law, the Markov Property, and Observability** Bayes’ Law tells us how to update the conditional probability density function (PDF) of  $\mathbf{X}(t_i)$ , given a realization  $\mathbf{Y}(t_i)$  of the random process  $\mathbf{Y}(t_i)$ :

$$p_{\mathbf{X}_i|\mathbf{Y}_i}(\mathbf{X}_i|\mathbf{Y}_i) = p_{\mathbf{Y}_i|\mathbf{X}_i}(\mathbf{Y}_i|\mathbf{X}_i) \frac{p_{\mathbf{X}_i}(\mathbf{X}_i)}{p_{\mathbf{Y}_i}(\mathbf{Y}_i)} \quad (1.11)$$

If all the PDFs in (1.11) were known, it would be relatively simple to use (1.11) to estimate the state vector from a single measurement; our best estimate of the state would simply be the mean of  $p_{\mathbf{X}_i|\mathbf{Y}_i}(\mathbf{X}_i|\mathbf{Y}_i)$ . But to apply (1.11) to the navigation problem, where we have a time sequence of measurements,  $\mathbb{Y}_i = \{\mathbf{Y}_i, \mathbf{Y}_{i-1}, \dots, \mathbf{Y}_1\}$ , we need to consider how the state dynamics evolve.

Unlike (1.1), our dynamics model, given by (1.3), only runs forward in time. Hence, the state at any future time depends only on its history. Also, because the non-homogeneous inputs to (1.3) are uncorrelated by the Dirac function in (1.5), the value of the state at

any particular time in the future depends only on its present value, and its accumulated diffusion due to the process noise over the interval between now and the future time of interest. Random processes such as this are said to possess the *Markov Property*. Using this property, we can write (1.11) in terms of the measurement history as follows:

$$p_{\mathbf{X}_i|\mathbb{Y}_i}(\mathbf{X}_i|\mathbb{Y}_i) = p_{\mathbf{Y}_i|\mathbf{X}_i}(\mathbf{Y}_i|\mathbf{X}_i) \frac{p_{\mathbf{X}_i|\mathbb{Y}_{i-1}}(\mathbf{X}_i|\mathbb{Y}_{i-1})}{p_{\mathbf{Y}_i|\mathbb{Y}_{i-1}}(\mathbf{Y}_i|\mathbb{Y}_{i-1})} \quad (1.12)$$

Even if we could compute all of the PDFs in (1.12), we are not guaranteed that the sequence of measurements provide sufficient information to reduce the initial uncertainty of all the modes of (1.1). If the system given by (1.1) and (1.10) is such that use of (1.12) results in uncertainty in all the modes going asymptotically to zero in finite time, from any initial condition, then we say the system is *globally asymptotically observable*. If at least all of the unstable modes are observable, then we say the system is *detectable*. Unfortunately, for nonlinear systems, there is no known way to compute global observability. At best, under certain restrictions on (1.1) and (1.10), we can in principal establish local observability, in the neighborhood of a particular initial condition. However, this is a laborious calculation, often numerically unstable to evaluate. Also, note that observability is a property of the structure of (1.1) and (1.10), and hence is dependent on how one chooses to represent the navigation problem. Hence, a system that is observable with one representation may be unobservable with a different representation.

Kalman's original filter, which we now usually call the *linear Kalman filter* (LKF), is the result when the dynamics and measurement models are linear, Markov, Gaussian, and observable. An appreciation of the linear Kalman filter is essential to understanding the strengths and weaknesses of the EKF, although it is almost never the case that such assumptions are valid for real-world navigation problems.

**1.1.3. The Linear Kalman Filter** Suppose the dynamics and measurements are given by the following discrete-time linear models:

$$\mathbf{x}_i = \Phi_{i,i-1}\mathbf{x}_{i-1} + \Gamma_i\mathbf{u}_i \quad (1.13)$$

$$\mathbf{y}_i = \mathbf{H}_i\mathbf{x}_i + \mathbf{v}_i \quad (1.14)$$

with

$$E[\mathbf{x}_o] = \bar{\mathbf{x}}_o \quad \text{and} \quad E[(\mathbf{x}_o - \bar{\mathbf{x}}_o)(\mathbf{x}_o - \bar{\mathbf{x}}_o)^\top] = \mathbf{P}_o \quad (1.15)$$

$$E[\mathbf{u}_i] = \mathbf{0} \quad \text{and} \quad E[\Gamma_i\mathbf{u}_i\mathbf{u}_j^\top\Gamma_i^\top] = \mathbf{S}_i\delta_{ij} \quad (1.16)$$

and the moments of  $\mathbf{v}_i$  as given by (1.8). This system will be globally observable if the *observability Gramian* is strictly positive definite,

$$\mathbf{W}_k = \sum_{i=1}^k \Phi'_{i,1} \mathbf{H}_i^\top \mathbf{H}_i \Phi_{i,1} > 0 \quad (1.17)$$

i.e. it has full rank.

With such assumptions, Kalman showed [36] that Algorithm 1.1 provides an optimal (both minimum variance and maximum likelihood) estimate of the moments of the PDFs appearing in (1.11). Note that in Algorithm 1.1, the covariance recursion given by (1.19) and (1.22) does not depend on the measurement history, and hence one may compute the gain sequence,  $\mathbf{K}_i$ , off-line and store it as a time-indexed table or *schedule*, along with  $\Phi_{i,i-1}$  and  $\mathbf{H}_i$ . Also note that because the system is globally observable, there is no chance that it

**Algorithm 1.1** The Linear Kalman Filter

$$\hat{\mathbf{x}}_i^- = \Phi_{i,i-1} \hat{\mathbf{x}}_{i-1}^+, \quad \hat{\mathbf{x}}_0^+ = \bar{\mathbf{x}}_o \quad (1.18)$$

$$\mathbf{P}_i^- = \Phi_{i,i-1} \mathbf{P}_{i-1}^+ \Phi_{i,i-1}^\top + \mathbf{S}_i, \quad \mathbf{P}_0^+ = \mathbf{P}_o \quad (1.19)$$

$$\mathbf{K}_i = \mathbf{P}_i^- \mathbf{H}_i^\top (\mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top + \mathbf{R}_i)^{-1} \quad (1.20)$$

$$\hat{\mathbf{x}}_i^+ = \hat{\mathbf{x}}_i^- + \mathbf{K}_i (\mathbf{y}_i - \mathbf{H}_i \hat{\mathbf{x}}_i^-) \quad (1.21)$$

$$\mathbf{P}_i^+ = \mathbf{P}_i^- - \mathbf{K}_i \mathbf{H}_i \mathbf{P}_i^- \quad (1.22)$$

will fail to converge from any initial condition, except perhaps due to build up of numerical truncation and/or roundoff error.

If we further suppose the dynamics and measurements are given by linear time-invariant (LTI) models,

$$\mathbf{x}_i = \Phi \mathbf{x}_{i-1} + \Gamma \mathbf{u}_i \quad (1.23)$$

$$\mathbf{y}_i = \mathbf{H} \mathbf{x}_i + \mathbf{v}_i \quad (1.24)$$

then we may test its global observability using a somewhat simpler calculation than (1.17), as follows:

$$\mathbf{W} = \begin{bmatrix} \mathbf{H} \\ \mathbf{H}\Phi \\ \mathbf{H}\Phi^2 \\ \vdots \\ \mathbf{H}\Phi^{n-1} \end{bmatrix} > 0 \quad (1.25)$$

If the system is detectable, then it turns out that the covariance recursion given by (1.19) and (1.22) reaches a steady-state, which we denote  $\mathbf{P}_\infty$ . The corresponding gain is  $\mathbf{K}_\infty = \mathbf{P}_\infty \mathbf{H}^\top \mathbf{R}^{-1}$ . There exist numerous software packages that will compute such quantities, e.g. the *Matlab Control Systems Toolbox*, which may unfortunately lead to their misuse in inappropriate contexts. Perhaps worse, experts from other domains, who are familiar with techniques such as pole placement for control of LTI systems, may recognize that the steady-state linear Kalman filter is “just a pole placement algorithm,” and may infer that the EKF is not much more than a clever pole placement algorithm as well. As we shall show below, this is far from being the case; the EKF operates directly on the nonlinear system of interest, for which such LTI concepts have dubious applicability.

**1.1.4. The Linearized Kalman Filter** An immediately apparent generalization of the linear Kalman Filter is to use it to solve for small corrections to a nonlinearly propagated reference trajectory. While such an approach may have certain applications over limited time horizons, and/or for ground-based applications where an operator can periodically intervene, experience with onboard navigation systems has shown that such corrections can fail to remain small enough to justify the required approximations.

**1.1.5. The Extended Kalman Filter** There are a number of ways to proceed from Algorithm 1.1 to “derive” the EKF, but all contain a variety of *ad hoc* assumptions that are not guaranteed to hold in all circumstances. Most weaknesses and criticisms of the EKF arise from such assumptions. Rather than reproduce one or more of such derivations, we will simply point out that if one replaces (1.18) with an integral of (1.1) over the time between measurements, and computes the coefficient matrices appearing in (1.13) and (1.14)

as Jacobians evaluated at the current solution of (1.1), then the result is Algorithm 1.2, which bears more than a passing resemblance to the Kalman filter.

---

**Algorithm 1.2** A Naive Extension of the Kalman Filter

---

$$\hat{\mathbf{X}}_i^- = \int_{t_{i-1}}^{t_i} \mathbf{f}(\mathbf{X}(\tau), \tau) d\tau, \quad \mathbf{X}(t_{i-1}) = \hat{\mathbf{X}}_{i-1}^+, \quad \hat{\mathbf{X}}_0^+ = \bar{\mathbf{X}}_o \quad (1.26)$$

$$\mathbf{A}(t) = \left. \frac{\partial \mathbf{f}}{\partial \mathbf{X}} \right|_{\hat{\mathbf{X}}(t)}, \quad \mathbf{H}_i = \left. \frac{\partial \mathbf{h}_i}{\partial \mathbf{X}} \right|_{\hat{\mathbf{X}}_i^-} \quad (1.27)$$

$$\Phi(t_i, t_{i-1}) = \int_{t_{i-1}}^{t_i} \mathbf{A}(\tau) \Phi(t_i, \tau) d\tau, \quad \Phi(t_{i-1}, t_{i-1}) = \mathbf{I} \quad (1.28)$$

$$\mathbf{S}_i = \int_{t_{i-1}}^{t_i} \int_{t_{i-1}}^{t_i} \Phi(t_i, \tau) \mathbf{B}(\tau) \mathbb{E}[\mathbf{w}(\tau) \mathbf{w}^\top(\sigma)] \mathbf{B}^\top(\sigma) \Phi^\top(t_i, \sigma) d\tau d\sigma \quad (1.29)$$

$$\mathbf{P}_i^- = \Phi(t_i, t_{i-1}) \mathbf{P}_{i-1}^+ \Phi^\top(t_i, t_{i-1}) + \mathbf{S}_i, \quad \mathbf{P}_0^+ = \mathbf{P}_o \quad (1.30)$$

$$\mathbf{K}_i = \mathbf{P}_i^- \mathbf{H}_i^\top (\mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top + \mathbf{R}_i)^{-1} \quad (1.31)$$

$$\hat{\mathbf{X}}_i^+ = \hat{\mathbf{X}}_i^- + \mathbf{K}_i (\mathbf{Y}_i - \mathbf{h}_i(\hat{\mathbf{X}}_i^-)) \quad (1.32)$$

$$\mathbf{P}_i^+ = \mathbf{P}_i^- - \mathbf{K}_i \mathbf{H}_i \mathbf{P}_i^- \quad (1.33)$$


---

Several observations are in order regarding Algorithm 1.2.

- One might infer from (1.26) that  $\hat{\mathbf{X}}_i^- = \mathbb{E}[\mathbf{X} | \mathbb{Y}_{i-1}]$ . This would be a somewhat problematic inference however, since in general

$$\int_{t_{i-1}}^{t_i} \mathbf{f}(\mathbb{E}[\mathbf{X}(\tau), \tau]) d\tau \neq \mathbb{E} \left[ \int_{t_{i-1}}^{t_i} \mathbf{f}(\mathbf{X}(\tau), \tau) d\tau \right] \quad (1.34)$$

This implies that an initially Gaussian distribution for the state cannot in general remain Gaussian. At best, all we can hope is that  $\hat{\mathbf{X}}_i^- \approx \mathbb{E}[\mathbf{X} | \mathbb{Y}_{i-1}]$ .

- Let us define the *estimation error* as  $\mathbf{e}(t) = \mathbf{X}(t) - \hat{\mathbf{X}}(t)$ . Then since  $\hat{\mathbf{X}}(t) \neq \mathbb{E}[\mathbf{X}(t) | \mathbb{Y}_{i-1}]$ ,

$$\mathbf{P}(t) \neq \mathbb{E}[\mathbf{e}(t) \mathbf{e}^\top(t) | \mathbb{Y}_{i-1}] \quad (1.35)$$

At best, all we can hope is that  $\mathbf{P}(t) \approx \mathbb{E}[\mathbf{e}(t) \mathbf{e}^\top(t) | \mathbb{Y}_{i-1}]$ .

- Let us define the *innovation* as  $\mathbf{r}_i = \mathbf{Y}_i - \mathbf{h}(\hat{\mathbf{X}}_i^-)$ . Then since  $\mathbf{h}(\hat{\mathbf{X}}_i^-) \neq \mathbb{E}[\mathbf{Y}]$ ,

$$\mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top + \mathbf{R}_i \neq \mathbb{E}[\mathbf{r}_i \mathbf{r}_i^\top] \quad (1.36)$$

At best, all we can hope is that the above will hold approximately.

- Taken together, the approximations listed above imply that (1.32) and (1.33) can at best satisfy (1.12) only approximately, not only because the mean and covariance are approximations, but also because the PDFs fail to remain Gaussian, and hence fail to be characterized completely by only their first two moments.
- Even if all of the above are reasonable approximations, there is a problem with (1.33). The *posterior covariance* should be approximated by

$$\mathbf{P}_i^+ \approx \mathbb{E}[\mathbf{e}_i^+ (\mathbf{e}_i^+)^\top | \mathbb{Y}_i] \quad (1.37)$$

Let us assume that

$$\mathbf{Y}_i - \mathbf{h}_i(\hat{\mathbf{X}}_i^-) \approx \mathbf{H}_i \mathbf{e}_i^- + \mathbf{v}_i \quad (1.38)$$

Then (1.32) implies that

$$\mathbf{e}_i^+ = \mathbf{e}_i^- - \mathbf{K}_i \mathbf{H}_i \mathbf{e}_i^- - \mathbf{K}_i \mathbf{v}_i \quad (1.39)$$

and by our prior assumption that  $\mathbb{E}[\mathbf{e}_i^- \mathbf{v}_i^\top] = \mathbf{0}$ ,

$$\mathbf{P}_i^+ \approx \mathbb{E}[\mathbf{e}_i^+ (\mathbf{e}_i^+)^\top | \mathbb{Y}_i] \quad (1.40)$$

$$= (\mathbf{I} - \mathbf{K}_i \mathbf{H}_i) \mathbf{P}_i^- (\mathbf{I} - \mathbf{K}_i \mathbf{H}_i)^\top + \mathbf{K}_i \mathbf{R}_i \mathbf{K}_i^\top \quad (1.41)$$

Equation (1.41) is *Joseph's Formula*, and it holds for any gain  $\mathbf{K}_i$ . Only for the optimal gain and true covariance does (1.41) reduce to (1.33). Since (1.31) was computed with only an approximate covariance, and due to the various other approximations listed above as well,  $\mathbf{K}_i$  cannot be the optimal gain, so at best, (1.33) will only hold approximately. At worst, such approximations may lead to  $\mathbf{P}_i^+$  becoming non-positive definite, which is a significant issue. Because of its symmetric and additive form, (1.41) is much less likely (but not impossible!) to produce a non-positive definite  $\mathbf{P}_i^+$ .

- By our assumption that  $\mathbb{E}[\mathbf{w}(t) \mathbf{w}^\top(\tau)] = \mathbf{Q}(t) \delta(t - \tau)$ , one of the integrals in (1.29) should be annihilated by the Dirac function, resulting in

$$\mathbf{S}_i = \int_{t_{i-1}}^{t_i} \Phi(t_i, \tau) \mathbf{B}(\tau) \mathbf{Q}(\tau) \mathbf{B}^\top(\tau) \Phi^\top(t_i, \tau) d\tau \quad (1.42)$$

In any case, unlike for a discrete time dynamics model,

$$\begin{aligned} & \mathbb{E} \left[ \int_{t_{i-1}}^{t_i} \int_{t_{i-1}}^{t_i} \Phi(t_i, \tau) \mathbf{B}(\tau) \mathbf{w}(\tau) \mathbf{w}^\top(\sigma) \mathbf{B}^\top(\sigma) \Phi^\top(t_i, \sigma) d\tau d\sigma \right] \\ & \neq \mathbb{E} \left[ \int_{t_{i-1}}^{t_i} \Phi(t_i, \tau) \mathbf{B}(\tau) \mathbf{w}(\tau) d\tau \int_{t_{i-1}}^{t_i} \mathbf{w}^\top(\tau) \mathbf{B}^\top(\tau) \Phi^\top(t_i, \tau) d\tau \right] \end{aligned} \quad (1.43)$$

- In general, one would need to simultaneously integrate (1.26) and (1.28) due to their interdependence via the Jacobian  $\mathbf{A}$ . If the time between measurements is small enough, then if one were to employ a suitable approximation for (1.28), perhaps as simple as

$$\Phi(t_i, t_{i-1}) \approx \mathbf{I} + \mathbf{A}(t_i) (t_i - t_{i-1}) \quad (1.44)$$

then one may reasonably expect that a carefully chosen approximation would be no worse than the many other approximations inherent in the EKF. One may also consider the same or simpler approximations when considering approximations to (1.42).

- Because there no way to prove global observability for a nonlinear system, the EKF may fail to converge from some initial conditions, even if the system is locally observable in particular neighborhoods.

In light of the above observations, we conclude this section by presenting a slightly improved version of the EKF as Algorithm 1.3. In subsequent chapters, we shall describe additional improvements to the EKF.

---

**Algorithm 1.3** A Slightly Improved Extension of the Kalman Filter
 

---

$$\hat{\mathbf{X}}_i^- = \int_{t_{i-1}}^{t_i} \mathbf{f}(\mathbf{X}(\tau), \tau) d\tau, \quad \mathbf{X}(t_{i-1}) = \hat{\mathbf{X}}_{i-1}^+, \quad \hat{\mathbf{X}}_0^+ = \bar{\mathbf{X}}_o \quad (1.45)$$

$$\mathbf{A}(t) = \left. \frac{\partial \mathbf{f}}{\partial \mathbf{X}} \right|_{\hat{\mathbf{X}}(t)}, \quad \mathbf{H}_i = \left. \frac{\partial \mathbf{h}_i}{\partial \mathbf{X}} \right|_{\hat{\mathbf{X}}_i^-} \quad (1.46)$$

$$\Phi(t_i, t_{i-1}) = \text{a suitable approximation to (1.28)} \quad (1.47)$$

$$\mathbf{S}_i = \text{a suitable approximation to (1.42)} \quad (1.48)$$

$$\mathbf{P}_i^- = \Phi(t_i, t_{i-1}) \mathbf{P}_{i-1}^+ \Phi^\top(t_i, t_{i-1}) + \mathbf{S}_i, \quad \mathbf{P}_0^+ = \mathbf{P}_o \quad (1.49)$$

$$\mathbf{K}_i = \mathbf{P}_i^- \mathbf{H}_i^\top (\mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top + \mathbf{R}_i)^{-1} \quad (1.50)$$

$$\hat{\mathbf{X}}_i^+ = \hat{\mathbf{X}}_i^- + \mathbf{K}_i (\mathbf{Y}_i - \mathbf{h}_i(\hat{\mathbf{X}}_i^-)) \quad (1.51)$$

$$\mathbf{P}_i^+ = (\mathbf{I} - \mathbf{K}_i \mathbf{H}_i) \mathbf{P}_i^- (\mathbf{I} - \mathbf{K}_i \mathbf{H}_i)^\top + \mathbf{K}_i \mathbf{R}_i \mathbf{K}_i^\top \quad (1.52)$$


---

### 1.2. The Multiplicative Extended Kalman Filter

An interesting variation on the EKF is possible in the context of estimating attitude parameters. An attitude correction may be viewed as a small-angle rotation from a frame associated with the previous estimate to a frame associated with a current estimate. In this context, one may use the previous attitude estimate as a linearization reference for a linearized Kalman Filter's Jacobian matrices, and estimate the small-angle correction as the filter state. After each state update, one performs a rectification of the attitude reference by applying the small-angle correction. Since for many attitude representations, a frame rotation is multiplicative operation, this procedure has become known as the multiplicative EKF. Chapter 8 covers this subject.



## CHAPTER 2

# The Covariance Matrix

Contributed by J. Russell Carpenter

As Chapter 1 pointed out, the EKF estimate  $\hat{\mathbf{X}}(t)$  is at best an approximation for  $E[\mathbf{X}(t)|\mathbb{Y}]$ , and hence the EKF symbol  $\mathbf{P}(t)$  is at best an approximation for  $E[\mathbf{e}(t)\mathbf{e}^\top(t)|\mathbb{Y}]$ . In the present Chapter we discuss best practices for maintaining such approximations, and henceforth we will simply refer to  $\mathbf{P}(t)$  as the covariance matrix.

### 2.1. Metrics for Orbit State Covariances

To discuss what makes one covariance approximate better or worse than another, we must be able to compare matrices to one another, and hence we must adopt metrics. For matrices that serve as coefficients, there exist various matrix norms that serve. For covariance matrices, we are usually more interested in measures of the range of possible realizations that could be drawn from the probability distribution characterized by the covariance. The square root of the trace of the covariance is often a reasonable choice, since it is the root sum squared (RSS) formal estimation error. A drawback to the trace for orbital applications is that coordinates and their derivatives typically differ by several orders of magnitude, so that for example the RSS position error will dominate the RSS velocity error if the trace is taken over a  $6 \times 6$  Cartesian position and velocity state error covariance, unless some scaling is introduced. Although arbitrary scalings are possible, we discuss several metrics herein that have been found to be especially suitable to space applications.

Orbit determination is distinguishable from other types of positioning and navigation not only by the use of dynamics suitable to orbiting bodies, but also by a fundamental need to produce states that predict accurately. This need arises because spacecraft operations require accurate predictions for acquisition by communications assets, for planning future activities such as maneuvers and observations, for predicting conjunctions with other space objects, etc. For closed, i.e. elliptical, orbits about most planetary bodies, the two-body potential dominates all other forces by several orders of magnitude. Thus, in most cases, the ability of an orbit estimate to predict accurately is dominated by semi-major axis (SMA) error,  $\delta a$ . This is because SMA error translates into period error through Kepler's third law, and an error in orbit period translates into a secularly increasing error in position along the orbit track. As Reference [6] shows, the along-track drift per orbit revolution,  $\delta s$ , for an elliptical orbit with eccentricity  $e$  is bounded by

$$\delta s = -3\pi\sqrt{\frac{1+e}{1-e}}\delta a \quad \text{from periapse to periapse} \quad (2.1)$$

$$\delta s = -3\pi\sqrt{\frac{1-e}{1+e}}\delta a \quad \text{from apoapse to apoapse} \quad (2.2)$$

This phenomenon is especially significant for rendezvous and formation flying applications, where relative positions must be precisely controlled.

For a central body whose gravitational constant is  $\mu$ , the SMA of a closed Keplerian orbit,  $a$ , may be found from the *vis viva* equation,

$$-\frac{\mu}{2a} = -\frac{\mu}{r} + \frac{v^2}{2} \quad (2.3)$$

from which one can see that achieving SMA accuracy requires good knowledge of both radius,  $r$ , and speed,  $v$ . What is less obvious from (2.3) is that radius and speed errors must also be both well-balanced and well-correlated to maximize SMA accuracy [6, 9, 27], as Figure 1 illustrates. In this figure, radius error,  $\sigma_r$ , has been normalized by the squared

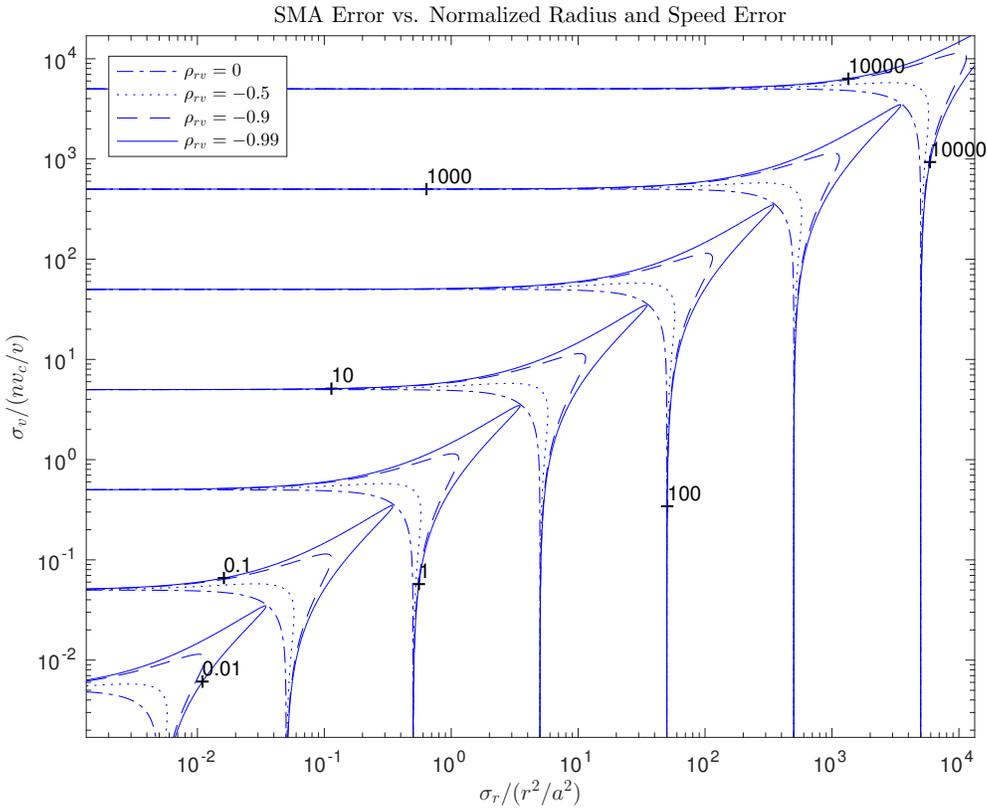


FIGURE 1. Semi-major axis accuracy depends on radius error, speed error, their correlation, and their balance. All scales are in units of position.

ratio of radius to SMA, and speed error,  $\sigma_v$ , has been normalized by  $nv_c/v$ , where the orbital rate is  $n = \sqrt{\mu/a^3}$ , and the circular speed is  $v_c = \sqrt{\mu/a}$ , to make the relationships illustrated be independent of any particular point in any particular closed orbit. Figure 1's contours of constant SMA error,  $\sigma_a$ , show that  $\sigma_a$  is dominated by radius error below a diagonal region, and dominated by speed error above the diagonal. When radius and speed errors are balanced, along the diagonal, SMA accuracy can be substantially improved by increasing (negative) correlation. Experience has shown that  $\sigma_a$  is one of the more useful

figures of merit for evaluating orbit determination performance, particularly for relative navigation applications.

In fact, it is easy to show that any function of two (scalar) random variables possesses a similar correlation and balance structure, at least to first order. For example, navigation requirements for atmospheric entry are often stated in terms of flight-path angle error,  $\delta\gamma$ . Since  $\sin \gamma = \vec{r}'/\|\vec{r}'\| \cdot \vec{v}/\|\vec{v}\|$ , then from geometrical considerations we should expect that  $\delta\gamma$  depends on the component of position error which is in the local horizontal plane, in the direction of the velocity vector, and on the component of velocity error that is normal to both velocity and angular momentum, i.e. binormal to the velocity vector. These are of course the in-plane components of position and velocity that are normal to radius and speed. Thus by using the pair  $\delta a$  and  $\delta\gamma$  as metrics, we can fully characterize the correlation and balance of the in-plane covariance components. The following subsections derive these relationships.

**2.1.1. Variance of an Arbitrary Function of Two Random Variables** Suppose there exists a random variable  $z$  which is a possibly nonlinear function of two other random variables,  $x$  and  $y$ , such that

$$z = f(x, y) \quad (2.4)$$

and let the joint covariance of  $x$  and  $y$  be given by

$$\mathbf{P} = \begin{bmatrix} \sigma_x^2 & \rho_{xy}\sigma_x\sigma_y \\ \rho_{xy}\sigma_x\sigma_y & \sigma_y^2 \end{bmatrix} \quad (2.5)$$

The variance of  $z$  is given by

$$\sigma_z^2 = \mathbb{E}[(z - \mathbb{E}[z])^2] \quad (2.6)$$

where

$$\mathbb{E}[z] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(\xi, \eta) p_z(f(\xi, \eta)) d\xi d\eta \quad (2.7)$$

Let  $\hat{x} = \mathbb{E}[x]$  and  $\hat{y} = \mathbb{E}[y]$ . Then  $\mathbb{E}[x - \hat{x}] = 0$  and  $\mathbb{E}[y - \hat{y}] = 0$ . Expanding  $f(x, y)$  around  $f(\hat{x}, \hat{y})$  in a Taylor series to first order, we find that

$$f(x, y) \approx f(\hat{x}, \hat{y}) + f_x \cdot (x - \hat{x}) + f_y \cdot (y - \hat{y}) \quad (2.8)$$

where  $f_x$  and  $f_y$  are the partials of  $f$  with respect to  $x$  and  $y$ , respectively; so, to first order,

$$\hat{z} = \mathbb{E}[z] = \mathbb{E}[f(x, y)] \approx f(\hat{x}, \hat{y}) \quad (2.9)$$

Now let us similarly expand  $(z - \mathbb{E}[z])^2$ :

$$(z - \mathbb{E}[z])^2 = (f(x, y) - f(\hat{x}, \hat{y}))^2 \quad (2.10)$$

$$\approx f_x^2 \cdot (x - \hat{x})^2 + 2f_x f_y \cdot (x - \hat{x})(y - \hat{y}) + f_y^2 \cdot (y - \hat{y})^2 \quad (2.11)$$

Taking expectations on both sides yields

$$\mathbb{E}[(z - \mathbb{E}[z])^2] = f_x^2 \mathbb{E}[(x - \hat{x})^2] + 2f_x f_y \mathbb{E}[(x - \hat{x})(y - \hat{y})] + f_y^2 \mathbb{E}[(y - \hat{y})^2] \quad (2.12)$$

$$\sigma_z^2 = f_x^2 \sigma_x^2 + 2f_x f_y \rho_{xy} \sigma_x \sigma_y + f_y^2 \sigma_y^2 \quad (2.13)$$

$$= \mathbf{F} \mathbf{P} \mathbf{F}^T \quad (2.14)$$

where  $\mathbf{F} = [f_x, f_y]$ . Since  $-1 < \rho_{xy} < 1$ , it is clear that a high negative correlation between  $x$  and  $y$  will minimize  $\sigma_z$  for given values of  $\sigma_x$  and  $\sigma_y$ , but if either  $f_x \sigma_x \gg f_y \sigma_y$  or  $f_x \sigma_x \ll f_y \sigma_y$ , the impact of the negative correlation will be insignificant. Thus, the only way to simultaneously achieve  $\sigma_z \ll f_x \sigma_x$  and  $\sigma_z \ll f_y \sigma_y$  is when  $\rho_{xy} \approx -1$  and

$f_x\sigma_x \approx f_y\sigma_y$ , which are the correlation and balance conditions mentioned above, and which occur along the diagonal of Figure 1.

Note also that by defining new variables scaled by their respective partial derivatives,  $\tilde{x} = xf_x$  and  $\tilde{y} = yf_y$ , and correspondingly  $\tilde{\sigma}_x = f_x\sigma_x$  and  $\tilde{\sigma}_y = f_y\sigma_y$ , then a normalization of the fashion described above is also possible:

$$\sigma_z = \sqrt{\tilde{\sigma}_x^2 + 2\rho_{xy}\tilde{\sigma}_x\tilde{\sigma}_y + \tilde{\sigma}_y^2} \quad (2.15)$$

**2.1.2. Semi-Major Axis Variance** To derive a relationship for semi-major axis variance, let us take variations on (2.3), which results in

$$\frac{\delta a}{a^2} = \frac{2\delta r}{r^2} + \frac{2v\delta v}{\mu} \quad (2.16)$$

If we replace the variations with deviations of random variables from their expectations, and the non-deviated terms with their expected values, we find that

$$(a - \hat{a}) = 2\hat{a}^2 \left( \frac{(r - \hat{r})}{\hat{r}^2} + \frac{\hat{v}(v - \hat{v})}{\mu} \right) \quad (2.17)$$

which by squaring and taking expectation yields the following equation for the SMA variance:

$$\sigma_a^2 = 4\hat{a}^4 \left\{ \frac{1}{\hat{r}^4} \sigma_r^2 + 2 \frac{\hat{v}}{\mu \hat{r}^2} \rho_{rv} \sigma_r \sigma_v + \frac{\hat{v}^2}{\mu^2} \sigma_v^2 \right\} \quad (2.18)$$

For the normalization used in Figure 1, rewrite (2.18) as

$$\sigma_a = 2\sqrt{\left( \frac{\sigma_r}{\hat{r}^2/\hat{a}^2} \right)^2 + 2\rho_{rv} \left( \frac{\sigma_r}{\hat{r}^2/\hat{a}^2} \right) \left( \frac{\sigma_v}{\mu/(\hat{a}^2\hat{v})} \right) + \left( \frac{\sigma_v}{\mu/(\hat{a}^2\hat{v})} \right)^2} \quad (2.19)$$

and note that  $\mu/(\hat{a}^2\hat{v}) = \hat{n}\hat{v}_c/\hat{v}$ . As mentioned above, normalizing radius and speed standard deviation in this manner permits comparison of data across all points in all closed orbits.

If the orbit is exactly circular, then further simplification of (2.18) is possible. In this case,  $a = r$  and  $v/\mu = T_p/2\pi$ , where  $T_p$  is the orbit period. Then (2.18) may be rewritten as

$$\sigma_a = 2\sqrt{\sigma_r^2 + 2 \left( \frac{T_p}{2\pi} \right) \rho_{rv} \sigma_r \sigma_v + \left( \frac{T_p}{2\pi} \right)^2 \sigma_v^2} \quad (2.20)$$

For orbit determination applications, the state representation most often chosen is a Cartesian inertial state vector,  $\mathbf{x} = [\vec{r}^T, \vec{v}^T]^T$ . Rewriting (2.3) as

$$a(\mathbf{x}) = \left( \frac{2}{\|\vec{r}\|} - \frac{\|\vec{v}\|^2}{\mu} \right)^{-1} \quad (2.21)$$

and taking partials yields

$$\left. \frac{\partial a}{\partial \mathbf{x}} \right|_{\hat{\mathbf{x}}} = \mathbf{F}_a(\hat{\mathbf{x}}) = 2\hat{a}^2 \begin{bmatrix} \frac{\hat{r}^T}{\hat{r}^3} & \frac{\hat{v}^T}{\mu} \end{bmatrix} \quad (2.22)$$

so that

$$\sigma_a^2 = \mathbf{F}_a(\hat{\mathbf{x}}) \mathbf{P}_x \mathbf{F}_a^T(\hat{\mathbf{x}}) \quad (2.23)$$

where  $\mathbf{P}_x$  is the state error covariance.

**2.1.3. Flight-Path Angle Variance** The flight-path angle,  $\gamma$ , is the angle between the velocity vector and the local horizontal plane; it is therefore the complement of the angle between the position and velocity vectors, so that

$$\gamma = \arcsin(\mathbf{u}_{\vec{r}} \cdot \mathbf{u}_{\vec{v}}) \quad (2.24)$$

where  $\mathbf{u}_{\vec{r}} = \vec{r}/r$  and  $\mathbf{u}_{\vec{v}} = \vec{v}/v$ . Taking partials with respect to  $\mathbf{x}$ , we find that

$$F_\gamma(\hat{\mathbf{x}}) = \frac{\partial \gamma}{\partial \mathbf{x}} \Big|_{\hat{\mathbf{x}}} = \frac{1}{\sqrt{1 - \sin^2 \hat{\gamma}}} \begin{bmatrix} \mathbf{u}_{\hat{v}}^\top - \sin \hat{\gamma} \mathbf{u}_{\hat{r}}^\top & \mathbf{u}_{\hat{r}}^\top - \sin \hat{\gamma} \mathbf{u}_{\hat{v}}^\top \end{bmatrix}, \quad -\frac{\pi}{2} < \hat{\gamma} < \frac{\pi}{2} \quad (2.25)$$

so that

$$\sigma_\gamma^2 = \mathbf{F}_\gamma(\hat{\mathbf{x}}) \mathbf{P}_\mathbf{x} \mathbf{F}_\gamma^\top(\hat{\mathbf{x}}) \quad (2.26)$$

which is a form suitable for use in an OD filter estimating a Cartesian inertial state.

For analysis, a simpler form of (2.26) is as follows. Let us define two vectors that are normal to both the position vector and a vector normal to the orbit plane,  $\vec{\mathbf{n}}$ : the unit in-track vector,

$$\mathbf{u}_{\Delta\vec{v}} = \mathbf{u}_{\vec{\mathbf{n}}} \times \mathbf{u}_{\vec{r}} = \frac{\mathbf{u}_{\vec{v}} - \mathbf{u}_{\vec{r}} \sin \gamma}{\sqrt{1 - \sin^2 \gamma}} \quad (2.27)$$

which defines a unit vector in the orbit plane that is along the orbit track at apoapsis and periapsis, and the unit bi-normal vector,

$$\mathbf{u}_{\vec{b}} = \mathbf{u}_{\vec{v}} \times \mathbf{u}_{\vec{\mathbf{n}}} = \frac{\mathbf{u}_{\vec{r}} - \mathbf{u}_{\vec{v}} \sin \gamma}{\sqrt{1 - \sin^2 \gamma}} \quad (2.28)$$

which defines a unit vector in the orbit plane that is along the position vector at apoapsis and peripasis. Let us next define a composite transformation matrix as follows. Let

$$\mathbf{M}_{rtn} = [\mathbf{u}_{\vec{r}_I} \quad \mathbf{u}_{\Delta\vec{v}_I} \quad \mathbf{u}_{\vec{\mathbf{n}}_I}] \quad (2.29)$$

where the subscript  $I$  indicates that the specified vector is expressed in an inertial basis. The unitary matrix  $\mathbf{M}_{rtn}^\top$  thus transforms coordinates of vectors in physical space, which are given in the frame  $I$ , to a coordinates defined in the basis given by the the radial, transverse, and normal unit vectors. Similarly, define  $\mathbf{M}_{vnb}$  as

$$\mathbf{M}_{vnb} = [\mathbf{u}_{\vec{v}_I} \quad \mathbf{u}_{\vec{\mathbf{n}}_I} \quad \mathbf{u}_{\vec{b}_I}] \quad (2.30)$$

Now define the block diagonal transformation matrix  $\mathbf{M}$  as

$$\mathbf{M} = \begin{bmatrix} \mathbf{M}_{rtn} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{M}_{vnb} \end{bmatrix} \quad (2.31)$$

Using the matrix  $\mathbf{M}$ , we can transform the state error covariance, given in inertial coordinates, such that its position error covariance is expressed in the ‘‘RTN’’ frame, and its velocity error covariance is in the ‘‘VNB’’ frame. Using the preceding results in (2.26) results in considerable simplification:

$$\sigma_\gamma^2 = [\mathbf{u}_{\hat{t}}^\top / \hat{r} \quad \mathbf{u}_{\hat{b}}^\top / \hat{v}] \begin{bmatrix} \mathbf{M}_{rtn}^\top & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{M}_{vnb}^\top \end{bmatrix} \mathbf{P}_\mathbf{x} \begin{bmatrix} \mathbf{M}_{rtn} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{M}_{vnb} \end{bmatrix} \begin{bmatrix} \mathbf{u}_{\hat{t}} / \hat{r} \\ \mathbf{u}_{\hat{b}} / \hat{v} \end{bmatrix} \quad (2.32)$$

$$= \left( \frac{\sigma_{\mathbf{r}_t}}{\hat{r}} \right)^2 + 2\rho_{\mathbf{r}_t \mathbf{v}_b} \left( \frac{\sigma_{\mathbf{r}_t}}{\hat{r}} \right) \left( \frac{\sigma_{\mathbf{v}_b}}{\hat{v}} \right) + \left( \frac{\sigma_{\mathbf{v}_b}}{\hat{v}} \right)^2 \quad (2.33)$$

which demonstrates the aforementioned assertion that flight-path angle error depends on the in-track component of position error, and the bi-normal component of velocity error. Note that (2.33) possesses the desirable feature that the relevant covariance information (in-track position variance and bi-normal velocity variance) is normalized by radius and

speed, allowing differing orbital conditions to be readily compared with one another. In most applications,  $\sigma_{r_t} \ll \hat{r}$  and  $\sigma_{v_b} \ll \hat{v}$  so that these ratios can be taken as small angles and expressed in angular measures commensurate with the units chosen for flight-path angle itself. Figure 2 employs this convention.

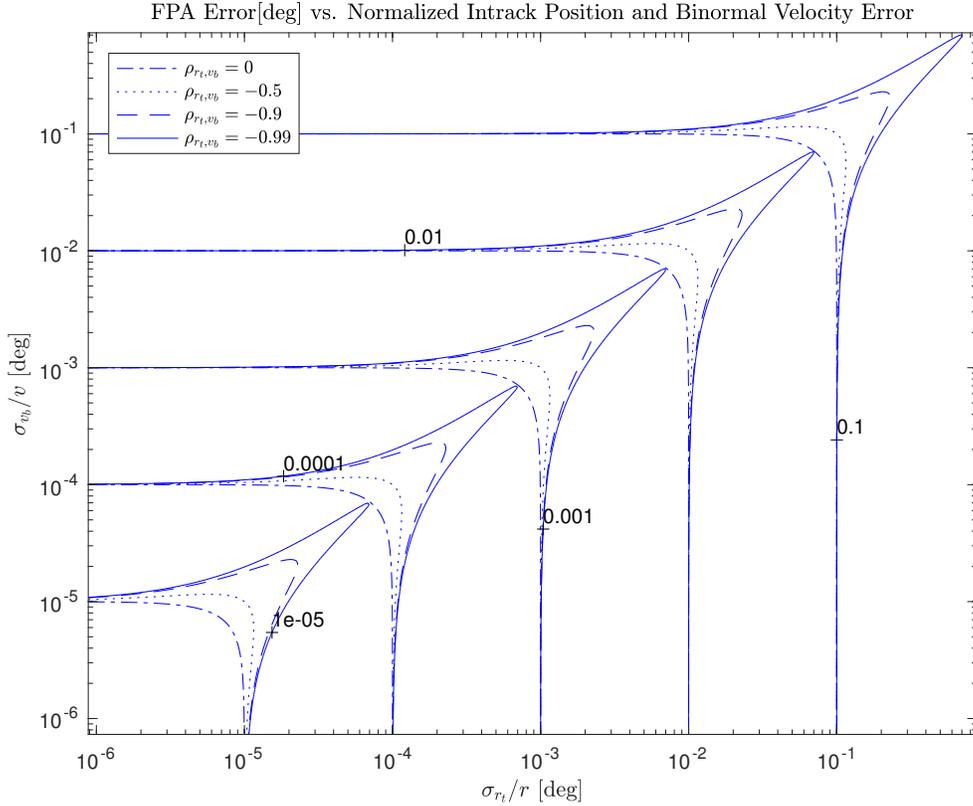


FIGURE 2. Flight-path angle accuracy depends on in-track position error and bi-normal velocity error, their correlation, and their balance.

**2.1.4. Summary of Orbit Determination Covariance Metrics** A recommended best practice for comparison of OD covariances is to use the semi-major axis standard deviation as a metric for most applications, with a secondary emphasis on flight-path angle standard deviation. For entry applications, a best practice is to use flight-path angle standard deviation at entry interface as the primary metric. A summary of these metrics is as follows.

For orbits that are very close to circular:

$$\sigma_a = 2\sqrt{\sigma_r^2 + 2\left(\frac{T_p}{2\pi}\right)\rho_{rv}\sigma_r\sigma_v + \left(\frac{T_p}{2\pi}\right)^2\sigma_v^2} \quad (2.34)$$

For elliptical orbits:

$$\sigma_a = 2\hat{a}^2\sqrt{\frac{1}{\hat{r}^4}\sigma_r^2 + 2\frac{\hat{v}}{\mu\hat{r}^2}\rho_{rv}\sigma_r\sigma_v + \frac{\hat{v}^2}{\mu^2}\sigma_v^2} \quad (2.35)$$

For normalization of radius and speed standard deviations across points in closed orbits:

$$\sigma_a = 2\sqrt{\left(\frac{\sigma_r}{\hat{r}^2/\hat{a}^2}\right)^2 + 2\rho_{rv}\left(\frac{\sigma_r}{\hat{r}^2/\hat{a}^2}\right)\left(\frac{\sigma_v}{\hat{n}\hat{v}_c/\hat{v}}\right) + \left(\frac{\sigma_v}{\hat{n}\hat{v}_c/\hat{v}}\right)^2} \quad (2.36)$$

For entry applications, and as a secondary metric:

$$\sigma_\gamma = \sqrt{\left(\frac{\sigma_{rt}}{\hat{r}}\right)^2 + 2\rho_{rtvb}\left(\frac{\sigma_{rt}}{\hat{r}}\right)\left(\frac{\sigma_{vb}}{\hat{v}}\right) + \left(\frac{\sigma_{vb}}{\hat{v}}\right)^2} \quad (2.37)$$

For use in an OD filter that is estimating a Cartesian inertial state vector:

$$\sigma_a = \sqrt{\mathbf{F}_a(\hat{\mathbf{x}})\mathbf{P}_x\mathbf{F}_a^\top(\hat{\mathbf{x}})}, \quad \mathbf{F}_a(\hat{\mathbf{x}}) = 2\hat{a}^2\begin{bmatrix} \hat{\mathbf{r}}^\top \\ \hat{r}^3 \\ \hat{\mathbf{v}}^\top \\ \mu \end{bmatrix} \quad (2.38)$$

$$\sigma_\gamma = \sqrt{\mathbf{F}_\gamma(\hat{\mathbf{x}})\mathbf{P}_x\mathbf{F}_\gamma^\top(\hat{\mathbf{x}})}, \quad \mathbf{F}_\gamma(\hat{\mathbf{x}}) = \begin{bmatrix} \mathbf{u}_t^\top & \mathbf{u}_b^\top \\ \hat{r} & \hat{v} \end{bmatrix} \quad (2.39)$$

## 2.2. Covariance Propagation

This section discusses best practices for implementing the covariance propagation recursion

$$\mathbf{P}_i^- = \Phi(t_i, t_{i-1})\mathbf{P}_{i-1}^+\Phi^\top(t_i, t_{i-1}) + \mathbf{S}_i, \quad \mathbf{P}_0^+ = \mathbf{P}_o \quad (2.40)$$

As Chapter 1 mentions, to achieve this goal we need suitable approximations for the state transition matrix,  $\Phi(t_i, t_{i-1})$ , and the process noise covariance,  $\mathbf{S}_i$ . However, the suitability of a given set of approximations is strongly dependent on specifics of the application. For example, if a set of measurements from which the state is fully observable are available at an interval that is a small fraction of the orbit period, without significant drop-outs, and prediction of the covariance far into the future of the time of availability of the measurements is not required, then simple models such as those to which Chapter 1 alluded have been successfully employed. Selection of appropriate covariance propagation approximations also depends strongly on the choice of state representation, which is the subject of Chapter 6. Therefore, this section will discuss the merits of some of the more common and generally applicable approaches, in the context of orbit determination.

**2.2.1. Matrix Ricatti Equation** Many textbooks on Kalman filtering derive (2.40) as the solution of a matrix Ricatti equation; using the notation of Chapter 1, this takes the following form:

$$\dot{\mathbf{P}}(t) = \mathbf{A}(t)\mathbf{P}(t) + \mathbf{P}(t)\mathbf{A}^\top(t) + \mathbf{Q}(t) \quad (2.41)$$

Use of (2.41) would seem to avoid the need to perform the integrations required to compute the state transition matrix (STM) and process noise covariance (PNC). In orbit determination practice however, (2.40) has been found to be more numerically stable and also, despite

the need to compute or approximate the state transition and process noise integrals, more efficient than (2.41).

**2.2.2. State Transition Matrix** A common approach in ground-based OD, especially in the batch least squares context, is to simultaneously integrate the STM along with the state vector,

$$\dot{\mathbf{X}}(t) = \mathbf{f}(\mathbf{X}(t), t), \quad \mathbf{X}(t_o) = \mathbf{X}_o \quad (2.42)$$

$$\dot{\Phi}(t, t_o) = \mathbf{A}(t)\Phi(t, t_o), \quad \Phi(t_o, t_o) = \mathbf{I} \quad (2.43)$$

When coupled with a good numerical integration algorithm, this method has excellent fidelity, which is rarely necessary in onboard OD applications. Lear [40] studied a number of practical methods for computing the STM in the onboard OD context. As his report is not widely available, we will summarize some key findings here.

Lear’s approach was to compare various orders of truncated Taylor series and Runge-Kutta approximations to the solution of (2.43). He used these STM approximations to propagate an initially diagonal covariance for one revolution in a two-body circular orbit around a point mass with the  $GM$  of Earth. By comparing these results to those he obtained using an analytic STM, Lear could compute the maximum step size that would result in a given relative accuracy for radius and speed formal standard deviations. Table 1 lists a few of Lear’s results. Notably, Method H has the desirable feature that simply saving the value

	Description	Method	Max. Step [sec]
A.	1 <sup>st</sup> -order Taylor	$\mathbf{I} + \mathbf{A}_i \Delta t$	0.125
B.	2 <sup>nd</sup> -order Taylor, ignoring $\dot{\mathbf{A}}$	$\mathbf{I} + \mathbf{A}_i \Delta t + \mathbf{A}_i^2 \Delta t^2 / 2$	1.0
C.	2 <sup>nd</sup> -order Taylor	$\mathbf{I} + \mathbf{A}_i \Delta t + (\dot{\mathbf{A}}_i + \mathbf{A}_i^2) \Delta t^2 / 2$	16
F.	1 <sup>st</sup> -order Runge-Kutta	$\mathbf{I} + \mathbf{A}_{i+.5} \Delta t$	0.14
G.	2 <sup>nd</sup> -order Runge-Kutta, with one evaluation of $\mathbf{A}$	$\mathbf{I} + \mathbf{A}_{i+.5} \Delta t + \mathbf{A}_{i+.5}^2 \Delta t^2 / 2$	14
H.	2 <sup>nd</sup> -order Runge-Kutta, with two evaluations of $\mathbf{A}$	$\mathbf{I} + (\mathbf{A}_i + \mathbf{A}_{i+1}) \Delta t + \mathbf{A}_i \mathbf{A}_{i+1} \Delta t^2 / 2$	16

TABLE 1. A few of Lear’s STM Comparison Results, for 1% relative error.

of the state Jacobian from the previous propagation step allows for more than an order of magnitude increase in allowable time step, with essentially the same computational burden as Method B. If it is not too burdensome to compute  $\mathbf{A}$  at the midpoint of the propagation step, then Method G offers nearly equivalent performance without the need to retain the previous value of  $\mathbf{A}$ . For higher-rate propagations, Method B offers far more accuracy than Method A with only a small additional computational burden. While Method A appears to be a poor choice for many applications, it does play a central role in some useful approximations to the process noise covariance, as the sequel shows.

**2.2.3. Process Noise Covariance** As stated in Chapter 1, nearly all practical methods for computing the process noise covariance assume that  $E[\mathbf{w}(t)\mathbf{w}^T(\tau)] = \mathbf{Q}(t)\delta(t - \tau)$ , so that (1.29) simplifies to a single integral for the process noise covariance<sup>1</sup>, given by (1.42),

<sup>1</sup>A notable exception is the work of Wright [76], which describes a correlated process noise model that is intended to account for gravity modeling errors in a “physically realistic” manner. Although this method

and repeated here:

$$\mathbf{S}_i = \int_{t_{i-1}}^{t_i} \boldsymbol{\Phi}(t_i, \tau) \mathbf{B}(\tau) \mathbf{Q}(\tau) \mathbf{B}^\top(\tau) \boldsymbol{\Phi}^\top(t_i, \tau) d\tau \quad (2.44)$$

Tapley, Schutz, and Born [69] describe two approximations to the portion of (2.44) which corresponds to position and velocity state noise, which have proven useful in both ground- and onboard-OD applications. Reference 69 refers to these methods as “State Noise Compensation” (SNC) and “Dynamic Model Compensation” (DMC). Before describing SNC and DMC however, we will consider some inappropriate models.

Chapter 1 pointed out that

$$\begin{aligned} & \mathbb{E} \left[ \int_{t_{i-1}}^{t_i} \int_{t_{i-1}}^{t_i} \boldsymbol{\Phi}(t_i, \tau) \mathbf{B}(\tau) \mathbf{w}(\tau) \mathbf{w}^\top(\sigma) \mathbf{B}^\top(\sigma) \boldsymbol{\Phi}^\top(t_i, \sigma) d\tau d\sigma \right] \\ & \neq \mathbb{E} \left[ \int_{t_{i-1}}^{t_i} \boldsymbol{\Phi}(t_i, \tau) \mathbf{B}(\tau) \mathbf{w}(\tau) d\tau \int_{t_{i-1}}^{t_i} \mathbf{w}^\top(\tau) \mathbf{B}^\top(\tau) \boldsymbol{\Phi}^\top(t_i, \tau) d\tau \right] \end{aligned} \quad (2.45)$$

Let us explore the implications of assuming equality of the expression above. Suppose we assume that the process noise increments are approximately constant over some particular interval  $\Delta t = t_i - t_{i-1}$ , and that  $\mathbb{E}[\mathbf{w}(t_i) \mathbf{w}^\top(t_j)] = \mathbf{W}(t_i) \delta_{ij}$ , where  $\delta_{ij}$  denotes the Kronecker delta function. Then,

$$\begin{aligned} & \mathbb{E} \left[ \int_{t_{i-1}}^{t_i} \boldsymbol{\Phi}(t_i, \tau) \mathbf{B}(\tau) \mathbf{w}(\tau) d\tau \int_{t_{i-1}}^{t_i} \mathbf{w}^\top(\tau) \mathbf{B}^\top(\tau) \boldsymbol{\Phi}^\top(t_i, \tau) d\tau \right] \\ & = \int_{t_{i-1}}^{t_i} \boldsymbol{\Phi}(t_i, \tau) \mathbf{B}(\tau) d\tau \mathbb{E}[\mathbf{w}(t_i) \mathbf{w}(t_i)^\top] \int_{t_{i-1}}^{t_i} \mathbf{B}^\top(\tau) \boldsymbol{\Phi}^\top(t_i, \tau) d\tau \\ & = \boldsymbol{\Gamma}_i \mathbf{W}_i \boldsymbol{\Gamma}_i^\top \end{aligned} \quad (2.46)$$

There is a subtlety with (2.46) that can lead to issues: if the time interval associated with the assumption that  $\mathbb{E}[\mathbf{w}(t_i) \mathbf{w}^\top(t_j)] = \mathbf{W}(t_i) \delta_{ij}$  is not the same as the time interval associated with the integral  $\boldsymbol{\Gamma}_i = \int_{t_{i-1}}^{t_i} \boldsymbol{\Phi}(t_i, \tau) \mathbf{B}(\tau) d\tau$ , then the process noise will not be consistently applied. For example, if the EKF is tuned at a particular time step using a particular noise covariance  $\mathbf{W}$ , and then for some reason the time step is changed, then one must retune the value of  $\mathbf{W}$ .

A similar issue occurs when the process noise covariance is chosen without regard for the dynamics, e.g. by setting it equal to a diagonal matrix of user-specified parameters. Whatever careful tuning has been done to choose such parameters will be invalidated by a change in the time step.

**2.2.3.1. State Noise Compensation** For SNC, as applied to OD, we assume velocity error is an uncorrelated random walk with fixed intensity in orbit-fixed coordinates, such as the RTN or VNB coordinates described above. Thus, assuming RTN coordinates without loss of generality, the process noise spectral density matrix becomes

$$\mathbf{Q}(t) = \mathbf{Q}_{rtn} = \begin{bmatrix} q_r & 0 & 0 \\ 0 & q_t & 0 \\ 0 & 0 & q_n \end{bmatrix} \quad (2.47)$$

---

has had occasional onboard application, it is more widely known for its inclusion in commercial-off-the-shelf software for ground-based OD.

We assume the transformation from orbit-fixed coordinates to the coordinates used for navigation, which are typically inertial coordinates, is approximately constant over the interval  $\Delta t = t_i - t_{i-1}$ , and ignore the correlation-inducing dependence of this transformation on the estimated position and velocity. This results in

$$\mathbf{B}(t) = \mathbf{B}_{rtn} = \begin{bmatrix} \mathbf{0}_{3 \times 3} \\ \mathbf{M}_{rtn} \end{bmatrix} \quad (2.48)$$

We also assume that  $\Delta t$  is small enough that a 1<sup>st</sup>-order Taylor series truncation (Lear's Model A) is adequate for modeling the STM  $\Phi(t_i, \tau)$  in the integrand of (2.44), and that  $\mathbf{M}_{rtn}$  is constant. With these assumptions, (2.44) becomes

$$\mathbf{S}_i = \begin{bmatrix} \tilde{\mathbf{Q}} \frac{\Delta t^3}{3} & \tilde{\mathbf{Q}} \frac{\Delta t^2}{2} \\ \tilde{\mathbf{Q}} \frac{\Delta t^2}{2} & \tilde{\mathbf{Q}} \Delta t \end{bmatrix} \quad (2.49)$$

where  $\tilde{\mathbf{Q}} = \mathbf{M}_{rtn} \mathbf{Q}_{rtn} \mathbf{M}_{rtn}^T$ . Note that with the SNC model, velocity covariance grows linearly with time, as expected for a random walk model, and hence we should expect the units of  $\sqrt{q_i}$ ,  $i = r, t, n$  to be meters per second<sup>3/2</sup>.

**2.2.3.2. State Noise Compensation for Maneuvers** During powered flight, it is often necessary to include additional process noise to accommodate maneuver magnitude and direction errors. One approach is to simply define an additional SNC process noise covariance, with intensities that are sized to the maneuvering errors. While this works fine for modeling maneuver magnitude errors, direction errors may be more accurately modeled by recognizing that a misaligned maneuver vector may be represented by  $(\mathbf{I}_3 - \delta\boldsymbol{\theta}^\times) \Delta \vec{v}_{nom}$ , where  $\delta\boldsymbol{\theta}^\times$  represents a skew-symmetric matrix of small angle misalignments, and  $\Delta \vec{v}_{nom}$  is the nominal maneuver vector. Thus,  $\Delta \vec{v}_{nom}^\times \delta\boldsymbol{\theta}$  is the error in the velocity increment due to maneuver direction errors, or if sensed accelerations are being fed-forward into the dynamics, due to IMU misalignments. To model these direction errors as a process noise term, let

$$\mathbf{B}(t) = \begin{bmatrix} \mathbf{0}_{3 \times 3} \\ \Delta \vec{v}_{nom}^\times \end{bmatrix} \quad (2.50)$$

and let  $q_\theta$  be the intensity of the maneuver direction noise. Then the SNC-style process noise for accommodating maneuver direction errors becomes

$$\mathbf{S}_i = q_\theta \begin{bmatrix} -(\Delta \vec{v}_{nom}^\times)^2 \frac{\Delta t^3}{3} & -(\Delta \vec{v}_{nom}^\times)^2 \frac{\Delta t^2}{2} \\ -(\Delta \vec{v}_{nom}^\times)^2 \frac{\Delta t^2}{2} & -(\Delta \vec{v}_{nom}^\times)^2 \Delta t \end{bmatrix} \quad (2.51)$$

since  $\Delta \vec{v}_{nom}^\times \Delta \vec{v}_{nom}^{\times T} = -(\Delta \vec{v}_{nom}^\times)^2$ . A version of this method was used by the Space Shuttle during powered flight with IMU-sensed accelerations.

**2.2.3.3. Dynamic Model Compensation** The DMC approach assumes the presence of exponentially-correlated acceleration biases, which are included as additional solve-fors in the filter state. As Chapter 5 and Appendix A discuss, a model for such biases is given by

$$b(t + \Delta t) = e^{-\frac{\Delta t}{\tau}} b(t) + \varpi(t) \quad (2.52)$$

where  $b(t_o) \sim N(0, p_{bo})$ , and  $\varpi(t) \sim N(0, \frac{q\tau}{2} (1 - e^{-\frac{2\Delta t}{\tau}}))$ . As Chapter 5 discusses,  $\tau$  is a time constant controlling the ‘‘smoothness’’ of the random process, and  $q$  is a power

spectral density that describes the intensity of the random input. While Chapter 5 discusses a variety of other bias models that might be used, the exponentially-correlated model has proved to be a *best practice* for applications in which there are measurements continually available to persistently excite it. Refer to Chapter 5 for a fuller discussion of the relative merits of various bias modeling approaches.

As above, without loss of generality we can assume the acceleration biases are aligned with the RTN frame, and again assume that  $\Delta t$  is small enough that a 1<sup>st</sup>-order Taylor series truncation is adequate for modeling the portion of the STM corresponding to position and velocity errors. With these assumptions, the terms appearing the integrand of (2.44) become

$$\mathbf{B}(t) = \mathbf{B}_{rtn} = \begin{bmatrix} \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} \\ \mathbf{M}_{rtn} \end{bmatrix} \quad (2.53)$$

and

$$\Phi(t + \Delta t, t) = \begin{bmatrix} \mathbf{I}_3 & \Delta t \mathbf{I}_3 & \{\tau \Delta t - \tau^2(1 - e^{-\Delta t/\tau})\} \mathbf{I}_3 \\ \mathbf{0}_{3 \times 3} & \mathbf{I}_3 & \tau(1 - e^{-\Delta t/\tau}) \mathbf{I}_3 \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} & e^{-\Delta t/\tau} \mathbf{I}_3 \end{bmatrix} \quad (2.54)$$

and the process noise covariance becomes<sup>2</sup>

$$\mathbf{S}_i = \begin{bmatrix} \gamma_{pp} \tilde{\mathbf{Q}} & \gamma_{pv} \tilde{\mathbf{Q}} & \gamma_{pa} \tilde{\mathbf{Q}} \\ \gamma_{pv} \tilde{\mathbf{Q}} & \gamma_{vv} \tilde{\mathbf{Q}} & \gamma_{va} \tilde{\mathbf{Q}} \\ \gamma_{pa} \tilde{\mathbf{Q}} & \gamma_{va} \tilde{\mathbf{Q}} & \gamma_{aa} \tilde{\mathbf{Q}} \end{bmatrix} \quad (2.55)$$

with

$$\gamma_{pp} = \frac{\tau^5}{2} \left\{ \left(1 - e^{-2\Delta t/\tau}\right) + \frac{2\Delta t}{\tau} \left(1 - 2e^{-\Delta t/\tau}\right) - 2 \left(\frac{\Delta t}{\tau}\right)^2 + \frac{2}{3} \left(\frac{\Delta t}{\tau}\right)^3 \right\} \quad (2.56)$$

$$\gamma_{pv} = \frac{\tau^4}{2} \left\{ \left(e^{-2\Delta t/\tau} - 1\right) - 2 \left(e^{-\Delta t/\tau} - 1\right) + \frac{2\Delta t}{\tau} \left(e^{-\Delta t/\tau} - 1\right) + \left(\frac{\Delta t}{\tau}\right)^2 \right\} \quad (2.57)$$

$$\gamma_{pa} = \frac{\tau^3}{2} \left\{ \left(1 - e^{-2\Delta t/\tau}\right) - \frac{2\Delta t}{\tau} e^{-\Delta t/\tau} \right\} \quad (2.58)$$

$$\gamma_{vv} = \frac{\tau^3}{2} \left\{ \left(1 - e^{-2\Delta t/\tau}\right) - 4 \left(1 - e^{-\Delta t/\tau}\right) + 2\Delta t/\tau \right\} \quad (2.59)$$

$$\gamma_{va} = \frac{\tau^2}{2} \left(1 - e^{-\Delta t/\tau}\right)^2 \quad (2.60)$$

$$\gamma_{aa} = \frac{\tau}{2} \left(1 - e^{-2\Delta t/\tau}\right) \quad (2.61)$$

---

<sup>2</sup>In (2.55), the matrix  $\tilde{\mathbf{Q}}$  has the same form as it does for the SNC method, but with  $\sqrt{q_i}$ ,  $i = r, t, n$  now representing acceleration intensities, with units of meters per second<sup>5/2</sup>. Also note that (2.55) assumes the same time constant is applicable to all three acceleration channels. While this is usually sufficient, it is straightforward to extend (2.55) to accommodate separate time constants for each channel.

2.2.3.4. *Explicit Dynamic Biases* While the DMC approach allows for quite general estimation of otherwise unmodeled forces on the spacecraft, it is often the case that the domain of application provides context that can narrow the filter designer’s focus. For example, it may be the case that the only under-modeled force of appreciable significance on the spacecraft is drag, or perhaps solar radiation pressure, within the context of the application. Alternatively, the application may require much higher resolution models than DMC, which might necessitate estimation of smaller forces with larger uncertainties such as Earth radiation pressure, spacecraft thermal emission, etc., or panel-based modeling of drag and/or solar radiation pressure, etc. In such cases, it is often useful to tailor the DMC approach so that it estimates model-specific biases, such a drag or SRP corrections, rather than modeling three general RTN biases. Similarly, during powered flight, maneuver magnitude and direction errors might be more successfully modeled explicitly.

As an example of an explicit bias, consider estimating a multiplicative correction to the density; a similar approach may be used for drag or solar radiation pressure coefficients. Let  $t$  denote geocentric coordinate time. Let  $\mathcal{R}$  denote a planetary-body-fixed, body-centric system of coordinates, aligned with the central body’s rotation axis. Let  $\mathcal{I}$  denote a body-centric, celestially-referenced system of coordinates, aligned with  $\mathcal{R}$  at an epoch  $t_o$ . Let  $\vec{r}$  represent the position of the center of gravity of a satellite, expressed in  $\mathcal{I}$ . Let  $\vec{v}$  represent the satellite’s velocity within  $\mathcal{I}$ . Let  $\vec{v}_r$  represent the satellite’s velocity within frame  $\mathcal{R}$ . Assume that  $\vec{r}$  evolves with respect to  $t$  and  $\mathcal{I}$  in the vicinity of  $t_o$  according to

$$\frac{\mathcal{I}d^2}{dt^2}\vec{r} = -\frac{\mu}{r^3}\vec{r} - \frac{1}{2}C_D\frac{A}{m}\rho\left(1 + \frac{\delta\rho}{\rho}\right)v_r\vec{v}_r \quad (2.62)$$

where  $r = \|\vec{r}\|$ ,  $v_r = \|\vec{v}_r\|$ ,  $\delta\rho$  is an atmospheric density disturbance,  $\rho$  is the undisturbed atmospheric density,  $A$  is the area of the satellite in a plane normal to  $\vec{v}_r$ ,  $m$  is the satellite mass, and  $C_D$  is the satellite’s coefficient of drag. Assume that  $\delta\rho/\rho$  is a random process that formally evolves as a first-order Gauss-Markov process, similar to a DMC bias:

$$\frac{d}{dt}\left(\frac{\delta\rho}{\rho}\right) = -\frac{1}{\tau}\left(\frac{\delta\rho}{\rho}\right) + w_\rho \quad (2.63)$$

where  $q_\rho$  is the intensity of  $w_\rho$ . Let the state vector be  $\mathbf{x} = [\vec{r}', \vec{v}', (\delta\rho/\rho)']'$ . With these assumptions, the state dynamics and noise input partials are

$$\mathbf{A}(t) = \begin{bmatrix} \mathbf{0}_{3\times 3} & \mathbf{I}_3 & \mathbf{0}_{3\times 1} \\ \mathbf{G}(t) + \mathbf{D}_r(t) & \mathbf{D}_v(t) & \vec{d}(t) \\ \mathbf{0}_{1\times 3} & \mathbf{0}_{1\times 3} & -1/\tau \end{bmatrix}, \mathbf{B}(t) = \begin{bmatrix} \mathbf{0}_{3\times 1} \\ \mathbf{0}_{3\times 1} \\ 1 \end{bmatrix} \quad (2.64)$$

where  $\mathbf{G}(t)$  is the gravity gradient matrix,  $\vec{d}(t) = -\frac{1}{2}C_D\frac{A}{m}\rho v_r\vec{v}_r$  is the nominal drag acceleration, and  $\mathbf{D}_r$  and  $\mathbf{D}_v$  are partials of the drag acceleration with respect to position and velocity, respectively<sup>3</sup>.

Using the DMC results from above, with the assumption that the nominal drag acceleration is approximately constant over the integration time,  $\Delta t$ , the term  $\Phi(t + \Delta t, t)\mathbf{B}(t)$

<sup>3</sup>The sensitivity  $\mathbf{D}_r$  contains terms that are roughly proportional to the drag acceleration magnitude divided by the atmospheric scale height, and to the product of drag acceleration magnitude and the ratio of planetary rotation rate to speed relative to the atmosphere. For nearly all spacecraft, these terms will be many orders magnitude smaller than the gravity gradient, which is proportional to the gravity acceleration divided by the radius. So  $\mathbf{D}_r$  can usually be neglected. For reference,  $\mathbf{D}_v = -\frac{1}{2}C_D\frac{A}{m}\rho v_r(\vec{u}_{v_r}\vec{u}'_{v_r} + \mathbf{I}_3)$ , and  $\mathbf{D}_r = d/R_s(\vec{u}_r\vec{u}'_{v_r}) - \mathbf{D}_v\vec{\omega}^\times$ , where  $\vec{\omega}^\times$  is the skew-symmetric “cross-product” matrix formed from the central body’s rotation rate vector, and the notation  $\vec{u}_{(\cdot)}$  indicates the unit vector of its subscript.

appearing in the integrand of (2.44) becomes

$$\Phi(t + \Delta, t)\mathbf{B}(t) = \begin{bmatrix} \{\tau\Delta t - \tau^2(1 - e^{-\Delta t/\tau})\} \vec{\mathbf{d}} \\ \tau(1 - e^{-\Delta t/\tau}) \vec{\mathbf{d}} \\ e^{-\Delta t/\tau} \end{bmatrix} \quad (2.65)$$

and the process noise for a proportional density bias becomes

$$\mathbf{S}_i = q_\rho \begin{bmatrix} \gamma_{pp} \vec{\mathbf{d}} \vec{\mathbf{d}}^T & \gamma_{pv} \vec{\mathbf{d}} \vec{\mathbf{d}}^T & \gamma_{pa} \vec{\mathbf{d}} \\ \gamma_{pv} \vec{\mathbf{d}} \vec{\mathbf{d}}^T & \gamma_{vv} \vec{\mathbf{d}} \vec{\mathbf{d}}^T & \gamma_{va} \vec{\mathbf{d}} \\ \gamma_{pa} \vec{\mathbf{d}} & \gamma_{va} \vec{\mathbf{d}} & \gamma_{aa} \end{bmatrix} \quad (2.66)$$

**2.2.3.5. Episodic Dynamic Biases** It has sometimes been found to be the case, particularly for crewed missions, that episodic spacecraft activities can produce un-modeled accelerations. In the early days of manned spaceflight, events such as vents, momentum unloads, RCS firings, etc. that may perturb a spacecraft's trajectory were not well modeled, and came to be described as FLAK, which was supposed to be an acronym for (un)-Fortunate Lack of Acceleration Knowledge.

If the mean time between such activities can be characterized, along with the expected intensity of the acceleration, a compound Poisson-Gaussian process noise model may be effective. Conveniently, it turns out that the covariance of a linear system driven by a train of Gaussian-distributed impulses whose arrival times follow a Poisson distribution is the same as the covariance of the same system driven by a white noise input process, except for the scaling of the process noise covariance by the Poisson process rate parameter [24].

To understand this result, consider a linear model of the error in a spacecraft trajectory as follows

$$\dot{\mathbf{x}}(t) = \mathbf{A}(t)\mathbf{x}(t) + \mathbf{B}(t)\mathbf{u}(t) \quad (2.67)$$

where  $\mathbf{x}$  represents the deviation of the actual position/velocity state from its estimated or nominal value, and  $\mathbf{u}$  represents the FLAK. Then, if we make use of inertial coordinates,

$$\mathbf{A}(t) = \begin{bmatrix} \mathbf{0} & \mathbf{I} \\ \mathbf{G}(t) & \mathbf{0} \end{bmatrix}, \quad \mathbf{B}(t) = \begin{bmatrix} \mathbf{0} \\ \mathbf{M}(t) \end{bmatrix} \quad (2.68)$$

where  $\mathbf{G}$  represents the gravity gradient matrix, and  $\mathbf{M}$  is the direction cosine matrix rotating the supposed body-fixed FLAK into the inertial frame. There is no general solution to this differential equation, but over short time intervals, we can assume that

$$\mathbf{x}(t_k) = \Phi(t_k, t_{k-1})\mathbf{x}_{k-1} + \int_{t_{k-1}}^{t_k} \Phi(t_k, \tau)\mathbf{B}(\tau)\mathbf{u}(\tau) d\tau \quad (2.69)$$

where

$$\Phi(t_k, t_{k-1}) = \mathbf{I} + \mathbf{A}(t_k)(t_k - t_{k-1}) \quad (2.70)$$

Since the input is a sequence of impulses of (random) length  $n_k$ ,

$$\mathbf{u}(t) = \sum_{i=1}^{n_k} \mathbf{u}_i \delta(t - t_i) \quad (2.71)$$

with  $t_k < t_i < t_{k+1}$ , the delta functions annihilate the integral and our model becomes (with  $\Delta t_i = t_k - t_i$ ):

$$\mathbf{x}(t_k) = \Phi(t_k, t_{k-1})\mathbf{x}_{k-1} + \sum_{i=1}^{n_k} \begin{bmatrix} \vec{\mathbf{k}}_i \Delta t_i \\ \vec{\mathbf{k}}_i \end{bmatrix} \mathbf{u}_i \quad (2.72)$$

where  $\vec{\kappa}_i$  is the unit inertial direction vector for the FLAK event. Note that the input response is in some sense fundamentally an increment to the entire state vector at each  $t_k$ ; we can however compute an equivalent zero-order hold acceleration by dividing the velocity increment by the time step  $t_k - t_{k-1}$ .

Since we assume each impulse is Gaussian, the input response has zero mean. To find a tractable form for the covariance, assume that the direction of the FLAK event is constant over each interval  $t_k - t_{k-1}$ . This assumption assures that the impulses are identically distributed over each sampling interval. Then, the process noise covariance is given by Reference 24:

$$\mathbf{S}(t_k) = q\lambda \int_{t_{k-1}}^{t_k} \begin{bmatrix} \vec{\kappa}_k(t_k - \tau) \\ \vec{\kappa}_k \end{bmatrix} [\vec{\kappa}'_k(t_k - \tau), \vec{\kappa}'_k] d\tau \quad (2.73)$$

where  $q$  is the intensity of the Gaussian impulses and  $\lambda$  is the rate parameter of the Poisson process. Carrying out the integration results in

$$\mathbf{S}(t_k) = q\lambda \begin{bmatrix} \vec{\kappa}_k \vec{\kappa}'_k \Delta t^3 / 3 & \vec{\kappa}_k \vec{\kappa}'_k \Delta t^2 / 2 \\ \vec{\kappa}_k \vec{\kappa}'_k \Delta t^2 / 2 & \vec{\kappa}_k \vec{\kappa}'_k \Delta t \end{bmatrix} \quad (2.74)$$

**2.2.3.6. Computational Considerations** The primary computational issues that can affect covariance propagation may be broadly characterized as underflow and overflow. Underflow can occur especially because many of the process noise parameters described above may often have values that approach computational truncation limits, which can lead to non-positive-definite process noise covariances. Overflow can similarly occur when truncation limits are approached by very large covariance values such as can occur with long propagation times. Because the orbital dynamics are at best marginally stable, even propagating without process noise can result in differences of many orders of magnitude between largest and smallest eigenvalues. This problem will be exacerbated if process noise is present, since all of the process noise models described above introduce unbounded position covariance error growth<sup>4</sup>.

Simple tricks like enforcing symmetry, or adding a small positive diagonal matrix, will not always ensure positive eigenvalues in such cases. A better solution is to maintain the covariance in factorized form, for example as Chapter 7 describes. In lieu of a fully factorized filtering approach, process noise factors may be computed from their factorizations. Chapter 5 shows a few examples of Cholesky factorizations that may be employed in this fashion.

**2.2.4. Tuning the Covariance Propagation** Since even the best practices this Chapter has discussed are at best approximations, it is inevitable that EKF designers must perform some artful tuning of the free parameters to achieve acceptable results. Furthermore, computational limitations of flight computers often lead to the need for compromises in modeling fidelity. What one generally hopes to accomplish via tuning of the covariance propagation is that any approximations or compromises the EKF has had to endure to be implementable have not impaired its covariance's accuracy too much. In particular, one would like to compute an idealized "truth" covariance matrix, based on the best-available models and data, and adjust the EKF's "formal" covariance via the tuning process to yield some semblance of a match.

---

<sup>4</sup>Reference 8 proposes an approximate "solution" to this problem via a Floquet analysis of a modified set of covariance propagation dynamics that include artificially-introduced damping.

In some cases, it is possible to compute the true covariance. In particular, if we are studying a linear system, and the random components have zero-mean Gaussian distributions, then the mean errors will be zero, and we can use linear covariance analysis [20, 48, 50] to compute the covariances. It is often possible to approximate the performance for a nonlinear system with this technique by linearization. This is often a first step in early conceptual design studies.

One may divide tuning of the covariance propagation into those activities a designer performs (1) during the detailed development of a system, prior to the collection of any flight data, and (2) during the commissioning of a new system or an existing system in a new application, when flight data are available. For pre-flight detailed design studies, one generally simulates the system, so one has access to truth data. One can also run the mission simulation many times, generating an ensemble of parallel results, performing a Monte Carlo analysis. During and after the actual mission, we never have access to truth data. At best, we can reconstruct the trajectory after the fact using more sophisticated processing and additional data that were not available in real-time. For near-realtime analysis, we can compare current definitive states to predictions from previous epochs. These predictions can come from either mission products generated in real-time at a past epoch, or past reconstructions of the trajectory. In all cases, the best we have are differences between estimates, not errors from the truth. There are several empirical approximations to the true covariance that one might use in these situations.

2.2.4.1. *Empirical Approximations of the True Covariance* Let  $\mathbf{e}_j(t_i)$  represent the random error vector at time  $t_i$  for case  $j$  from a Monte Carlo simulation, and  $\{\mathbf{e}(t_i)\}$  be the set of all cases at time  $t_i$ . Assume the total number of cases is  $K$ , and the total number of time samples is  $N$ .

**Time Series Expectation:** We can take statistics of the error realizations across time for each case,  $\mathbf{e}_j(t_i)$ , to get the time series expectations for each case:

$$\hat{\mathbf{E}}_t(\mathbf{e}_j) = \frac{1}{N} \sum_{i=1}^N \mathbf{e}_j(t_i) \quad (2.75)$$

$$\hat{\mathbf{E}}_t(\mathbf{e}_j \mathbf{e}_j^\top) = \frac{1}{N-1} \sum_{i=1}^N \mathbf{e}_j(t_i) \mathbf{e}_j^\top(t_i) \quad (2.76)$$

If the data are stationary<sup>5</sup>, the time series statistics are usually an adequate approximation, if we have considered a long enough time span. In some systems, described as *ergodic*, a long time series is in some sense equivalent to a large number of shorter Monte Carlo cases.

**Ensemble Expectation:** We can take statistics of the error realizations over all the cases at each time sample to get the ensemble expectations:

$$\hat{\mathbf{E}}_e\{\mathbf{e}(t_i)\} = \frac{1}{K} \sum_{j=1}^K \mathbf{e}_j(t_i) \quad (2.77)$$

$$\hat{\mathbf{E}}_e\{\mathbf{e}(t_i) \mathbf{e}^\top(t_i)\} = \frac{1}{K-1} \sum_{j=1}^K \mathbf{e}_j(t_i) \mathbf{e}_j^\top(t_i) \quad (2.78)$$

---

<sup>5</sup>Stationary data are those for which the statistics do not change when the time origin shifts.

The ensemble statistics will generally give the best indication of performance if the data are non-stationary, so long as we use an adequate number of Monte Carlo cases.

Since there is nothing analogous to an ensemble of Monte Carlo cases for flight data, we cannot use ensemble statistics as defined above. Let  $\mathbf{d}$  represent the random difference vector between the quantity of interest and its comparison value. We can apply time series statistics, but as the mission evolves, the span of the time series continually extends, so we have to decide which subsets of the entire mission span to use, e.g. the time series extending back over the entire history of the mission, extending back only over some shorter interval, etc., and also how frequently to recompute the time series statistics, e.g. continuously, once per day, etc. There are some other approximations to the expectation that we might use here.

**Sliding Window Time Series Expectation:** We can take statistics of the difference realizations,  $\mathbf{d}(t_i)$ , across a sliding window extending  $\Delta t$  into the past from each observation, for each case, to get the  $\Delta t$ -sliding window time series expectations:

$$\hat{\mathbf{E}}_{t,\Delta t}(\mathbf{d}) = \frac{1}{\Delta n} \sum_{i=0}^{\Delta n-1} \mathbf{d}(t_{N-i}) \quad (2.79)$$

$$\hat{\mathbf{E}}_{t,\Delta t}(\mathbf{d}\mathbf{d}^\top) = \frac{1}{\Delta n - 1} \sum_{i=0}^{\Delta n-1} \mathbf{d}(t_{N-i})\mathbf{d}(t_{N-i})^\top \quad (2.80)$$

where  $\Delta n$  is the number of time samples in the window  $\Delta t$ .

**Period-Folding Expectation:** If the data are periodic, we can break up the data into  $K$  spans of one period in duration each, and shift the time origin of each span so that the data are “folded” into the same, one-period-long interval. We can then take ensemble statistics over times at the same phase angle,  $t_{\phi i}$ , within each period.

$$\hat{\mathbf{E}}_f\{\mathbf{d}(t_{\phi i})\} = \frac{1}{K} \sum_{j=1}^K \mathbf{d}_j(t_{\phi i}) \quad (2.81)$$

$$\hat{\mathbf{E}}_f\{\mathbf{d}(t_{\phi i})\mathbf{d}(t_{\phi i})^\top\} = \frac{1}{K - 1} \sum_{j=1}^K \mathbf{d}_j(t_{\phi i})\mathbf{d}_j^\top(t_{\phi i}) \quad (2.82)$$

It is often useful to fold the data into bins of equal mean anomaly. This is especially useful for orbits with notable eccentricity, since it ensures that a roughly equal number of time points will be present in each bin.

**Sliding Window Period-Folding Expectation:** Period-folding can obviously be applied over a sliding window as well, with each window extending  $n$  periods into the past.

$$\hat{\mathbf{E}}_{f,n}\{\mathbf{d}(t_{\phi i})\} = \frac{1}{n} \sum_{j=0}^{n-1} \mathbf{d}_{K-j}(t_{\phi i}) \quad (2.83)$$

$$\hat{\mathbf{E}}_{f,n}\{\mathbf{d}(t_{\phi i})\mathbf{d}(t_{\phi i})^\top\} = \frac{1}{n - 1} \sum_{j=0}^{n-1} \mathbf{d}_{K-j}(t_{\phi i})\mathbf{d}_{K-j}^\top(t_{\phi i}) \quad (2.84)$$

This is especially useful for identifying secular trends in periodic data sets.

2.2.4.2. *Tuning for Along-track Error Growth* As described above, the position error component along the orbit track will dominate covariance propagation error, and so the most important step in tuning the covariance propagation is to ensure that this component grows no faster or slower than it should based on the truncations and approximations that the EKF design has employed. One may use any of the analytical or empirical methods described above to estimate the “true” covariance. For example, for preflight analysis, one may generate a time series or ensemble of time series of differences between states propagated using the formal models the filter employs, and a best available “truth” model of the system. One can then compare the appropriate empirical covariance computed from this data set to the filter’s formal covariance, and adjust the process noise intensities until a reasonable match occurs. For flight data analysis, one may similarly difference across overlaps between predictive and definitive states, and compare these empirical covariances of these differences to the sum of the predictive and definitive formal covariances from the filter.

If one uses the SMC method, the primary “knob” for tuning the alongtrack covariance growth rate is the corresponding alongtrack component of the process noise intensity  $q_T$  or  $q_V$ , depending on whether RTN or VNB components are used, respectively. Essentially, an impulse along the velocity vector, or change in speed, causes a change in SMA, corresponding to a change in period, and hence a secular growth in position error along the orbit, as discussed at the beginning of this Chapter. This mechanism is especially transparent for near-circular orbits, and some simple analysis yields a good starting point. One may find a fuller exposition of the following result in Reference **21**.

For near-circular orbits, the position components of the integrand in (2.44) become, in RTN coordinates,

$$\Phi_{rv}(\Delta t) \begin{bmatrix} q_R & 0 & 0 \\ 0 & q_T & 0 \\ 0 & 0 & q_N \end{bmatrix} \Phi_{rv}^T(\Delta t) \quad (2.85)$$

where  $\Phi_{rv}(\Delta t)$  is given, per Hill, Clohessy, and Wiltshire, by

$$\Phi_{rv}(\Delta t) = \begin{bmatrix} \sin(n\Delta t)/n & 2(1 - \cos(n\Delta t))/n & 0 \\ 2(\cos(n\Delta t) - 1)/n & 4\sin(n\Delta t)/n - 3\Delta t & 0 \\ 0 & 0 & \sin(n\Delta t)/n \end{bmatrix} \quad (2.86)$$

Retaining only secular terms and carrying out the integral, the along-track component of the process noise covariance becomes

$$S_T(\Delta t) \approx 3\Delta t^3 q_T \quad (2.87)$$

an approximation which holds for  $\Delta t > T_p$ . Thus, one may use an empirical covariance of the along-track error after one orbit period, such as  $\hat{\sigma}_{\delta_s}^2 = \hat{E}_f\{\delta s^2\}$ , to derive a starting point from which to tune  $q_T$ , as

$$q_T = \frac{\hat{\sigma}_{\delta_s}^2}{3T_p^3} \quad (2.88)$$

### 2.3. Covariance Measurement Update

This section discusses methods for implementing the covariance measurement update. Some of the most important of these best practices are related to factorization methods and underweighting, which are topics of enough significance to warrant their own chapters.

**2.3.1. “Stable Form” of non-Joseph Covariance Update** As Chapter 1 pointed out, only for the optimal gain and true covariance does the Joseph form of the covariance measurement update, (1.41),

$$\mathbf{P}_i^+ = (\mathbf{I} - \mathbf{K}_i \mathbf{H}_i) \mathbf{P}_i^- (\mathbf{I} - \mathbf{K}_i \mathbf{H}_i)^\top + \mathbf{K}_i \mathbf{R}_i \mathbf{K}_i^\top \quad (2.89)$$

reduce to (1.33),

$$\mathbf{P}_i^+ = \mathbf{P}_i^- - \mathbf{K}_i \mathbf{H}_i \mathbf{P}_i^- \quad (2.90)$$

While this assertion is strictly true, the cancellations that produce the above results will still incur so long as the EKF algorithm is internally consistent with truncating and approximating its various terms. The resulting “covariance” will not accurately represent  $E[\mathbf{e}_i^+ (\mathbf{e}_i^+)^\top | \mathbb{Y}_i]$ , but the fact that these truncations and approximations have produced a suboptimal gain will, in themselves, provide no computational issues. In effect, the resulting suboptimal gain remains “optimal” with respect to the internally consistent set of approximations and truncations internal to the filter.

However, even if the gain is optimal, the stability of the non-Joseph form depends on the order of multiplication, as Schmidt points out [64]. He describes a “stable form” of the non-Joseph update, given by Algorithm 2.1, which was successfully used by the Space Shuttle. Algorithm 2.1 processes each  $j$ th scalar element of the measurement vector one at a time, using only the  $j$ th column of the measurement partials matrix,  $\mathbf{h}_j$ , and the  $(j, j)$  diagonal element of the measurement noise covariance,  $r_j$ , assuming  $\mathbf{R}_i$  is a diagonal matrix. In comparison with  $\mathbf{P}_i^+ = (\mathbf{I} - \mathbf{K}_i \mathbf{H}_i) \mathbf{P}_i^-$ , use of Algorithm 2.1 also reduces the

---

**Algorithm 2.1** “Stable form” of the non-Joseph Covariance Measurement Update

---

```

 $\mathbf{P}_1 = \mathbf{P}_i^-$ 
for each scalar measurement  $j = 1$  through  $k$  do
     $\mathbf{h}_j = j$ th column of  $\mathbf{H}_i$ 
     $r_j = (j, j)$  element of  $\mathbf{R}_i$ 
     $\mathbf{b}_j = \mathbf{P}_j \mathbf{h}_j^\top$ 
     $\mathbf{k}_j = \mathbf{b}_j / (\mathbf{h}_j \mathbf{b}_j + r_j)$ 
     $\mathbf{P}_j \leftarrow \mathbf{P}_j - \mathbf{k}_j \mathbf{b}_j^\top$ 
end for
 $\mathbf{P}_i^+ = \mathbf{P}_k$ 

```

Only the non-redundant (upper or lower triangular) portions of the covariance should be updated, and then the other redundant elements set equal to the ones that have been computed.

---

computational burden from  $O(n^3)$  to  $O(n^2/2)$ , where  $n$  is the state dimension.

Although Algorithm 2.1 does not show the state update, it may also be sequentially updated as part of the iteration. However, the order in which the scalar measurements update the state can affect the outcome, if the measurement partials are computed one column at a time, corresponding with each scalar update. This may produce undesirable or even unstable outcomes. Chapter 3 will discuss such issues further.

Despite the extensive and successful flight heritage of Algorithm 2.1, it cannot guarantee numerical stability and positive definiteness of the covariance. Therefore, the recommended best practice for the covariance update is to utilize the  $UD$ -factorization, which Chapter 7 describes.

**2.3.2. Use of Consider States** It may often be the case that unobservable states are present in the system being estimated. Most commonly, such states will be parameters whose values are unknown or uncertain. Inclusion of such parameters as solve-for states in the EKF is a not a recommended practice. However, if the EKF completely ignores the uncertainty that such parameters introduce, its covariance can become overly optimistic, a condition sometimes known as “filter smugness.” One approach to addressing this problem was introduced by Schmidt [64], originally in the context of reducing the computational burden that the EKF imposed on flight computers of the 1960’s. Schmidt’s idea is essentially for the EKF to maintain a covariance containing all of the states whose uncertainties are significant enough to affect filter performance, but only to update a subset of those states. The states which are not updated in this framework are typically known as “consider” parameters, and such a filter has been called a “consider filter” or a “Schmidt-Kalman” filter. Although most commonly the state space is simply partitioned by selecting states as either solve-for or consider states, Reference 48 points out that partitioning using linear combinations of the full state space is also possible.

Following Reference 48, suppose the filter produces estimates for a subset of  $n_s$  solve-for states, out of the full state of size  $n$ . The filter does not estimate the remaining  $n_c = n - n_s$  consider states. Denote the true solve-for vector by  $\mathbf{s}(t)$ , and the true consider vector by  $\mathbf{c}(t)$ . Assume that these are linear combinations of the true states, according to the following:

$$\mathbf{s}(t) = \mathbf{S}(t)\mathbf{x}(t) \quad \text{and} \quad \mathbf{c}(t) = \mathbf{C}(t)\mathbf{x}(t) \quad (2.91)$$

where the  $n_s \times n$  matrix  $\mathbf{S}(t)$  and the  $n_c \times n$  matrix  $\mathbf{C}(t)$  are such that the matrix

$$\mathbf{M} = \begin{bmatrix} \mathbf{S} \\ \mathbf{C} \end{bmatrix} \quad (2.92)$$

is non-singular. The inverse of  $\mathbf{M}$  is partitioned into an  $n \times n_s$  matrix  $\tilde{\mathbf{S}}$  and an  $n \times n_c$  matrix  $\tilde{\mathbf{C}}$ :

$$\mathbf{M}^{-1} = \begin{bmatrix} \tilde{\mathbf{S}} & \tilde{\mathbf{C}} \end{bmatrix}. \quad (2.93)$$

The properties of the matrix inverse then lead immediately to the identities

$$\tilde{\mathbf{S}}\tilde{\mathbf{S}} = \mathbf{I}_{n_s}, \quad \tilde{\mathbf{C}}\tilde{\mathbf{C}} = \mathbf{I}_{n_c}, \quad \tilde{\mathbf{S}}\tilde{\mathbf{C}} = \mathbf{0}_{n_s \times n_c}, \quad \tilde{\mathbf{C}}\tilde{\mathbf{S}} = \mathbf{0}_{n_c \times n_s}, \quad (2.94)$$

and

$$\tilde{\mathbf{S}}\mathbf{S} + \tilde{\mathbf{C}}\mathbf{C} = \mathbf{I}_n. \quad (2.95)$$

In the usual case that the elements of the solve-for and consider vectors are merely selected and possibly permuted components of the state vector, the matrix  $\mathbf{M}$  is an orthogonal permutation matrix. In this case, and in any case for which  $\mathbf{M}$  is orthogonal, the matrices  $\tilde{\mathbf{S}}$  and  $\tilde{\mathbf{C}}$  are just the transposes of  $\mathbf{S}$  and  $\mathbf{C}$ , respectively, which makes inversion of  $\mathbf{M}$  unnecessary and simplifies many of the following equations.

It follows from Eqs. 2.91 and 2.95 that

$$\mathbf{x}(t) = \tilde{\mathbf{S}}(t)\mathbf{s}(t) + \tilde{\mathbf{C}}(t)\mathbf{c}(t). \quad (2.96)$$

Relations similar to Eq. 2.91 give the estimated solve-for vector  $\hat{\mathbf{s}}(t)$  and the assumed consider vector  $\hat{\mathbf{c}}(t)$  in terms of the estimated state  $\hat{\mathbf{x}}(t)$ . Thus, errors in the solve-for and consider states are given by

$$\mathbf{e}_s(t) = \mathbf{s}(t) - \hat{\mathbf{s}}(t) = \mathbf{S}(t)\mathbf{e}(t) \quad (2.97)$$

$$\mathbf{e}_c(t) = \mathbf{c}(t) - \hat{\mathbf{c}}(t) = \mathbf{C}(t)\mathbf{e}(t) \quad (2.98)$$

and the true error may be written in terms of the solve-for and consider errors by

$$\mathbf{e}(t) = \tilde{\mathbf{S}}(t)\mathbf{e}_s(t) + \tilde{\mathbf{C}}(t)\mathbf{e}_c(t). \quad (2.99)$$

In terms of this notation, the EKF update has the form

$$\hat{\mathbf{s}}_i^+ = \hat{\mathbf{s}}_i^- + \mathbf{K}_i\mathbf{r}_i^- \quad (2.100)$$

where

$$\mathbf{r}_i^- = \mathbf{y}_i - \hat{\mathbf{y}}_i^- = \mathbf{H}_i\mathbf{e}_i^- + \mathbf{v}_i, \quad (2.101)$$

and the subscript  $i$  is a shorthand for the time argument  $t_i$ . The usual EKF will not contain the full covariance, but only its solve-for part

$$\mathbf{P}_{ss}(t) = \mathbf{E}[\mathbf{e}_s(t)\mathbf{e}_s^T(t)] \quad (2.102)$$

By contrast, the Schmidt-Kalman filter will use the full covariance,  $\mathbf{P}(t)$ . In the usual case, the Kalman gain is given by

$$\mathbf{K}_i = \mathbf{P}_{ssi}^- \mathbf{H}_{si}^T [\mathbf{H}_{si} \mathbf{P}_{ssi}^- \mathbf{H}_{si}^T + \mathbf{R}_i]^{-1} \quad (2.103)$$

where

$$\mathbf{H}_{si} = \mathbf{H}_i \tilde{\mathbf{S}}_i \quad (2.104)$$

In the Schmidt-Kalman case,

$$\mathbf{K}_i = \mathbf{S}_i \mathbf{P}_i^- \mathbf{H}_i^T [\mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^T + \mathbf{R}_i]^{-1} \quad (2.105)$$

Thus, the Schmidt-Kalman gain matrix is computed from the full covariance, but only applies measurement innovations to the solve-for states.

#### 2.4. Sigma-Point Methods

Before concluding this Chapter, it is worth noting that a promising method for propagating and updating the covariance that is coming into greater use and acceptance within the navigation community is the use of sigma-point methods, also known as “unscented” transforms. A subsequent chapter covering Advanced Topics will cover these topics in a section on sigma-point filtering.

## CHAPTER 3

# Processing Measurements

Contributed by Christopher N. D’Souza and J. Russell Carpenter

This chapter will discuss how to handle real-world measurement processing issues. In particular, being able to handle measurements that aren’t synchronous is of paramount importance to running filters in a real-time environment. As well, the performance of navigation filters which have nonlinear measurement models are susceptible to divergence depending on the order of processing of measurements which occur at the same epoch. Therefore, a technique which provides invariance to measurement processing is detailed. A technique for processing correlated measurements is presented, and brief comments on filter cascading and processing of inertial data are offered.

### 3.1. Measurement Latency

In general, the measurement time tags are not going to be equal to the current filter epoch time,  $t_k$ . To state it another way, the measurements do not come in at the current filter time. Rather, they may be latent by up to  $p$  seconds. Thus, a situation will arise where the filter has propagated its state and covariance to time  $t = t_k$  from time  $t = t_{k-1}$ , and is subsequently given a measurement to be filtered (denoted by subscript  $m$ ) that corresponds to the time  $t = t_m$ , where

$$t_m \leq t_k \tag{3.1}$$

If  $\Delta t = t_m - t_k$  is not insignificant, the time difference between the measurement and the filter state and covariance will need to be accounted for during filtering in order to accurately process the measurement. This can be done in much the same way a batch filter operates (see pages 196-197 of Tapley [69]). If the measurement at time  $t = t_m$  is denoted as  $y_m$ , the nominal filter state at that time is given by  $\mathbf{X}_m^* \equiv \mathbf{X}^*(t_m)$  (\* denotes the nominal), and the measurement model is denoted as  $h_m(\mathbf{X}_m, t_m)$ , then one can expand the measurement model to first order about the nominal filter state to get

$$h_m(\mathbf{X}_m, t_m) = h_m(\mathbf{X}_m^*, t_m) + \mathbf{H}_m \mathbf{x}_m + \nu_m \tag{3.2}$$

where  $\mathbf{x}_m = \mathbf{X}_m - \mathbf{X}_m^*$  and  $\mathbf{H}_m$  is defined as

$$\mathbf{H}_m \triangleq \left( \frac{\partial h_m(\mathbf{X}, t_m)}{\partial \mathbf{X}} \right)_{\mathbf{x}=\mathbf{x}_m^*} \tag{3.3}$$

The perturbed state at time  $t_m$  can be written in terms of the state at time  $t_k$  as follows

$$\mathbf{x}_m = \Phi(t_m, t_k) \mathbf{x}_k + \Gamma_m \mathbf{w}_m \tag{3.4}$$

so we can compute the measurement as

$$h_m(\mathbf{X}_m, t_m) = h_m(\mathbf{X}_m^*, t_m) + \mathbf{H}_m \Phi(t_m, t_k) \mathbf{x}_k - \mathbf{H}_m \Gamma_m \mathbf{w}_m + \nu_m \tag{3.5}$$

where the measurement noise has the characteristics  $E[\nu_m] = 0$  and  $E[\nu_m^2] = R_m$ , the state process noise from  $t = t_m$  to  $t = t_k$  has the characteristics  $E[\mathbf{w}_m] = \mathbf{0}$  and  $E[\mathbf{w}_m \mathbf{w}_m^T] = \mathbf{Q}_m$ , and the state deviation is given by

$$\mathbf{x}_k = \delta \mathbf{X}_k = \mathbf{X}_k - \mathbf{X}_k^* \quad (3.6)$$

(Note that the effect of the state process noise would be to increase the measurement noise variance. However, because the process noise term is very small over time periods of a few seconds, it can safely be neglected for the remainder of this analysis.)

Upon taking the conditional expectation of the measurement equation and rearranging, the scalar residual of the measurement is given by

$$y_m - \tilde{\mathbf{H}}_m \hat{\mathbf{x}}_k(-) = y_m - h_m(\mathbf{X}_m^*, t_m) - \mathbf{H}_m \Phi(t_m, t_k) \hat{\mathbf{x}}_k(-) \quad (3.7)$$

where  $\hat{\cdot}$  denotes an estimated value,

$$\begin{aligned} y_m &= y_m - h_m(\mathbf{X}_m^*, t_m) \\ \hat{\mathbf{x}}_k(-) &= \tilde{\mathbf{X}}_k(-) - \mathbf{X}_k^* \end{aligned} \quad (3.8)$$

The measurement partials that are used in the update, which map the measurement to the state at time  $t = t_k$ , are given by

$$\tilde{\mathbf{H}}_m = \mathbf{H}_m \Phi(t_m, t_k) \quad (3.9)$$

Eq. 3.9 was derived by noting that

$$\begin{aligned} \mathbf{H}_m \hat{\mathbf{x}}_m &= \mathbf{H}_m \Phi(t_m, t_k) \hat{\mathbf{x}}_k \\ &= \tilde{\mathbf{H}}_m \hat{\mathbf{x}}_k \end{aligned} \quad (3.10)$$

From the above discussion, it is evident that the unknown quantities needed to update the state at time  $t = t_k$  with a measurement from time  $t = t_m$  are the nominal state at the measurement time,  $\mathbf{X}_m^*$ , and the state transition matrix relating the two times,  $\Phi(t_m, t_k)$ . Given those values,  $h_m(\mathbf{X}_m^*, t_m)$  and  $\tilde{\mathbf{H}}_m$  can be calculated.

Thus the nominal state at the measurement time is calculated by back-propagating the filter state from time  $t_k$  to time  $t_m$  using buffered IMU data. The same thing is done to calculate the required state transition matrix. The same propagation algorithms used in forward propagation ought to be utilized for the back-propagation, with the exception that the smaller time step allows for a 1<sup>st</sup>-order approximation of the matrix exponential used to update the state transition matrix.

### 3.2. Invariance to the Order of Measurement Processing

It has long been known that the performance of an EKF is dependent on the order in which one processes measurements. This is of particular import in the case when there is powerful measurement coupled with a large *a priori* error. The state (and covariance) update will be large, very likely out of the linear range. Subsequent measurements which are processed may well be outside the residual edit thresholds, and hence will be rejected. In order to remedy this, we employ a hybrid Linear/Extended Kalman Filter measurement update. Recall that in an Extended Kalman Filter, the state is updated / relinearized / rectified after each measurement is processed, regardless of whether the measurements occur at the same time. Hence, the solution is highly dependent on the *order* in which the measurements are processed. This is not a desirable situation in which to be.

We obviate this difficulty simply by not updating the state until *all* the measurements at a given time are processed. We accumulate the state updates in state deviations  $\mathbf{x}$ , using

Algorithm 3.1. This algorithm makes use of the fact that, in the absence of process noise, a batch/least squares algorithm is mathematically equivalent to a linear Kalman Filter [25]. Algorithm 3.1 is a recommended *best practice*. Algorithm 3.1 may readily be combined with

---

**Algorithm 3.1** Measurement Update Invariant to Order of Processing

---

```

for each (scalar) measurement  $j = 1$  through  $k$  do
   $y_j = Y_j - h_j(\mathbf{X}_m^*, t_m)$ 
   $\mathbf{H}_j = \frac{\partial h_j}{\partial \mathbf{X}}(\mathbf{X}_m^*, t_m)$ 
   $\mathbf{K}_j = \mathbf{P}_j \mathbf{H}_j^\top (\mathbf{H}_j \mathbf{P}_j \mathbf{H}_j^\top + R_j)^{-1}$ 
   $\hat{\mathbf{x}}_j \leftarrow \hat{\mathbf{x}}_j - \mathbf{K}_j (y_j - \mathbf{H}_j \hat{\mathbf{x}}_j)$ 
   $\mathbf{P}_j \leftarrow (\mathbf{I} - \mathbf{K}_j \mathbf{H}_j) \mathbf{P}_j (\mathbf{I} - \mathbf{K}_j \mathbf{H}_j)^\top + \mathbf{K}_j R_j \mathbf{K}_j^\top$ 
end for
 $\mathbf{X}_m \leftarrow \mathbf{X}_m^* + \hat{\mathbf{x}}_j$ 

```

---

the residual mapping approach described above when the measurements are asynchronous. Algorithm 3.1 may also be readily combined with Algorithm 2.1, for cases in which the preferred factorized covariance methods are precluded.

### 3.3. Processing Vector Measurements

If the UDU factorization is used, measurements need to be processed as scalars. If the vector measurements are correlated, one option is to assume they are uncorrelated and ignore the correlations between the measurements.

However, there is a better alternative. Given the measurement equation ( $\mathbf{Y}_j = \mathbf{H}_j(\mathbf{X}_j, t_j) + \boldsymbol{\nu}_j$ ) with measurement error covariance matrix,  $\mathbf{R}_i$ , first decompose the matrix with a Cholesky factorization as

$$\mathbf{R}_i = \text{E}[\boldsymbol{\nu}_j \boldsymbol{\nu}_j^\top] = \mathbf{R}_i^{1/2} \mathbf{R}_i^{\top/2} \quad (3.11)$$

and premultiply the measurement equation by  $\mathbf{R}_i^{-1/2}$  to yield

$$\tilde{\mathbf{Y}}_j = \tilde{\mathbf{H}}_j(\mathbf{X}_j, t_j) + \tilde{\boldsymbol{\nu}}_j \quad (3.12)$$

with

$$\tilde{\mathbf{Y}}_j = \mathbf{R}_i^{-1/2} \mathbf{Y}_j \quad (3.13)$$

$$\tilde{\mathbf{H}}_j(\mathbf{X}_j, t_j) = \mathbf{R}_i^{-1/2} \mathbf{H}_j(\mathbf{X}_j, t_j) \quad (3.14)$$

$$\tilde{\boldsymbol{\nu}}_j = \mathbf{R}_i^{-1/2} \boldsymbol{\nu}_j \quad (3.15)$$

so that  $\text{E}[\tilde{\boldsymbol{\nu}}_j \tilde{\boldsymbol{\nu}}_j] = \mathbf{I}$ . Thus, the new measurement equation has errors which are now decorrelated.

Alternatively, one can decompose the  $m \times m$  measurement error covariance matrix with a UDU decomposition as  $\mathbf{R}_i = \mathbf{U}_{R_i} \mathbf{D}_{R_i} \mathbf{U}_{R_i}^\top$  so that using a similar reasoning, we premultiply the measurement equation by  $\mathbf{U}_{R_i}^{-1}$  so that in this case

$$\tilde{\mathbf{Y}}_j = \mathbf{U}_{R_i}^{-1} \mathbf{Y}_j \quad (3.16)$$

$$\tilde{\mathbf{H}}_j(\mathbf{X}_j, t_j) = \mathbf{U}_{R_i}^{-1} \mathbf{H}_j(\mathbf{X}_j, t_j) \quad (3.17)$$

$$\tilde{\boldsymbol{\nu}}_j = \mathbf{U}_{R_i}^{-1} \boldsymbol{\nu}_j \quad (3.18)$$

so that  $\text{E}[\tilde{\boldsymbol{\nu}}_j \tilde{\boldsymbol{\nu}}_j] = \mathbf{D}_{R_i}$  where  $\mathbf{D}_{R_i}$  is a diagonal matrix and, as in the case of the Cholesky decomposition, the new measurement model has decorrelated measurement errors.

### 3.4. Filter Cascades

It is often tempting to ingest the output of one navigation filter as a measurement into another, “downstream” filter. The problem with this approach is that the output of the upstream filter will be time-correlated. This requires the downstream filter to model the correlation structure induced by the upstream filter. Such an exercise has rarely led to useful and robust designs, and therefore such filter cascades are strongly discouraged. Nonetheless, some forms of pre-filtering noisy high-rate data, such as carrier-smoothing, have been usefully employed.

### 3.5. Use of Data from Inertial Sensors

Inertial measurement units (IMUs), consisting of gyros and accelerometers, sense rotational and translational accelerations. While in principle these high-rate data could be processed as observations in the navigation filter, it is often sufficient instead to use this data in *model replacement mode*, which Brown and Hwang [3] compare to *complementary* filtering. In this approach, the sensed accelerations are fed forward as deterministic inputs to the navigation filter’s dynamics model. Biases affecting the IMU data usually must be estimated by the filter. Thresholding the IMU accelerations is also usually necessary to avoid the introduction of unfiltered IMU noise into the filter state propagation. Chapter 6 describes best practices for modeling the structures associated with IMU data.

## Measurement Underweighting

Contributed by Renato Zanetti

### 4.1. Introduction

Given an  $m$ -dimensional random measurement  $\mathbf{y}$  which is somehow related to an unknown,  $n$ -dimensional random vector  $\mathbf{x}$  the family of affine estimators of  $\mathbf{x}$  from  $\mathbf{y}$  is

$$\hat{\mathbf{x}} = \mathbf{a} + \mathbf{K} \mathbf{y} \quad (4.1)$$

where  $\mathbf{a} \in \mathfrak{R}^n$  and  $\mathbf{K} \in \mathfrak{R}^{n \times m}$ . The optimal, in a Minimum Mean Square Error sense, affine estimator has

$$\mathbf{K} = \mathbf{P}_{xy} \mathbf{P}_{yy}^{-1} \quad (4.2)$$

$$\mathbf{a} = \mathbf{E}[\mathbf{x}] - \mathbf{K} \mathbf{E}[\mathbf{y}] \quad (4.3)$$

where

$$\mathbf{P}_{xy} = \mathbf{E} \left[ \left( \mathbf{x} - \mathbf{E}[\mathbf{x}] \right) \left( \mathbf{y} - \mathbf{E}[\mathbf{y}] \right)^{\top} \right] \quad (4.4)$$

$$\mathbf{P}_{yy} = \mathbf{E} \left[ \left( \mathbf{y} - \mathbf{E}[\mathbf{y}] \right) \left( \mathbf{y} - \mathbf{E}[\mathbf{y}] \right)^{\top} \right] \quad (4.5)$$

In the presence of nonlinear measurements of the state,

$$\mathbf{y} = \mathbf{h}(\mathbf{x}) + \mathbf{v} \quad (4.6)$$

(where  $\mathbf{v}$  is zero-mean measurement noise) the extended Kalman filter (EKF) [20] approximates all moments of  $\mathbf{y}$  by linearization of the measurement function centered on the mean of  $\mathbf{x}$ . This methodology has proven very effective and produces very satisfactory results in most cases. Approaches other than the EKF exist, for example the Unscented Kalman Filter [32] approximates the same quantities via stochastic linearization using a deterministic set of Sigma Points. High order truncations of the Taylor series are also possible. Underweighting [38, 40] is an ad-hoc technique to compensating for nonlinearities in the measurement models that are neglected by the EKF and successfully flew on the Space Shuttle and on Orion Exploration Flight Test 1.

The commonly implemented method for the underweighting of measurements for human space navigation was introduced by Lear [39] for the Space Shuttle navigation system. In 1966 Denham and Pines showed the possible inadequacy of the linearization approximation when the effect of measurement nonlinearity is comparable to the measurement error [12]. To compensate for the nonlinearity Denham and Pines proposed to increase the measurement noise covariance by a constant amount. In the early seventies, in anticipation of Shuttle flights, Lear and others developed relationships which accounted for the second-order effects in the measurements [40]. It was noted that in situations involving large state

errors and very precise measurements, application of the standard extended Kalman filter mechanization leads to conditions in which the state estimation error covariance decreases more rapidly than the actual state errors. Consequently the extended Kalman filter begins to ignore new measurements even when the measurement residual is relatively large. Underweighting was introduced to slow down the convergence of the state estimation error covariance thereby addressing the situation in which the error covariance becomes overly optimistic with respect to the actual state errors. The original work on the application of second-order correction terms led to the determination of the underweighting method by trial-and-error [39].

More recently, studies on the effects of nonlinearity in sensor fusion problems with application to relative navigation have produced a so-called “bump-up” factor. [16, 44, 56, 58]. While Ferguson [16] seems to initiate the use of the bump-up factor, the problem of mitigating filter divergence was more fully studied by Plinval [58] and subsequently by Mandic [44]. Mandic generalized Plinval’s bump-up factor to allow flexibility and notes that the value selected influences the steady-state convergence of the filter. In essence, it was found that a larger factor keeps the filter from converging to the level that a lower factor would permit. This finding prompted Mandic to propose a two-step algorithm in which the bump-up factor is applied for a certain number of measurements only, upon which the factor was completely turned off. Finally, Perea, et al. [56] summarize the findings of the previous works and introduce several ways of computing the applied factor. In all cases, the bump-up factor amounts in application to the underweighting factor introduced in Lear [39]. Save for the two-step procedure of Mandic, the bump-up factor is allowed to persistently affect the Kalman gain which directly influences the obtainable steady-state covariance. Effectively, the ability to remove the underweighting factor autonomously and under some convergence condition was not introduced.

The work of Lear is not well known as it is only documented in internal NASA memos [39, 40]. Kriegsman and Tau [38] mention underweighting in their 1975 Shuttle navigation paper without a detailed explanation of the technique.

## 4.2. Nonlinear Effects and the Need for Underweighting

We review briefly the three state estimate update approaches assuming a linear time-varying measurement model leading to the classical Kalman filter, a nonlinear measurement model with first-order linearization approximations leading the widely used extended Kalman filter, and a nonlinear model with second-order approximations leading to the second-order extended Kalman filter.

### 4.2.1. Linear Measurement Model and the Classical Kalman Filter Update

Let’s briefly recap the linear Kalman filter, the measurement model is

$$\mathbf{y}_i = \mathbf{H}_i \mathbf{x}_i + \mathbf{v}_i, \quad (4.7)$$

where  $\mathbf{y}_i \in \mathbb{R}^m$  are the  $m$  measurements at each time  $t_i$ ,  $\mathbf{x}_i \in \mathbb{R}^n$  is the  $n$ -dimensional state vector,  $\mathbf{H}_i \in \mathbb{R}^{m \times n}$  is the known measurement mapping matrix,  $\mathbf{v}_i \in \mathbb{R}^m$  is modeled as a zero-mean white-noise sequence with  $E[\mathbf{v}_i] = 0, \forall k$  and  $E[\mathbf{v}_i \mathbf{v}_j^T] = \mathbf{R}_i \delta_{kj}$  where  $\mathbf{R}_i > 0 \forall k$  and  $\delta_{kj} = 1$  when  $k = j$  and  $\delta_{kj} = 0$  when  $k \neq j$ . The Kalman filter state update algorithm provides an optimal blending of the *a priori* estimate  $\hat{\mathbf{x}}_i^-$  and the measurement  $\mathbf{y}_i$  at time  $t_i$  to obtain the *a posteriori* state estimate  $\hat{\mathbf{x}}_i^+$  via

$$\hat{\mathbf{x}}_i^+ = \hat{\mathbf{x}}_i^- + \mathbf{K}_i [\mathbf{y}_i - \mathbf{H}_i \hat{\mathbf{x}}_i^-], \quad (4.8)$$

where the superscript  $-$  denotes *a priori* and  $+$  denotes *a posteriori*.

Defining the *a priori* estimation error as  $\mathbf{e}_i^- = \mathbf{x}_i - \hat{\mathbf{x}}_i^-$  and the *a posteriori* estimation error as  $\mathbf{e}_i^+ = \mathbf{x}_i - \hat{\mathbf{x}}_i^+$  and assuming these errors to be zero mean, the associated symmetric, positive definite *a priori* and *a posteriori* estimation error covariances are  $\mathbf{P}_i^- = \mathbb{E}[\mathbf{e}_i^- (\mathbf{e}_i^-)^\top]$  and  $\mathbf{P}_i^+ = \mathbb{E}[\mathbf{e}_i^+ (\mathbf{e}_i^+)^\top]$ , respectively. Using Eq. (4.8) and the definitions of the state estimation errors and error covariances, we obtain the *a posteriori* state estimation error covariance via the well-known Joseph formula

$$\mathbf{P}_i^+ = [\mathbf{I} - \mathbf{K}_i \mathbf{H}_i] \mathbf{P}_i^- [\mathbf{I} - \mathbf{K}_i \mathbf{H}_i]^\top + \mathbf{K}_i \mathbf{R}_i \mathbf{K}_i^\top, \quad (4.9)$$

which is valid for any  $\mathbf{K}_i$ . If the gain  $\mathbf{K}_i$  is chosen so as to minimize the trace of the *a posteriori* estimation error, we call that gain the Kalman gain which is given by

$$\mathbf{K}_i = \mathbf{P}_i \mathbf{H}_i^\top [\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^\top + \mathbf{R}_i]^{-1}. \quad (4.10)$$

The trace of the state estimation error covariance is generally not a norm but is equivalent to the nuclear norm (the matrix Shatten 1-norm) for symmetric semi-positive matrices. If the gain given in Eq. (4.10) is applied to the state estimation error covariance of Eq. (4.9), then the update equation can be rewritten after some manipulation as

$$\mathbf{P}_i^+ = [\mathbf{I} - \mathbf{K}_i \mathbf{H}_i] \mathbf{P}_i^-, \quad (4.11)$$

or equivalently,

$$\mathbf{P}_i^+ = \mathbf{P}_i^- - \mathbf{K}_i [\mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top + \mathbf{R}_i] \mathbf{K}_i^\top. \quad (4.12)$$

Under the assumptions of the Kalman filter development (linear, time-varying measurement model with a zero-mean white-noise sequence corrupting the measurements, unbiased *a priori* estimation errors, known dynamics and measurement models, etc.), the state estimate and state estimation error covariance updates are optimal and we expect no filter divergence issues. The estimation error covariance will remain positive definite for all  $t_i$  and the estimation error covariance will be consistent with the true errors. In practice, the measurements are usually nonlinear functions of the state leading to a variety of engineering solutions, that must be carefully designed to ensure acceptable state estimation performance. Underweighting is one such method to improve the performance of the extended Kalman filter in practical settings.

**4.2.2. Nonlinear Measurement Model and the Extended Kalman Filter Update** In the nonlinear setting, consider the measurement model given by

$$\mathbf{y}_i = \mathbf{h}(\mathbf{x}_i, t_i) + \mathbf{v}_i, \quad (4.13)$$

where  $\mathbf{h}(\mathbf{x}_i) \in \mathbb{R}^m$  is a vector-valued differentiable nonlinear function of the state vector  $\mathbf{x}_i \in \mathbb{R}^n$ . The idea behind the extended Kalman filter (EKF) is to utilize Taylor series approximations to obtain linearized models in such a fashion that the EKF state update algorithm has the same general form as the Kalman filter.

$$\mathbf{y}_i \simeq \mathbf{h}(\hat{\mathbf{x}}_i^-, t_i) + \mathbf{H}_i \mathbf{e}_i^- + \mathbf{v}_i, \quad (4.14)$$

where

$$\mathbf{H}_i \triangleq \left[ \frac{\partial \mathbf{h}(\mathbf{x}_i, t_i)}{\partial \mathbf{x}_i} \Big|_{\mathbf{x}_i = \hat{\mathbf{x}}_i^-} \right]. \quad (4.15)$$

Since the estimation error is (approximately) zero mean and the measurement noise is zero mean, it follows that

$$\mathbb{E}[\mathbf{y}_i] \simeq \mathbf{h}(\hat{\mathbf{x}}_i^-, t_i), \quad (4.16)$$

all expectations are conditioned on past measurements and we find that the state estimate update is given by [20]

$$\hat{\mathbf{x}}_i^+ = \hat{\mathbf{x}}_i^- + \mathbf{K}_i[\mathbf{y}_i - \mathbf{h}(\hat{\mathbf{x}}_i^-)]. \quad (4.17)$$

Similarly, the measurement residual is given by

$$\mathbf{r}_i = \mathbf{y}_i - \mathbf{h}(\hat{\mathbf{x}}_i^-) \simeq \mathbf{H}_i \mathbf{e}_i^- + \mathbf{v}_i, \quad (4.18)$$

Computing the measurement residual covariance  $\mathbb{E}[\mathbf{r}_i \mathbf{r}_i^\top]$  yields

$$\mathbf{W}_i = \mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top + \mathbf{R}_i. \quad (4.19)$$

The state estimation error covariance and Kalman gain are the same as in Eqs. (4.9) and (4.10), respectively, with  $\mathbf{H}_i$  given as in Eq. (4.15). The state estimation error covariances in the forms shown in Eqs. (4.11) and (4.12) also hold in the nonlinear setting with  $\mathbf{H}_i$  as in Eq. (4.15).

From Eqs. (4.12) and (4.17), it is seen that reducing the Kalman gain leads to a smaller update in both the state estimation error covariance and the state estimate, respectively. Reducing the gain and hence the update is the essence of underweighting and the need for this adjustment is illuminated in the following discussion.

Adopting the viewpoint that the state estimation error covariance matrix represents the level of uncertainty in the state estimate, we expect that when we process a measurement (adding new information) that the uncertainty would decrease (or at least, not increase). This is, in fact, the case and can be seen in Eq. (4.12). Under the assumption that the symmetric matrices  $\mathbf{P}_i^- > 0$  and  $\mathbf{R}_i > 0$ , it follows that

$$\mathbf{K}_i[\mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top + \mathbf{R}_i] \mathbf{K}_i^\top \geq 0, \quad (4.20)$$

and we can find a number  $\alpha_i \geq 0$  at each time  $t_i$  such that

$$\mathbf{P}_i^- - \mathbf{P}_i^+ \geq \alpha_i \mathbf{I}, \quad (4.21)$$

which shows that the  $\mathbf{P}_i^- - \mathbf{P}_i^+$  is non-negative definite. The same argument can be made from the viewpoint of comparing the trace (or the matrix norm) of the *a posteriori* and *a priori* state estimation error covariances. As each new measurement is processed by the EKF, we expect the uncertainty in the estimation error to decrease. The question is, does the *a posteriori* uncertainty as computed by the EKF represent the actual uncertainty, or in other words, is the state estimation error covariance matrix always consistent with the actual state errors? In the nonlinear setting when there is a large *a priori* uncertainty in the state estimate and a very accurate measurement, it can happen that the state estimation error covariance reduction at the measurement update is too large. Underweighting is a method to address this situation by limiting the magnitude of the state estimation error covariance update with the goal of retaining consistency of the filter covariance and the actual state estimation error through situations of high nonlinearity of the measurements.

Pre- and post-multiplying the *a posteriori* state estimation error covariance in Eq. (4.12) by  $\mathbf{H}_i$  and  $\mathbf{H}_i^\top$ , respectively, yields (after some manipulation)

$$\mathbf{H}_i \mathbf{P}_i^+ \mathbf{H}_i^\top = \mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top (\mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top + \mathbf{R}_i)^{-1} \mathbf{R}_i. \quad (4.22)$$

In Eq. (4.22), we see that if  $\mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top \gg \mathbf{R}_i$  then it follows that

$$\mathbf{H}_i \mathbf{P}_i^+ \mathbf{H}_i^\top \simeq \mathbf{R}_i. \quad (4.23)$$

The result in Eq. (4.23) is of fundamental importance and is the motivation behind underweighting. What this equations express is the fact that when the *a priori* estimated

state uncertainty  $\mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top$  is much larger than the measurement error covariance  $\mathbf{R}_i$ , the Kalman filter largely neglects the prior information and relies heavily on the measurement. Therefore the *a posteriori* estimated state uncertainty  $\mathbf{H}_i \mathbf{P}_i^+ \mathbf{H}_i^\top$  is approximately equal to  $\mathbf{R}_i$ .

**4.2.3. Nonlinear Measurement Model and the 2nd-Order Kalman Filter Update** Eq. (4.14) truncates the Taylor series expansion of the nonlinear measurement function to first order, carrying it to second order we obtain

$$\mathbf{y}_i \simeq \mathbf{h}(\hat{\mathbf{x}}_i^-, t_i) + \mathbf{H}_i \mathbf{e}_i^- + \sum_{k=1}^m \left( \left. \frac{\partial^2 \mathbf{h}(\mathbf{x}_i, t_i)}{\partial \mathbf{x}_i \partial x_i(k)} \right|_{\mathbf{x}_i = \hat{\mathbf{x}}_i^-} \mathbf{e}_i(k) \mathbf{e}_i \right) + \mathbf{v}_i \quad (4.24)$$

$$= \mathbf{h}(\hat{\mathbf{x}}_i^-, t_i) + \mathbf{H}_i \mathbf{e}_i^- + \mathbf{b}_i + \mathbf{v}_i, \quad (4.25)$$

where  $e_i(k)$  and  $x_i(k)$  are the  $k$ -th elements of vectors  $\mathbf{e}_i$  and  $\mathbf{x}_i$ , respectively. The expected value of the measurement now includes contributions from the second order terms, denoted as  $\hat{\mathbf{b}}_i$

$$\mathbb{E}[\mathbf{y}_i] \simeq \mathbf{h}(\hat{\mathbf{x}}_i^-, t_i) + \hat{\mathbf{b}}_i \quad (4.26)$$

Define

$$\mathbf{H}_{i,k}^\top \triangleq \left[ \left. \frac{\partial^2 h_i(\mathbf{x}_i)}{\partial \mathbf{x}_i \partial x_i^\top} \right|_{\mathbf{x}_i = \hat{\mathbf{x}}_i^-} \right],$$

where  $h_i(\mathbf{x}_i)$  is the  $k$ -th component of  $\mathbf{h}(\mathbf{x}_i)$ . Then the  $k$ -th component of  $\mathbf{b}_i$  is given by

$$b_{i,k} = \frac{1}{2} (\mathbf{e}_i^-)^\top \mathbf{H}_{i,k}^\top \mathbf{e}_i^- = \frac{1}{2} \text{tr}(\mathbf{H}_{i,k}^\top \mathbf{e}_i^- (\mathbf{e}_i^-)^\top). \quad (4.27)$$

where  $\text{tr}$  denotes the trace. To keep the filter unbiased the  $k$ -th component of  $\hat{\mathbf{b}}_i$  is given by

$$\hat{b}_{i,k} = 1/2 \text{tr}(\mathbf{H}_{i,k}^\top \mathbf{P}_i^-).$$

The measurement residual is

$$\mathbf{r}_i = \mathbf{y}_i - \mathbb{E}[\mathbf{y}_i] \quad (4.28)$$

Expanding Eq. (4.28), the  $k$ -th component of the residual is obtained to be

$$r_{i,k} = \mathbf{h}_{i,k}^\top \mathbf{e}_i^- + 1/2 \text{tr}(\mathbf{H}_i^\top(t_i) \mathbf{e}_i^- (\mathbf{e}_i^-)^\top) - 1/2 \text{tr}(\mathbf{H}_i^\top(t_i) \mathbf{P}_i^-) + v_{i,k}, \quad (4.29)$$

where  $\mathbf{h}_{i,k}^\top$  is the  $ik$ -th row of the measurement Jacobian and  $v_{i,k}$  is the  $k$ -th component of the measurement noise  $\mathbf{v}_i$ . Computing the measurement residual covariance  $\mathbb{E}[\mathbf{r}_i \mathbf{r}_i^\top]$  yields

$$\mathbf{W}_i = \mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top + \mathbf{B}_i + \mathbf{R}_i, \quad (4.30)$$

where matrix  $\mathbf{B}_i$  is the contribution of the second order effects and its  $(kj)$ -th component is given by

$$B_{i,kj} \triangleq 1/4 \mathbb{E}[\text{tr}(\mathbf{H}_i^\top(t_i) \mathbf{e}_i^- (\mathbf{e}_i^-)^\top) \text{tr}(\mathbf{H}_i^\top(t_i) \mathbf{e}_i^- (\mathbf{e}_i^-)^\top)] \\ - 1/4 \text{tr}(\mathbf{H}_i^\top(t_i) \mathbf{P}_i^-) \text{tr}(\mathbf{H}_i^\top(t_i) \mathbf{P}_i^-).$$

where it was assumed that the third order central moments are all zeros. Assuming the prior estimation error is distributed as a zero-mean gaussian distribution with covariance  $\mathbf{P}_i^-$ , the  $ij^{\text{th}}$  component of  $\mathbf{B}_i$  is given by

$$B_{i,kj} = \frac{1}{2} \text{tr}(\mathbf{H}_{i,k}^\top \mathbf{P}_i^- \mathbf{H}_{i,j}^\top \mathbf{P}_i^-). \quad (4.31)$$

Comparing the measurement residual covariance for the EKF in Eq. (4.19) with the measurement residual covariance for the second-order filter in Eq. (4.30), we see that when

the nonlinearities lead to significant second-order terms which should not be neglected, then the EKF tends to provide residual covariance estimates that are not consistent with the actual errors. Typically, we address this by tuning the EKF using  $\mathbf{R}_i$  as a parameter to be tweaked. If the contribution of the *a priori* estimation error  $\mathbf{H}_i\mathbf{P}_i^-\mathbf{H}_i^\top$  to the residuals covariance is much larger than the contribution of the measurement error  $\mathbf{R}_i$ , the EKF algorithm will produce  $\mathbf{H}_i\mathbf{P}_i^+\mathbf{H}_i^\top \simeq \mathbf{R}_i$ . If  $\mathbf{B}_i$  is of comparable magnitude to  $\mathbf{R}_i$  then the actual covariance of the posterior measurement estimate should be  $\mathbf{H}_i\mathbf{P}_i^+\mathbf{H}_i^\top \simeq \mathbf{R}_i + \mathbf{B}_i$ . Therefore, a large underestimation of the *a posteriori* covariance can occur in the presence of nonlinearities when the estimated measurement error covariance is much larger than the measurement error covariance.

The covariance update is given by the modified Gaussian second order filter update [30]

$$\mathbf{P}_i^+ = \mathbf{P}_i^- - \mathbf{H}_i\mathbf{P}_i^-\mathbf{W}_i^{-1}(\mathbf{H}_i\mathbf{P}_i^-)^\top, \quad (4.32)$$

where the residual covariance  $\mathbf{W}_i$  is given by Eq. (4.30).

### 4.3. Underweighting Measurements

In the prior section we saw that when “large” values of  $\mathbf{H}_i\mathbf{P}_i^-\mathbf{H}_i^\top$  exist (or similarly, large values of  $\mathbf{P}_i^-$ ), and possibly “small” values of  $\mathbf{R}_i$ , the EKF is at risk underestimating the posterior estimation error covariance matrix. We must repeat that this can only happen in the presence of “large” nonlinearities. The larger  $\mathbf{P}_i^-$ , the larger the domain of possible values of the true state  $\mathbf{x}$ , hence the more likely the higher order terms of the expansion of the nonlinear measurement functions will become relevant. If a measurement function is largely non-linear, but the prior estimate is very precise, the EKF algorithm and linearization are likely sufficiently accurate since:

- (1) The posterior measurement will rely heavily on the prior and rely less on the measurement
- (2) Since the error is small, while the Hessian matrix might be relatively large, the actual contributions of the second order effects is likely to remain small

Underweighting is the process of modifying the residual covariance to reduce the update and compensate for the second-order effects described above. In this section, we describe three common methods for performing underweighting with the EKF algorithm.

**4.3.1. Additive Compensation Method** The most straightforward underweighting scheme is to add an underweighting factor  $\mathbf{U}_i$  as

$$\mathbf{W}_i = \mathbf{H}_i\mathbf{P}_i^-\mathbf{H}_i^\top + \mathbf{R}_i + \mathbf{U}_i. \quad (4.33)$$

With the Kalman gain given by

$$\mathbf{K}_i = \mathbf{P}_i^-\mathbf{H}_i^\top\mathbf{W}_i^{-1}, \quad (4.34)$$

we see that the symmetric, positive-definite underweighting factor  $\mathbf{U}_i$  decreases the Kalman gain, thereby reducing the state estimate and state estimation error covariance updates. One choice is to select  $\mathbf{U}_i = \mathbf{B}_i$ , which is the contribution to the covariance assuming the prior distribution of the estimation error is Gaussian. The advantage of this choice is its theoretical foundation based on analyzing the second-order terms of the Taylor series expansions. The disadvantages include higher computational costs to calculate the second-order partials and the reliance on the assumption that the estimation errors possess Gaussian distributions. In practical applications, the matrix  $\mathbf{U}_i$  needs to be tuned appropriately for acceptable overall performance of the EKF. The process of tuning a positive definite matrix is less obvious than tuning a single scalar parameter.

**4.3.2. Scaling the Measurement Error Covariance** Another possible underweighting approach is to scale the measurement noise by choosing

$$\mathbf{U}_i = \gamma \mathbf{R}_i, \quad (4.35)$$

where  $\gamma > 0$  is a scalar parameter selected in the design process. This approach has been successfully used [31]; however, it is not recommended from both a conceptual and a practical reason. Recalling that the underweighting is necessary because of neglecting the second-order terms of the Taylor series expansion of the measurement function, it seems more natural to express the underweighting as a function of the *a priori* estimation error covariance. Choosing a constant coefficient to scale  $\mathbf{R}_i$  seems less practical and will probably lead to a more complicated tuning procedure. As long as the measurement noise is white the contributions of the second order effects are not a function of the measurement error covariance, therefore making them a fraction or a multiple of the measurement noise (an unrelated quantity) is likely not the best choice.

**4.3.3. Lear's Method** Lear's choice was to make  $\mathbf{U}_i$  a percentage of the *a priori* estimation error covariance via [39]

$$\mathbf{U}_i = \beta \mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top. \quad (4.36)$$

Let  $\bar{\mathbf{P}}_i^- \in \mathbb{R}^{3 \times 3}$  be the partition of the state estimation error covariance associated with the position error states. The Space Shuttle employs underweighting when  $\sqrt{\text{tr } \bar{\mathbf{P}}_i^-} > \alpha$ . The positive scalars  $\alpha$  and  $\beta$  are design parameters. For the Space Shuttle,  $\alpha$  is selected to be 1000 meters and  $\beta$  is selected to be 0.2 [39]. When  $\sqrt{\text{tr } \bar{\mathbf{P}}_i^-} > 1000$  m, then  $\beta = 0.2$ , otherwise  $\beta = 0$ .

Orion employs a slightly different approach, underweighting is applied when  $\mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top > \alpha$ , where  $\alpha$  is a tunable flight software parameter.

This choice of underweighting scheme is sounds since it assumes that the higher order effect are a fraction (or a multiple) of the first order effects, which are a related quantity. Some unusual nonlinear measurement cases where the measurement Jacobian evaluates to zero, or a small value, while the Hessian does not vanish are not appropriately handled by underweighting.

#### 4.4. Pre-Flight Tuning Aids

In this section, a technique to aid the tuning of the underweighting coefficient during pre-flight analysis is presented. When the nonlinearities lead to second-order terms that cannot be neglected, we find that the measurement residual covariance is more accurately given by (see Eq. (4.30))

$$\mathbf{W}_i = \mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top + \mathbf{R}_i + \mathbf{B}_i. \quad (4.37)$$

Following Lear's approach to underweighting, we have that

$$\mathbf{W}_{U,i} = (1 + \beta) \mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top + \mathbf{R}_i. \quad (4.38)$$

Comparing Eqs. (4.37) and (4.38), the desired effect is to have

$$\text{tr } \mathbf{W}_{U,i} \geq \text{tr } \mathbf{W}_i. \quad (4.39)$$

This leads us to choose the underweighting coefficient  $\beta_i$  such that

$$\beta \geq \text{tr } \mathbf{B}_i / \text{tr } \mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^\top \quad \forall i. \quad (4.40)$$

Running simulations pre-flight, the designer can calculate the time history of  $\text{tr } \mathbf{B}_i / \text{tr } \mathbf{H}_i \mathbf{P}_i^- \mathbf{H}_i^T$  and choose an appropriate value of  $\beta$ . It is unlikely that higher order terms than  $\mathbf{B}_i$  will need to be considered in designing the value of  $\beta$ .

## CHAPTER 5

# Bias Modeling

Contributed by J. Russell Carpenter

A general model for a measurement error is as follows:

$$\mathbf{e} = \mathbf{b} + \mathbf{v} \quad (5.1)$$

where  $\mathbf{b}$  models the systematic errors, and  $\mathbf{v}$  models the measurement noise. We assume that the measurement noise is a discrete sequence of uncorrelated random numbers. Variables such as  $\mathbf{v}$  are known as random variables, and Appendix A describes how to model them. This Chapter describe models for the systematic errors.

The discussion of systematic errors treats such errors as scalar quantities to simplify the exposition; generalization to the vector case is straightforward. Note that if the measurement is non-scalar, but the errors in the component measurements are independent of one another, then we can model each measurement independently, so modeling the biases as vector is not required. If the measurement errors are not independent, then many estimators require that we apply a transformation to the data prior to processing so that the data input to the estimator have independent measurement errors; Appendix A describes some ways to accomplish this transformation.

### 5.1. Zero-Input Bias State Models

The simplest non-zero measurement error consists only of measurement noise. The next simplest class of measurement errors consists of biases which are either themselves constant, or are the integrals of constants. We can view such biases as the output of a system which has zero inputs, and which may have internal states. In the sequel, we will consider cases where there are random inputs to the system.

In cases were the bias is the output of a system with internal states, the estimator may treat the internal states as solve-for or consider parameters. In such cases, the estimator requires a measurement partials matrix. Otherwise, the “measurement partial” is just  $H = \partial b / \partial b = 1$ .

**5.1.1. Random Constant** The simplest type of systematic error is a constant bias on the measurement. There are two types of such biases: deterministic constants, which are truly constant for all time, and random constants, which are constant or very nearly so over a particular time of interest. For example, each time a sensor is power-cycled, a bias associated with it may change in value, but so long as the sensor remains powered on, the bias will not change.

In some cases, we may have reason to believe that a particular systematic error source truly is a deterministic bias, but due to limited observability, we do not have knowledge of

its true value. In such cases, we may view our estimate of the bias as a random constant, and its variance as a measure of the imprecision of our knowledge.

Thus, we may view all constants that could be solve-for or consider parameters in orbit determination as random constants. A model for a random constant is

$$\dot{\mathbf{b}}(t) = 0, \mathbf{b}(t_o) \sim N(0, p_{\mathbf{b}o}). \quad (5.2)$$

Thus the unconditional mean of  $\mathbf{b}(t)$  is zero for all time, and its unconditional covariance is constant for all time as well. If  $\mathbf{b}(t)$  is a filter solve-for variable that is observable, then its covariance conditioned on the measurement sequence will reach zero in the limit as  $t \rightarrow \infty$ . This is an undesirable characteristic for application in a sequential navigation filter.

To simulate a realization of the random constant, we need only generate a random number according to  $N(0, p_{\mathbf{b}o})$ , as the previous subsection described.

**5.1.2. Random Ramp** The random ramp model assumes that the rate of change of the bias is itself a random constant; thus the random ramp model is

$$\ddot{\mathbf{b}}(t) = 0, \dot{\mathbf{b}}(t_o) \sim N(0, p_{\dot{\mathbf{b}}o}). \quad (5.3)$$

Thus, the initial condition  $\dot{\mathbf{b}}(t_o)$  is a random constant. For a pure random ramp, the initial condition on  $\mathbf{b}(t_o)$  and its covariance are taken to be zero, but an obvious and common generalization is to allow  $\mathbf{b}(t_o)$  to also be a random constant.

It is convenient to write this model as a first-order vector system as follows:

$$\begin{bmatrix} \dot{\mathbf{b}}(t) \\ \ddot{\mathbf{b}}(t) \end{bmatrix} = \begin{bmatrix} \dot{\mathbf{b}}(t) \\ \mathbf{d}(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{b}(t) \\ \mathbf{d}(t) \end{bmatrix} \quad (5.4)$$

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) \quad (5.5)$$

The resulting output equation for the total measurement error is

$$\mathbf{e} = \begin{bmatrix} 1 & 0 \end{bmatrix} \mathbf{x} + \mathbf{v} \quad (5.6)$$

$$= \mathbf{H}\mathbf{x} + \mathbf{v} \quad (5.7)$$

Note that the ensemble of realizations of  $\mathbf{x}(t)$  has zero-mean for all time. The unconditional covariance evolves in time according to

$$\mathbf{P}_{\mathbf{x}}(t) = \mathbf{\Phi}(t - t_o)\mathbf{P}_{\mathbf{x}o}\mathbf{\Phi}^T(t - t_o) \quad (5.8)$$

where

$$\mathbf{\Phi}(t) = \begin{bmatrix} 1 & t \\ 0 & 1 \end{bmatrix} \text{ and } \mathbf{P}_{\mathbf{x}o} = \begin{bmatrix} p_{\mathbf{b}o} & 0 \\ 0 & p_{\dot{\mathbf{b}}o} \end{bmatrix} \quad (5.9)$$

which we can also write in recursive form as

$$\mathbf{P}_{\mathbf{x}}(t + \Delta t) = \mathbf{\Phi}(\Delta t)\mathbf{P}_{\mathbf{x}}(t)\mathbf{\Phi}^T(\Delta t) \quad (5.10)$$

Thus, we can generate realizations of the random ramp with either  $\mathbf{x}(t) \sim N(\mathbf{0}, \mathbf{P}_{\mathbf{x}}(t))$  or recursively from

$$\mathbf{x}(t + \Delta t) = \mathbf{\Phi}(\Delta t)\mathbf{x}(t) \quad (5.11)$$

Note that the norm of the unconditional covariance becomes infinite as  $t^2$  becomes infinite. This could lead to an overflow of the representation of the covariance in a computer program if the propagation time between measurements is large, if the bias is unobservable, or if the bias is a consider parameter, and could also lead to the representation of the covariance losing either its symmetry and/or its positive definiteness due to roundoff and/or truncation. If the bias and drift are filter solve-for variables, then the norm of their

covariance conditioned on the measurement sequence will reach zero in the limit as  $t \rightarrow \infty$ . These are all undesirable characteristics for application in a sequential navigation filter.

**5.1.3. Higher-Order Derivatives of Random Constants** In principle, a random constant may be associated with any derivative of the bias in a straightforward extension of the models above. In practice, it is rare to need more than two or three derivatives. Conventional terminology does not appear in the literature for derivatives of higher order than the random ramp. The slope of the bias is most commonly described as the “bias drift,” so that a “drift random ramp” would be one way to describe a bias whose second derivative is a random constant. The measurement partials matrix needs to be accordingly padded with trailing zeros for the derivatives of the bias in such cases.

## 5.2. Single-Input Bias State Models

The simplest non-constant systematic errors are systems with a single input that is a random process. We can think of a random process as the result of some kind of limit in which the intervals between an uncorrelated sequence of random variables get infinitesimally small. In this limit, each random increment instantaneously perturbs the sequence, so that the resulting process is continuous but non-differentiable. We call this kind of a random input “process noise.”

Although such random processes are non-differentiable, there are various techniques for generalizing the concept of integration so that something like integrals of the process noise exist, and hence so do the differentials that appear under the integral signs. It turns out that so long as any coefficients of the process noise are non-random, these differentials behave for all practical purposes as if they were differentiable.

**5.2.1. Random Walk** The random walk is the simplest random process of the type described above. In terms of the “formal derivatives” mentioned above, the random walk model for a measurement bias is

$$\dot{\mathbf{b}}(t) = \mathbf{w}(t), \quad w(t) \sim N(0, q\delta(t-s)) \quad (5.12)$$

The input noise process on the right hand side is known as “white noise,” and the Dirac delta function that appears in the expression for its variance indicates that the white noise process consists of something like an infinitely-tightly spaced set of impulses. The term  $q$  that appears along with the delta function is the intensity of each impulse<sup>1</sup>. The initial condition  $\mathbf{b}(t_o)$  is an unbiased random constant. Since  $\mathbf{b}(t_o)$  and  $w(t)$  are zero-mean, then  $\mathbf{b}(t)$  is also zero-mean for all time. The unconditional variance of  $\mathbf{b}$  evolves in time according to

$$p_{\mathbf{b}}(t) = p_{\mathbf{b}o} + q(t - t_o) \quad (5.13)$$

which we can also write in recursive form as

$$p_{\mathbf{b}}(t + \Delta t) = p_{\mathbf{b}}(t) + q\Delta t \quad (5.14)$$

Thus, to generate a realization of the random walk at time  $t$ , we need only generate a random number according to  $N(0, p_{\mathbf{b}}(t))$ . Equivalently, we could also generate realizations of  $\varpi(t) \sim N(0, q\Delta t)$ , and recursively add these discrete noise increments to the bias as follows:

$$b(t + \Delta t) = b(t) + \varpi(t) \quad (5.15)$$

---

<sup>1</sup>Another way to imagine the input sequence, in terms of a frequency domain interpretation, is that it is a noise process whose power spectral density,  $q$ , is non-zero at all frequencies, which implies infinite bandwidth.

Note that the unconditional variance becomes infinite as  $t$  becomes infinite. This could lead to an overflow of the representation of  $p_b$  if  $q$  is large in the following circumstances: in a long gap between measurements, if the bias is unobservable, or if the bias is a consider parameter. These are all somewhat undesirable characteristics for application in a sequential navigation filter. However, because the process is persistently stimulated by the input, its variance conditioned on a measurement history will remain positive for all time. Hence the random walk finds frequent application in sequential navigation filters, particularly when continuous measurements from which the bias is observable are generally available, such as often occurs for GPS data.

**5.2.2. Random Run** The random run model assumes that the rate of change of the bias is itself a random walk; thus the random run model is

$$\ddot{\mathbf{b}}(t) = \mathbf{w}(t), \quad w(t) \sim N(0, q\delta(t-s)) \quad (5.16)$$

The initial condition  $\dot{\mathbf{b}}(t_o)$  is a random constant. For a pure random run, the initial condition on  $\mathbf{b}(t_o)$  and its covariance are taken to be zero, but an obvious and common generalization is to allow  $\mathbf{b}(t_o)$  to also be a random constant.

It is convenient to write this model as a first-order vector system as follows:

$$\begin{bmatrix} \dot{\mathbf{b}}(t) \\ \ddot{\mathbf{b}}(t) \end{bmatrix} = \begin{bmatrix} \dot{\mathbf{b}}(t) \\ \dot{\mathbf{d}}(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{b}(t) \\ \mathbf{d}(t) \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \mathbf{w}(t) \quad (5.17)$$

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{w}(t) \quad (5.18)$$

The measurement partial is the same as for the random ramp. The initial condition  $\mathbf{x}(t_o)$  is an unbiased random constant. Since  $\mathbf{x}(t_o)$  and  $\mathbf{w}(t)$  are zero-mean, then  $\mathbf{x}(t)$  is also zero-mean for all time. The covariance evolves in time according to

$$\mathbf{P}_{\mathbf{x}}(t) = \mathbf{\Phi}(t-t_o)\mathbf{P}_{\mathbf{x}o}\mathbf{\Phi}^T(t-t_o) + \mathbf{S}(t-t_o) \quad (5.19)$$

where

$$\mathbf{\Phi}(t-t_o) = \begin{bmatrix} 1 & t \\ 0 & 1 \end{bmatrix} \text{ and } \mathbf{P}_{\mathbf{x}o} = \begin{bmatrix} p_{b_o} & 0 \\ 0 & p_{\dot{b}_o} \end{bmatrix} \quad (5.20)$$

and

$$\mathbf{S}(t) = q \begin{bmatrix} t^3/3 & t^2/2 \\ t^2/2 & t \end{bmatrix} \quad (5.21)$$

which we can also write in recursive form as

$$\mathbf{P}_{\mathbf{x}}(t+\Delta t) = \mathbf{\Phi}(\Delta t)\mathbf{P}_{\mathbf{x}}(t)\mathbf{\Phi}^T(\Delta t) + \mathbf{S}(\Delta t) \quad (5.22)$$

Thus, we can generate realizations of the random run with either  $\mathbf{x}(t) \sim N(\mathbf{0}, \mathbf{P}_{\mathbf{x}}(t))$  or recursively from

$$\mathbf{x}(t+\Delta t) = \mathbf{\Phi}(\Delta t)\mathbf{x}(t) + \boldsymbol{\varpi}(t) \quad (5.23)$$

where  $\boldsymbol{\varpi}(t) \sim N(\mathbf{0}, \mathbf{S}(\Delta t))$  is a noise sample *vector* arising from formal integration of the scalar noise input process over the sample time. A Cholesky decomposition of  $\mathbf{S}(t)$  useful for sampling is

$$\sqrt{\mathbf{S}(t)} = \begin{bmatrix} \sqrt{3t^3/3} & 0 \\ \sqrt{3t}/2 & \sqrt{t}/2 \end{bmatrix} \quad (5.24)$$

Note that the norm of the unconditional covariance becomes infinite as  $t^3$  becomes infinite, and the process is persistently stimulated by the input, so its covariance conditioned on a measurement history will remain positive definite for all time. Hence, the random run

shares similar considerations with the random walk for application in sequential navigation filters.

**5.2.3. Higher-Order Derivatives of Random Walks** In principle, a random walk may be associated with any derivative of the bias in a straightforward extension of the models above. In practice, it is rare to need more than two or three derivatives. Conventional terminology does not appear in the literature for derivatives of higher order than the random run. A “drift random run” would be one way to describe a bias whose second derivative is a random walk. Below, we will refer to such a model as a “random zoom.”

**5.2.4. First-Order Gauss-Markov** The first-order Gauss-Markov (FOGM) process is one of the simplest random processes that introduces time correlation between samples. In terms of a frequency domain interpretation, we can view it as white noise passed through a low-pass filter. Since such noise, often called “colored noise,” has finite bandwidth, it is physically realizable, unlike white noise. In the notation of formal derivatives, the FOGM model is

$$\dot{\mathbf{b}}(t) = -\frac{1}{\tau}\mathbf{b}(t) + \mathbf{w}(t), \quad (5.25)$$

where, as with the random walk,  $b(t_o) \sim N(0, p_{bo})$ , and  $w(t) \sim N(0, q\delta(t-s))$ . The time constant,  $\tau$  gives the correlation time, or the time over which the intensity of the time correlation will fade to  $1/e$  of its prior value<sup>2</sup>.

Since  $\mathbf{b}(t_o)$  and  $\mathbf{w}(t)$  are zero-mean, then  $\mathbf{b}(t)$  is also zero-mean for all time. The covariance evolves in time according to

$$p_{\mathbf{b}}(t) = e^{-\frac{2}{\tau}(t-t_o)} p_{bo} + s(t-t_o) \quad (5.26)$$

where

$$s(t-t_o) = \frac{q\tau}{2} \left(1 - e^{-\frac{2}{\tau}(t-t_o)}\right) \quad (5.27)$$

which we can also write in recursive form as

$$p_{\mathbf{b}}(t + \Delta t) = e^{-\frac{2\Delta t}{\tau}} p_{\mathbf{b}}(t) + s(\Delta t) \quad (5.28)$$

Thus, to generate discrete samples of a particular realization of the FOGM, we can either generate samples from  $b(t) \sim N(0, p_{\mathbf{b}}(t))$ , or generate a realization of the initial bias value, and then at each sample time generate realizations of  $\varpi(t) \sim N(0, s(\Delta t))$ , and recursively add these discrete noise sample increments to the bias sample history as follows:

$$b(t + \Delta t) = e^{-\frac{\Delta t}{\tau}} b(t) + \varpi(t) \quad (5.29)$$

Note that  $p_{\mathbf{b}}$  approaches a finite steady-state value equal to  $q\tau/2$  as  $t$  becomes infinite. One can choose the parameters of the FOGM so that this steady-state value avoids any overflow of the representation of  $p_{\mathbf{b}}$  in a computer program, and such that the FOGM’s covariance evolution prior to reaching steady-state closely mimics that of a random walk. For these reasons, the FOGM is recommended as a *best practice* for bias modeling in sequential navigation filters.

---

<sup>2</sup>One sometimes sees  $\tau$  described as the “half-life,” but since  $1/e < 1/2$ , this is not an accurate label.

**5.2.5. Integrated First-Order Gauss-Markov Model** As with the random walk and random constant models, any number of derivatives of the bias may be associated with a FOGM process. However, integration of the FOGM destroys its stability. For example, the singly integrated first-order Gauss-Markov model is given by

$$\begin{bmatrix} \dot{\mathbf{b}}(t) \\ \dot{\mathbf{d}}(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & -1/\tau \end{bmatrix} \begin{bmatrix} \mathbf{b}(t) \\ \mathbf{d}(t) \end{bmatrix} + \begin{bmatrix} 0 \\ \mathbf{w}(t) \end{bmatrix}, \quad (5.30)$$

which leads to the following state transition matrix,

$$\Phi(t) = \begin{bmatrix} 1 & \tau(1 - e^{-t/\tau}) \\ 0 & e^{-t/\tau} \end{bmatrix}, \quad (5.31)$$

and process noise covariance<sup>3</sup>,

$$\mathbf{S}(t) = \frac{q\tau}{2} \begin{bmatrix} \tau^2 \{ (1 - e^{-2t/\tau}) - 4(1 - e^{-t/\tau}) + 2t/\tau \} & \tau(1 - e^{-t/\tau})^2 \\ \tau(1 - e^{-t/\tau})^2 & (1 - e^{-2t/\tau}) \end{bmatrix}. \quad (5.32)$$

Clearly, this is an unstable model, as the bias variance increases linearly with elapsed time. As an alternative, the following second-order model is available.

**5.2.6. Second-Order Gauss-Markov** The model for a second-order Gauss-Markov (SOGM) random process is

$$\ddot{\mathbf{b}}(t) = -2\zeta\omega_n\dot{\mathbf{b}}(t) - \omega_n^2\mathbf{b}(t) + \mathbf{w}(t), \quad \mathbf{w}(t) \sim N(0, q\delta(t-s)) \quad (5.33)$$

The initial conditions  $\mathbf{b}(t_o)$  and  $\dot{\mathbf{b}}(t_o)$  are random constants. It is convenient to write this model as a first-order vector system as follows:

$$\begin{bmatrix} \dot{\mathbf{b}}(t) \\ \dot{\mathbf{b}}(t) \end{bmatrix} = \begin{bmatrix} \dot{\mathbf{b}}(t) \\ \mathbf{b}(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -\omega_n^2 & -2\zeta\omega_n \end{bmatrix} \begin{bmatrix} \mathbf{b}(t) \\ \mathbf{d}(t) \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \mathbf{w}(t) \quad (5.34)$$

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{w}(t) \quad (5.35)$$

The measurement partial is the same as for the random ramp. The initial condition  $\mathbf{x}(t_o)$  is an unbiased random constant vector. Since  $\mathbf{x}(t_o)$  and  $\mathbf{w}(t)$  are zero-mean, then  $\mathbf{x}(t)$  is also zero-mean for all time.

The covariance evolves in time according to according to

$$\mathbf{P}_{\mathbf{x}}(t) = \Phi(t - t_o)\mathbf{P}_{\mathbf{x}o}\Phi^T(t - t_o) + \mathbf{S}(t - t_o) \quad (5.36)$$

which we can also write in recursive form as

$$\mathbf{P}_{\mathbf{x}}(t + \Delta t) = \Phi(\Delta t)\mathbf{P}_{\mathbf{x}}(t)\Phi^T(\Delta t) + \mathbf{S}(\Delta t) \quad (5.37)$$

Thus, we can generate realizations of the SOGM with either  $\mathbf{x}(t) \sim N(\mathbf{0}, \mathbf{P}_{\mathbf{x}}(t))$  or recursively from

$$\mathbf{x}(t + \Delta t) = \Phi(\Delta t)\mathbf{x}(t) + \boldsymbol{\varpi}(t) \quad (5.38)$$

where  $\boldsymbol{\varpi}(t) \sim N(0, \mathbf{S}(\Delta t))$ .

For the underdamped case ( $\zeta < 1$ ), the state transition matrix and discrete process noise covariance are given by Reference **73**:

$$\Phi(t) = \frac{e^{-\zeta\omega_n t}}{\omega_d} \begin{bmatrix} (\omega_d \cos \omega_d t + \zeta\omega_n \sin \omega_d t) & \sin \omega_d t \\ -\omega_n^2 \sin \omega_d t & (\omega_d \cos \omega_d t - \zeta\omega_n \sin \omega_d t) \end{bmatrix} \quad (5.39)$$

<sup>3</sup>Note that (5.32) corrects an error in Reference **10**.

and

$$Q_{\Delta}^{(1,1)}(t) = \frac{q}{4\zeta\omega_n^3} \left[ 1 - \frac{e^{-2\zeta\omega_n t}}{\omega_d^2} (\omega_d^2 + 2\zeta\omega_n\omega_d \cos \omega_d t \sin \omega_d t + 2\zeta^2\omega_n^2 \sin^2 \omega_d t) \right] \quad (5.40)$$

$$Q_{\Delta}^{(2,2)}(t) = \frac{q}{4\zeta\omega_n} \left[ 1 - \frac{e^{-2\zeta\omega_n t}}{\omega_d^2} (\omega_d^2 - 2\zeta\omega_n\omega_d \cos \omega_d t \sin \omega_d t + 2\zeta^2\omega_n^2 \sin^2 \omega_d t) \right] \quad (5.41)$$

$$Q_{\Delta}^{(1,2)}(t) = \frac{q}{2\omega_d^2} e^{-2\zeta\omega_n t} \sin^2 \omega_d t, \quad (5.42)$$

$$Q_{\Delta}^{(2,1)}(t) = Q_{\Delta}^{(1,2)}(t) \quad (5.43)$$

where  $\omega_d = \omega_n \sqrt{1 - \zeta^2}$ . In the over-damped case ( $\zeta > 1$ ), replace sin and cos with sinh and cosh, respectively. In the critically-damped case,

$$\Phi(t) = \begin{bmatrix} e^{-\omega_n t}(1 + \omega_n t) & t e^{-\omega_n t} \\ -\omega_n^2 t e^{-\omega_n t} & e^{-\omega_n t}(1 - \omega_n t) \end{bmatrix} \quad (5.44)$$

and

$$Q_{\Delta}^{(1,1)}(t) = \frac{q}{4\omega_n^3} [1 - e^{-2\omega_n t}(1 + 2\omega_n t + 2\omega_n^2 t^2)] \quad (5.45)$$

$$Q_{\Delta}^{(2,2)}(t) = \frac{q}{4\omega_n} [1 - e^{-2\omega_n t}(1 - 2\omega_n t + 2\omega_n^2 t^2)] \quad (5.46)$$

$$Q_{\Delta}^{(2,1)}(t) = Q_{\Delta}^{(1,2)}(t) = \frac{qt^2}{2} e^{-2\omega_n t} \quad (5.47)$$

Note that for any damping ratio,  $\|\mathbf{P}_{\mathbf{x}}\|$  remains finite, since as  $t \rightarrow \infty$ ,

$$\mathbf{P}_{\mathbf{x}}(t \rightarrow \infty) = \frac{q}{4\zeta\omega_n} \begin{bmatrix} 1/\omega_n^2 & 0 \\ 0 & 1 \end{bmatrix}. \quad (5.48)$$

Thus, the ratio of the steady-state standard deviations of the bias and drift will be

$$\frac{\sigma_d}{\sigma_b} = \omega_n, \quad (5.49)$$

and these are related to the power spectral density by

$$q = 4\zeta \frac{\sigma_d^3}{\sigma_b}. \quad (5.50)$$

Hence, we can choose the parameters of the SOGM so that we avoid any overflow, loss of symmetry and/or positive definiteness of  $\mathbf{P}_{\mathbf{x}}$  due to roundoff and/or truncation. For these reasons, the SOGM is recommended as a *best practice* for bias drift modeling in sequential navigation filters.

**5.2.7. Vasicek Model** A criticism of the FOGM process is that as  $t \rightarrow \infty$ ,  $E[b(t)] \rightarrow 0$ . In the filtering context, this implies that a data outage that is long relative to the time constant,  $\tau$ , can result in the filter's bias estimate decaying toward zero, which may be undesirable. To address this concern, Seago et al. [66] proposed that biases the filter should retain across such outages might be modeled instead with a model proposed by Vasicek [72] for modeling interest rates:

$$\dot{\mathbf{b}}(t) = -\frac{1}{\tau}(\mathbf{b}(t) - \mathbf{b}_{\infty}) + \mathbf{w}(t), \quad (5.51)$$

where, as previously,  $b(t_o) \sim N(0, \mathbf{P}_{b_o})$ , and  $w(t) \sim N(0, q\delta(t-s))$ . A formal solution to (5.51) gives

$$\mathbf{b}(t) = \mathbf{b}(t_o) e^{-\frac{t-t_o}{\tau}} + b_\infty(1 - e^{-\frac{t-t_o}{\tau}}) + \int_{t_o}^t e^{-\frac{s}{\tau}} \mathbf{w}(s) ds \quad (5.52)$$

Since  $\mathbf{b}(t_o)$  and  $\mathbf{w}(t)$  are zero-mean, then

$$\mathbb{E}[\mathbf{b}(t)] = b_\infty(1 - e^{-\frac{t-t_o}{\tau}}) \quad (5.53)$$

and  $\mathbb{E}[\mathbf{b}(t)] = b_\infty$  as  $t \rightarrow \infty$ . Since  $\mathbb{E}[\mathbf{b}(t)]^2$  is subtracted from  $\mathbb{E}[\mathbf{b}(t)^2]$  to get the covariance, the covariance evolves in time identically to the FOGM,

$$p_b(t) = e^{-\frac{2}{\tau}(t-t_o)} p_{b_o} + s(t-t_o) \quad (5.54)$$

where as before

$$s(t-t_o) = \frac{q\tau}{2} \left(1 - e^{-\frac{2}{\tau}(t-t_o)}\right) \quad (5.55)$$

Thus, to generate a realization of the Vasicek Model at particular time  $t$ , we could generate a realization of the initial bias value, and then at each sample time generate realizations of  $\varpi(t) \sim N(0, s(\Delta t))$ , and recursively add these discrete noise sample increments to the bias sample history as follows:

$$b(t + \Delta t) = b(t) e^{-\frac{\Delta t}{\tau}} + b_\infty(1 - e^{-\frac{\Delta t}{\tau}}) + \varpi(t) \quad (5.56)$$

or we could generate a random realization of  $N(0, p_b(t))$  and add this to  $b_\infty(1 - e^{-\frac{t-t_o}{\tau}})$ .

To configure or “tune” the Vasicek model, one chooses the time constant  $\tau$  and the noise PSD  $q$  in a manner analogous to the FOGM process; it is less clear how one might choose  $b_\infty$ . Seago et al. [66] proposed that  $b_\infty$  be estimated as a random constant filter state. Casting the Vasicek model into such a two-state form results in the following model:

$$\begin{bmatrix} \dot{\mathbf{b}}(t) \\ \dot{\mathbf{b}}_\infty \end{bmatrix} = \begin{bmatrix} -1/\tau & 1/\tau \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{b}(t) \\ \mathbf{b}_\infty \end{bmatrix} + \begin{bmatrix} \mathbf{w}(t) \\ 0 \end{bmatrix}, \quad (5.57)$$

which leads to the following state transition matrix and process noise covariance:

$$\Phi(t) = \begin{bmatrix} e^{-t/\tau} & 1 - e^{-t/\tau} \\ 0 & 1 \end{bmatrix}, \quad \mathbf{S}(t) = \begin{bmatrix} \frac{q\tau}{2} \left(1 - e^{-\frac{2t}{\tau}}\right) & 0 \\ 0 & 0 \end{bmatrix}. \quad (5.58)$$

While the Vasicek Model shares with the FOGM the desirable feature that  $p_b \rightarrow q\tau/2$  as  $t \rightarrow \infty$ , in the two-state form just described, it also has the undesirable feature that the variance of  $\mathbf{b}_\infty$  goes to zero as  $t \rightarrow \infty$ . Modeling  $\mathbf{b}_\infty$  with process noise, e.g. as a random walk with PSD of  $q_\infty$ , introduces an unstable integral of the process noise as occurs for the integrated FOGM:

$$\mathbf{S}(t) = \begin{bmatrix} q_\infty \left( t - \frac{3\tau}{2} + 2\tau e^{-\frac{t}{\tau}} - \frac{\tau}{2} e^{-\frac{2t}{\tau}} \right) + \frac{q\tau}{2} \left( 1 - e^{-\frac{2t}{\tau}} \right) & q_\infty t - q_\infty \tau \left( 1 - e^{-\frac{t}{\tau}} \right) \\ q_\infty t - q_\infty \tau \left( 1 - e^{-\frac{t}{\tau}} \right) & q_\infty t \end{bmatrix}, \quad (5.59)$$

although choosing  $q_\infty$  appropriately small may mitigate this concern. In any case, retaining a steady-state bias across long data gaps may not always be warranted, depending on the context. And if long measurement gaps are not present, the need to retain such a bias, with the accompanying complexity of maintaining an additional state, may not be necessary. We will consider further such multi-input bias models in the sequel.

### 5.3. Multi-Input Bias State Models

We may combine any of the above models to create multi-input bias models; for example the bias could be a second-order Gauss-Markov, and the bias rate could be a first-order Gauss-Markov. In practice, the most useful combinations have been found to be the following.

**5.3.1. Bias and Drift Random Walks (Random Walk + Random Run)** A common model for biases in clocks, gyros, and accelerometers is that the bias is driven by both its own white noise input, and also by the integral of the white noise of its drift. Such models derive from observations that the error magnitudes of these devices depend on the time scale over which the device is observed. They are often characterized by Allan deviation specifications, which may be heuristically associated with the white noise power spectral densities. The model is as follows:

$$\begin{bmatrix} \dot{\mathbf{b}}(t) \\ \dot{\mathbf{d}}(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{b}(t) \\ \mathbf{d}(t) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{w}_b(t) \\ \mathbf{w}_d(t) \end{bmatrix} \quad (5.60)$$

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{w}(t) \quad (5.61)$$

The measurement partial is the same as for the random ramp. The initial condition  $\mathbf{x}(t_o)$  is an unbiased random constant. Since  $\mathbf{x}(t_o)$  and  $\mathbf{w}(t)$  are zero-mean, then  $\mathbf{x}(t)$  is also zero-mean for all time. The covariance evolves in time according to

$$\mathbf{P}_x(t) = \mathbf{\Phi}(t - t_o)\mathbf{P}_{x_o}\mathbf{\Phi}^\top(t - t_o) + \mathbf{S}(t - t_o) \quad (5.62)$$

where

$$\mathbf{\Phi}(t) = \begin{bmatrix} 1 & t \\ 0 & 1 \end{bmatrix} \text{ and } \mathbf{P}_{x_o} = \begin{bmatrix} p_{b_o} & 0 \\ 0 & p_{d_o} \end{bmatrix} \quad (5.63)$$

and

$$\mathbf{S}(t) = \begin{bmatrix} q_b t + q_d t^3/3 & q_d t^2/2 \\ q_d t^2/2 & q_d t \end{bmatrix} \quad (5.64)$$

which we can also write in recursive form as

$$\mathbf{P}_x(t + \Delta t) = \mathbf{\Phi}(\Delta t)\mathbf{P}_x(t)\mathbf{\Phi}^\top(\Delta t) + \mathbf{S}(\Delta t) \quad (5.65)$$

Thus, we can generate realizations of the random run with either  $\mathbf{x}(t) \sim N(\mathbf{0}, \mathbf{P}_x(t))$  or recursively from

$$\mathbf{x}(t + \Delta t) = \mathbf{\Phi}(\Delta t)\mathbf{x}(t) + \boldsymbol{\varpi}(t) \quad (5.66)$$

where  $\boldsymbol{\varpi}(t) \sim N(0, \mathbf{S}(\Delta t))$ . Note that a Cholesky decomposition of  $\mathbf{S}(t)$  is

$$\sqrt{\mathbf{S}(t)} = \begin{bmatrix} \sqrt{q_b t + q_d t^3/12} & \sqrt{q_d t^3/2} \\ 0 & \sqrt{q_d t} \end{bmatrix} \quad (5.67)$$

As with its constituent models, the norm of the unconditional covariance becomes infinite as  $t^3$  becomes infinite, while the process is persistently stimulated by the input, so its covariance conditioned on a measurement history will remain positive definite for all time. Hence, the this model shares similar considerations with its constituents for application in sequential navigation filters.

**5.3.2. Bias, Drift, and Drift Rate Random Walks (Random Walk + Random Run + Random Zoom)** Another model for biases in very-high precision clocks, gyros, and accelerometers is that the bias is driven by two integrals of white noise in addition to its own white noise input. Such models are often characterized by Hadamard deviation specifications, which may be heuristically associated with the white noise power spectral densities. The model is as follows:

$$\begin{bmatrix} \dot{\mathbf{b}}(t) \\ \dot{\mathbf{d}}(t) \\ \ddot{\mathbf{d}}(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{b}(t) \\ \mathbf{d}(t) \\ \dot{\mathbf{d}}(t) \end{bmatrix} + \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{w}_b(t) \\ \mathbf{w}_d(t) \\ \mathbf{w}_{\dot{d}}(t) \end{bmatrix} \quad (5.68)$$

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{w}(t) \quad (5.69)$$

The resulting output equation is

$$\mathbf{e} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \mathbf{x} + \mathbf{v} \quad (5.70)$$

$$= \mathbf{H}\mathbf{x} + \mathbf{v} \quad (5.71)$$

The initial condition  $\mathbf{x}(t_o)$  is an unbiased random constant. Since  $\mathbf{x}(t_o)$  and  $\mathbf{w}(t)$  are zero-mean, then  $\mathbf{x}(t)$  is also zero-mean for all time. The covariance evolves in time according to

$$\mathbf{P}_x(t) = \mathbf{\Phi}(t - t_o)\mathbf{P}_{x_o}\mathbf{\Phi}^\top(t - t_o) + \mathbf{S}(t - t_o) \quad (5.72)$$

where

$$\mathbf{\Phi}(t) = \begin{bmatrix} 1 & t & t^2/2 \\ 0 & 1 & t \\ 0 & 0 & 1 \end{bmatrix} \text{ and } \mathbf{P}_{x_o} = \begin{bmatrix} p_{b_o} & 0 & 0 \\ 0 & p_{d_o} & 0 \\ 0 & 0 & p_{\dot{d}_o} \end{bmatrix} \quad (5.73)$$

and

$$\mathbf{S}(t) = \begin{bmatrix} q_b t + q_d t^3/3 + q_{\dot{d}} t^5/5 & q_d t^2/2 + q_{\dot{d}} t^4/8 & q_{\dot{d}} t^3/6 \\ q_d t^2/2 + q_{\dot{d}} t^4/8 & q_d t + q_{\dot{d}} t^3/3 & q_{\dot{d}} t^2/2 \\ q_{\dot{d}} t^3/6 & q_{\dot{d}} t^2/2 & q_{\dot{d}} t \end{bmatrix} \quad (5.74)$$

which we can also write in recursive form as

$$\mathbf{P}_x(t + \Delta t) = \mathbf{\Phi}(\Delta t)\mathbf{P}_x(t)\mathbf{\Phi}^\top(\Delta t) + \mathbf{S}(\Delta t) \quad (5.75)$$

Thus, we can generate realizations of the random run with either  $\mathbf{x}(t) \sim N(\mathbf{0}, \mathbf{P}_x(t))$  or recursively from

$$\mathbf{x}(t + \Delta t) = \mathbf{\Phi}(\Delta t)\mathbf{x}(t) + \boldsymbol{\varpi}(t) \quad (5.76)$$

where  $\boldsymbol{\varpi}(t) \sim N(0, \mathbf{S}(\Delta t))$ . Note that a Cholesky decomposition of  $\mathbf{S}(t)$  is

$$\sqrt{\mathbf{S}(t)} = \begin{bmatrix} \sqrt{q_b t + q_d t^3/12 + q_{\dot{d}} t^5/720} & t/2\sqrt{q_d t + q_{\dot{d}} t^3/12} & t^2/6\sqrt{q_{\dot{d}} t} \\ 0 & \sqrt{q_d t + q_{\dot{d}} t^3/12} & t/2\sqrt{q_{\dot{d}} t} \\ 0 & 0 & \sqrt{q_{\dot{d}} t} \end{bmatrix} \quad (5.77)$$

Similar to its constituent models, the norm of the unconditional covariance becomes infinite as  $t^5$  becomes infinite, while the process is persistently stimulated by the input, so its covariance conditioned on a measurement history will remain positive definite for all time. Hence, the this model shares similar considerations with its constituents for application in sequential navigation filters.

**5.3.3. Bias and Drift Coupled First- and Second-Order Gauss-Markov** The following model provides a stable alternative, developed in Reference **10**, to the “Random Walk + Random Run” model. Note that the following description corrects a sign error in the process noise cross-covariance results of the cited work. The transient response of the stable alternative can be tuned to approximate the Random Walk + Random Run model, and its stable steady-state response can be used to avoid computational issues with long propagation times, observability, consider states, etc. Although this model has received limited application as of the time of this writing, due to its stability, it shows promising potential to evolve into a *best practice* for sequential navigation filtering applications.

The coupled first- and second-order Gauss-Markov model is as follows.

$$\begin{bmatrix} \dot{\mathbf{b}}(t) \\ \dot{\mathbf{d}}(t) \end{bmatrix} = \begin{bmatrix} -1/\tau & 1 \\ -\omega_n^2 & -2\zeta\omega_n \end{bmatrix} \begin{bmatrix} \mathbf{b}(t) \\ \mathbf{d}(t) \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mathbf{w}_b(t) \\ \mathbf{w}_d(t) \end{bmatrix} \quad (5.78)$$

$$\dot{\mathbf{x}}(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{w}(t) \quad (5.79)$$

The measurement partial is the same as for the random ramp. The initial condition  $\mathbf{x}(t_o)$  is an unbiased random constant. Since  $\mathbf{x}(t_o)$  and  $\mathbf{w}(t)$  are zero-mean, then  $\mathbf{x}(t)$  is also zero-mean for all time. The covariance evolves in time according to

$$\mathbf{P}_x(t) = \mathbf{\Phi}(t - t_o)\mathbf{P}_{x_o}\mathbf{\Phi}^\top(t - t_o) + \mathbf{S}(t - t_o) \quad (5.80)$$

where

$$\mathbf{\Phi}(t) = \frac{e^{\eta t}}{\nu} \begin{bmatrix} \nu \cos \nu t + (\eta + 2\zeta\omega_n) \sin \nu t & \sin \nu t \\ -\omega_n^2 \sin \nu t & \nu \cos \nu t + (\eta + \beta) \sin \nu t \end{bmatrix} \quad (5.81)$$

with

$$\beta = 1/\tau, \quad (5.82)$$

$$\eta = -\frac{1}{2}(\beta + 2\zeta\omega_n), \quad (5.83)$$

$$\nu = \sqrt{\omega_d^2 + \beta\zeta\omega_n - \frac{1}{4}\beta^2}, \quad (5.84)$$

$$\omega_d = \omega_n\sqrt{1 - \zeta^2}, \quad (5.85)$$

and we assume that  $\nu^2 > 0$ . Let

$$\kappa = -\frac{\beta}{2} + \zeta\omega_n;$$

then, the process noise covariance is given by the following:

$$\begin{aligned} \mathbf{S}^{(1,1)}(t) = q_b & \left[ \frac{e^{2\eta t} - 1}{4\eta} \left( 1 + \frac{\kappa^2}{\nu^2} \right) + \frac{e^{2\eta t} \sin 2\nu t}{4(\eta^2 + \nu^2)} \left( \frac{\nu^2 - \kappa^2 + \eta\kappa}{\nu} \right) \right. \\ & \left. + \frac{e^{2\eta t} \cos 2\nu t - 1}{4(\eta^2 + \nu^2)} \left( \frac{\eta\nu^2 - \eta\kappa^2 + 2\nu^2\kappa}{\nu^2} \right) \right] \end{aligned} \quad (5.86)$$

$$\begin{aligned} \mathbf{S}^{(2,2)}(t) = q_d & \left[ \frac{e^{2\eta t} - 1}{4\eta} \left( 1 + \frac{\kappa^2}{\nu^2} \right) + \frac{e^{2\eta t} \sin 2\nu t}{4(\eta^2 + \nu^2)} \left( \frac{\nu^2 - \kappa^2 + \eta\kappa}{\nu} \right) \right. \\ & \left. + \frac{e^{2\eta t} \cos 2\nu t - 1}{4(\eta^2 + \nu^2)} \left( \frac{\eta\nu^2 - \eta\kappa^2 + 2\nu^2\kappa}{\nu^2} \right) \right] \\ & + \frac{q_b \omega_n^4}{\nu^2} \left( \frac{e^{2\eta t} - 1}{4\eta} - \frac{e^{2\eta t} (\nu \sin 2\nu t + \eta \cos 2\nu t) - \eta}{4(\eta^2 + \nu^2)} \right) \end{aligned} \quad (5.87)$$

$$\begin{aligned} \mathbf{S}^{(1,2)}(t) = \frac{q_b \omega_n^2}{\nu^2} & \left[ \frac{\kappa}{4\eta} (1 - e^{2\eta t}) + \frac{e^{2\eta t} [(\nu\kappa - \eta\nu) \sin 2\nu t + (\eta\kappa + \nu^2) \cos 2\nu t] - (\eta\kappa + \nu^2)}{4(\eta^2 + \nu^2)} \right] \\ & + \frac{q_d}{\nu^2} \left[ \frac{\kappa}{4\eta} (1 - e^{2\eta t}) + \frac{e^{2\eta t} [(\eta\nu + \nu\kappa) \sin 2\nu t + (\eta\kappa - \nu^2) \cos 2\nu t] - (\eta\kappa - \nu^2)}{4(\eta^2 + \nu^2)} \right]. \end{aligned} \quad (5.88)$$

$$\mathbf{S}^{(2,1)}(t) = \mathbf{S}^{(1,2)}(t) \quad (5.89)$$

Examining the solution given above, we see that the parameter  $\eta$  governs the rate of decay of all of the exponential terms. Therefore, we define the “rise time” as that interval within which the transient response of the covariance will reach a close approximation to the above steady-state value; thus, we define the rise time as follows:

$$t_r = -\frac{3}{\eta}. \quad (5.90)$$

Next, we note that all of the trigonometric terms are modulated by  $2\nu$ ; thus we may view this value as a characteristic damped frequency of the coupled system. The period of the oscillation,  $\Pi$ , is then

$$\Pi = \pi/\nu \quad (5.91)$$

In the limit as  $t \rightarrow \infty$ , all the exponential terms in the analytical solution die out, so that the steady-state value of the covariance simplifies to:

$$\mathbf{P}(\infty) = -\frac{1}{4\eta(\eta^2 + \nu^2)} \begin{bmatrix} q_d + (2\eta^2 + \nu^2 + \kappa^2 - \eta\kappa)q_b & q_b \omega_n^2 (\eta - \kappa) - q_d (\eta + \kappa) \\ q_b \omega_n^2 (\eta - \kappa) - q_d (\eta + \kappa) & (2\eta^2 + \nu^2 + \kappa^2 + \eta\kappa)q_d + q_b \omega_n^4 \end{bmatrix} \quad (5.92)$$

which may be expressed in terms of the original parameters as

$$\begin{aligned} \mathbf{P}(\infty) = & \frac{1}{4\omega_n(\omega_n + 2\beta\zeta)(\zeta\omega_n + \beta/2)} \\ & \cdot \begin{bmatrix} q_d + (\omega_n^2 + 2\beta\zeta\omega_n + 4\zeta^2\omega_n^2)q_b & q_d\beta - 2\zeta\omega_n^3q_b \\ q_d\beta - 2\zeta\omega_n^3q_b & (\omega_n^2 + 2\beta\zeta\omega_n + \beta^2)q_d + \omega_n^4q_b \end{bmatrix} \end{aligned} \quad (5.93)$$

## CHAPTER 6

# State Representations

Contributed by J. Russell Carpenter and Christopher N. D’Souza

This Chapter discusses state representation, primarily for translations; attitude representations are discussed in Chapter 8.

### 6.1. Selection of Solve-For State Variables for Estimation

As has been discussed in Chapter 2, it is not good practice to include unobservable states in the EKF solve-for vector, particularly if this introduces unstable dynamical modes. Nonetheless, during the early stages of designing a navigation filter, it may not be clear to the designer which states to include. There are essentially two approaches to addressing this question, which we may describe as the additive and subtractive methods. With the additive approach, one begins with the smallest possible set of states, adding additional models as one deems them necessary. The problem with this approach is essentially that it is not possible to foresee how additional states will affect the system until they are added; one cannot analyze the sensitivity of the filter’s performance to states which are not present in the analysis. The preferred, subtractive, approach is instead to start with a design of as high a fidelity as practical, including even modes which one may suspect are unobservable and possibly destabilizing. A designer may then readily perform sensitivity and covariance analysis [20, 50] to winnow the solve-for state to as parsimonious a set of observable states as needed to achieve design requirements.

### 6.2. Units and Precision

In the early days of ground-based orbit determination, canonical units were preferred due to the limited word lengths that were available for computation. Factorized filtering methods largely eliminated the need for canonical units even before “modern” double-precision word lengths became available in onboard processors. A renewed interest in single-precision computations has emerged however as the desire to utilize processors based on Field Programmable Gate Arrays has become widespread. Thus, the possibility of overflow, truncation, and roundoff errors must still be considered. Wherever possible, filter computations, especially those involving the covariance matrix (even when it is factorized!), should be done in double precision, and time should be maintained in either two double-precision variables, or in quadruple precision if available. For low-Earth orbit navigation in an Earth-centered frame, position/velocity units based on meters and seconds are often adequate; for applications that may reach into cislunar space and beyond, units based on kilometers and seconds are preferred.

### 6.3. Coordinate and Time Systems

For most orbital navigation applications, use of an “inertial” coordinate frame, such as the International Celestial Reference Frame, the “J2000” (FK5) frame, etc., will be desirable, since onboard computations utilizing navigation filter state estimates will typically most naturally occur in an inertial setting. It is usually convenient to choose a frame whose origin is at the center of the primary gravitational body. Some missions, such as cislunar and interplanetary missions, will occur within the Hill spheres of more than one celestial body, and some mechanism for changing the coordinate system origin without requiring reset of the filter should be considered.

In some cases, consideration may be given to central-body-fixed frames, such as the International Terrestrial Reference Frame, World Geodetic System of 1984, etc., particularly for applications that rely primarily on Global Navigation Satellites Systems (GNSS), and/or ground-based tracking systems. Although integrating the equations of motion in such systems necessitates additional calculations of Coriolis and centripetal acceleration, such calculations are relatively trivial in comparison to the computations required to accurately maintain a transformation between central-body-fixed and inertial frames. Computations of higher-order gravity acceleration are simplified, and maintenance of polar motion coefficients is also eliminated. Several of NASA’s early Global Position System relative navigation experiments used such a formulation successfully [55, 63]. If other onboard applications require inertial states, but are indifferent as to which inertial frame is provided, it may be prudent to consider defining a fixed, true-of-date inertial frame which is identical to the body-fixed frame at the initial power-up of the navigation system, and which is henceforth related to the body-fixed frame by a simple single-axis polar rotation. Such an approach will permit navigation in the body-fixed frame without the difficulties of maintaining a relationship onboard the spacecraft to one of the conventional inertial frames.

With regard to time systems, navigation filter designs should strongly avoid dependence upon discontinuous time scales, such as Coordinated Universal Time (“UTC”). While ground-based applications will generally prefer UTC, it is far easier for the mission’s ground system to manage leap seconds than it is to robustly test and maintain discontinuous time scales in an autonomous onboard navigation setting. The filter designer should strive to ensure that a misapplication of leap second logic can never affect filter performance. If requirements for maintenance of UTC onboard cannot be avoided, all such calculations should occur independently from the uniform continuous time scale that the filter uses internally. Time-tagged commands that affect filter performance should also utilize the same internal, continuous time scale.

### 6.4. Orbit Parameterizations

For orbital navigation applications, orbital elements are geometrically appealing as a state representation, and there exist various “semi-analytic” theories for improving their usefulness as ephemeris representations for real-world orbits, such as the GPS broadcast ephemeris model, two-line elements, etc. Furthermore, long-term evolution of the orbital error covariance more naturally occurs in element representations, which may be especially relevant to conjunction analysis. However, NASA’s experience has been that Cartesian parameters generally offer simpler computational efficiencies for high-fidelity measurement and dynamics models, including for the Jacobian matrices required in estimation algorithms.

### 6.5. Relative State Representations

Although the subject of this Section implies the need for relative state knowledge, it is not necessarily the case that this implies estimation of the relative states directly. For example, if each spacecraft’s only sensor is a GNSS receiver, and there is no method for exchanging the GNSS data between spacecraft, then each satellite’s measurement errors will be largely uncorrelated, assuming that errors in the GNSS constellation data are minimal. Furthermore, there may be no common sources of dynamical error, such as might arise from common yet imperfect models of atmospheric density for low Earth orbiters. In such cases, mission requirements may be met simply by performing isolated state estimation onboard each satellite, and simply differencing the estimated state vectors. In such cases, the covariance of the relative state error between any two spacecraft is given by

$$\begin{aligned} P_{\text{rel}} &= \mathbb{E}[(e_2 - e_1)(e_2 - e_1)^\top] \\ &= P_1 + P_2 \end{aligned} \tag{6.1}$$

where  $e_i$  denotes the estimation error for spacecraft  $i$ , since (by assumption)  $\mathbb{E}[e_1 e_2^\top] = 0$ . For many other applications, either the measurements or the dynamics or both will induce a correlation structure, making it necessary to simultaneously estimate some combination of Earth-centered (a.k.a. “absolute” or “inertial”) states and spacecraft-to-spacecraft relative states. The choice of state representation and associated dynamical model for each application can have significant impacts on efficiency and accuracy, and requires careful consideration.

Aside from the choice of orbit parameterization, there are at least three choices for estimating a relative orbit. The most obvious choice is to solve directly for the differences between the parameters chosen for the orbit representation; that is, to solve for relative position and relative velocity, or relative orbital elements. In some contexts, such as nearly circular orbits, efficient dynamics, such as the linear time-invariant Hill-Clohessy-Wilshire model [11, 26], become available with a choice to solve directly for relative Cartesian states. In many other contexts, higher fidelity may be required, and furthermore, models for drag, solar radiation pressure, the ephemerides of other gravitational bodies, etc. may require knowledge of the Earth-centered (Cartesian) state of one or more of the spacecraft. In such cases, the estimator may solve for a combination of pure absolute/inertial states, or some combination of absolute and relative states. The architecture originally developed for NASA’s *Apollo* missions was the former “dual-inertial” formulation [52]. While the absolute/relative formulation may appear to be mathematically equivalent, computational considerations may choose one or the other to be favored in various application contexts. A general observation is that the dual-inertial formulation may be favorable for computations involving the state and state error covariance, and for “absolute” measurements such as undifferenced GPS pseudorange, while the absolute/relative formulation may be favorable for computations involving satellite-to-satellite relative measurements. Reference [52] provides a comprehensive mathematical description of the dual-inertial formulation in the context of relative range, Doppler, and bearing measurements that one may easily adapt to an other measurement types.

**6.5.1. Dual Inertial State Representation** Here, we consider only two spacecraft, but the results are easily generalized. Let  $x_i = [r_i^\top, v_i^\top]^\top$ ,  $i = 1, 2$  denote the true state of spacecraft  $i$ , with  $r_i, v_i$  the position and velocity vectors expressed in non-rotating coordinates centered on the primary central gravitational body. Based on mission requirements,

any appropriate fidelity of dynamics may be directly utilized per the methods of Chapter 50, e.g.

$$\dot{x}_i = \begin{bmatrix} v_i \\ -\frac{\mu}{\|r_i\|^3} r_i + \sum_j f_j \end{bmatrix} \quad (6.2)$$

where the specific forces  $f_j$  may include thrust, higher-order gravity, drag, solar radiation pressure, gravity from non-central bodies such as the moon and the sun, etc.

Let  $e_i = \hat{x}_i - x_i$ , where  $\hat{x}_i$  is an estimate for the state of spacecraft  $i$ . Then, the error in the state estimate  $\hat{x} = [\hat{x}_1^\top, \hat{x}_2^\top]^\top$  is  $e = [e_1^\top, e_2^\top]^\top$ , and the error covariance is

$$P = \mathbf{E}[ee^\top] = \begin{bmatrix} P_1 & P_{12} \\ P_{12}^\top & P_2 \end{bmatrix} \quad (6.3)$$

Any linear unbiased estimate of  $x$  will have the following measurement update equation:

$$\hat{x}^+ = \hat{x}^- + K(y - h(\hat{x}^-)) \quad (6.4)$$

where  $\hat{x}^-$  is the value of  $\hat{x}$  immediately prior to incorporating the observation,  $y$ , and  $h(\hat{x}^-)$  is an unbiased prediction of the measurement's value. The optimal gain is

$$K = PH^\top(HPH^\top + R)^{-1} \quad (6.5)$$

where  $R$  is the measurement noise covariance and  $H = \partial h(x)/\partial x|_{\hat{x}^-}$ . Partition the update as follows:

$$\begin{bmatrix} \hat{x}_1^+ \\ \hat{x}_2^+ \end{bmatrix} = \begin{bmatrix} \hat{x}_1^- \\ \hat{x}_2^- \end{bmatrix} + \begin{bmatrix} K_1 \\ K_2 \end{bmatrix} (y - h(\hat{x}^-)) \quad (6.6)$$

$$= \begin{bmatrix} \hat{x}_1^- \\ \hat{x}_2^- \end{bmatrix} + \begin{bmatrix} P_1 H_1^\top + P_{12} H_2^\top \\ P_{12}^\top H_1^\top + P_2 H_2^\top \end{bmatrix} (HPH^\top + R)^{-1} (y - h(\hat{x}^-)) \quad (6.7)$$

from which it is clear that the optimal update for the relative state  $\hat{x}_{\text{rel}} = \hat{x}_2 - \hat{x}_1$  is

$$\hat{x}_{\text{rel}}^+ = \hat{x}_{\text{rel}}^- + (P_2 H_2^\top - P_1 H_1^\top - P_{12} H_2^\top + P_{12}^\top H_1^\top) (HPH^\top + R)^{-1} (y - h(\hat{x}^-)) \quad (6.8)$$

with corresponding relative error covariance

$$P_{\text{rel}} = P_1 + P_2 - P_{12} - P_{12}^\top \quad (6.9)$$

Noting that it must be true that  $h(x_{\text{rel}}) = h(x)$  and hence that  $\partial h(x_{\text{rel}})/\partial x_{\text{rel}} = \partial h(x_2)/\partial x_2 = -\partial h(x_1)/\partial x_1$ , let  $H_{\text{rel}} = H_2 = -H_1$ . Then it is clear that

$$P_{\text{rel}} H_{\text{rel}}^\top = P_2 H_2^\top - P_1 H_1^\top - P_{12} H_2^\top + P_{12}^\top H_1^\top \quad (6.10)$$

and that

$$H_{\text{rel}} P_{\text{rel}} H_{\text{rel}}^\top = HPH^\top \quad (6.11)$$

and hence

$$\hat{x}_{\text{rel}}^+ = \hat{x}_{\text{rel}}^- + P_{\text{rel}} H_{\text{rel}}^\top (H_{\text{rel}} P_{\text{rel}} H_{\text{rel}}^\top + R)^{-1} (y - h(\hat{x}_{\text{rel}}^-)) \quad (6.12)$$

Therefore, the dual inertial state update is mathematically (although perhaps not computationally) equivalent to a direct update of the relative state.

Appendix C reproduces a memorandum that further details the benefits of the dual inertial formulation.

**6.5.2. Linearized Relative State Representation** While it is sometimes useful to employ a linear model of the relative dynamics, especially for close proximity operations in near-circular orbits, there is significant benefit to casting the equations of motion in spherical coordinates. The following derivation of the Hill-Clohessy-Wiltshire equations in spherical coordinates is derived from notes from a lecture given by Robert H. Bishop. Let the position of a spacecraft be give by a set of right-handed spherical coordinates

$$r = \rho \begin{bmatrix} \cos \phi \sin \theta \\ \sin \phi \\ \cos \phi \cos \theta \end{bmatrix} \quad (6.13)$$

where  $\rho$  is the distance from the central body to the spacecraft,  $\theta$  is measured along some specified great circle of the central body, and  $\phi$  is measured along a great circle of the central body that is normal to the former great circle, and contains the position vector, as Figure 1 depicts. Define a state vector as follows:  $x = [\rho, \dot{\rho}, \theta, \dot{\theta}, \phi, \dot{\phi}]$ . If the only force on

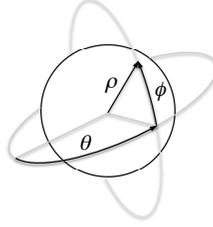


FIGURE 1. Spherical coordinates.

the spacecraft is point-mass gravity from the central body, then the equations of motion are given by

$$\dot{x} = f(x) = \begin{bmatrix} -\mu/\rho^2 + \rho\dot{\phi}^2 + \rho\dot{\theta}^2 \cos^2 \phi \\ \dot{\rho} \\ -2\rho\dot{\theta}/\rho + 2\dot{\phi}\dot{\theta} \tan \phi \\ \dot{\theta} \\ -2\rho\dot{\phi}/\rho - \dot{\theta}^2 \cos \phi \sin \phi \\ \dot{\phi} \end{bmatrix} \quad (6.14)$$

Now consider a circular reference orbit, with radius  $\rho_*$ , which is in the plane of the great circle containing the  $\theta$  coordinate. Let  $\omega_* = \sqrt{\mu/\rho_*^3}$ . Then, the state of an object following the circular reference orbit at any time  $t > t_o$  will be  $x_*(t) = [\rho_*, 0, \omega_*(t - t_o) - \theta_o, \omega_*, 0, 0]$ . Without loss of generality, take  $\theta_o = t_o = 0$ . Letting  $\delta x = x - x_*$ , linearization of (6.14) in the neighborhood of  $x_*$  yields

$$\delta \dot{x}(t) = \left. \frac{\partial f(x)}{\partial x} \right|_{x_*(t)} \delta x(t) = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 3\omega_*^2 & 0 & 0 & 2\omega_*\rho_* & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & -2\omega_*/\rho_* & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & -\omega_*^2 & 0 \end{bmatrix} \delta x(t) \quad (6.15)$$

In this context, it is useful to redefine the state vector  $\tilde{x} = [\rho, \dot{\rho}, \rho_*\theta, \rho_*\dot{\theta}, \rho_*\phi, \rho_*\dot{\phi}]$  so that angles are replaced by arc lengths. Then, the linearized equations of motion become

$$\delta\dot{\tilde{x}}(t) = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 \\ 3\omega_*^2 & 0 & 0 & 2\omega_* & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & -2\omega_* & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & -\omega_*^2 & 0 \end{bmatrix} \delta\tilde{x}(t) \quad (6.16)$$

If linearized relative dynamics are to be used for relative navigation in near-circular orbits, then interpreting the motion along the orbit track, and normal to the orbit track, as arc lengths, per the derivation above, is desirable since it will preserve the linearity of the approximation over a much wider range than if the along-track and cross-track coordinates are taken as rectilinear tangents to the reference orbit position.

### 6.6. Modeling Inertial Components

Many onboard navigation systems employ inertial components consisting of gyros and/or accelerometers. In some applications, an external algorithm will process the increments of rotational and/or translational motion that these devices inherently measure, and produce an acceleration vector that the filter can directly incorporate into its computation of the equations of motion. However, it will often be the case that biases affecting these devices should be estimated by the filter. This section describes some recommended models for such biases, as well as the computations that need to be performed to accumulate the inertial measurement unit's (IMU) angle and velocity increments.

**6.6.1. The Gyro Model** The gyro is modeled in terms of the bias, scale factor and non-orthogonality. The IMU *case frame* is defined such that the  $x$ -axis of the gyro is the reference direction with the  $x - y$  plane being the reference plane; the  $y$ - and  $z$ -axes are not mounted perfectly orthogonal to it (this is why we don't have a full misalignment/nonorthogonality matrix as we will in the accelerometer model). The errors in determining these misalignments are the so-called *non-orthogonality errors*, expressed as a matrix  $\mathbf{\Gamma}$ , as

$$\mathbf{\Gamma} \triangleq \begin{bmatrix} 0 & 0 & 0 \\ \gamma_{yx} & 0 & 0 \\ \gamma_{zx} & \gamma_{zy} & 0 \end{bmatrix}$$

The gyro scale factor represents the error in conversion from raw sensor outputs (gyro digitizer pulses) to useful units. In general we model the scale-factor error as a first-order Markov (or a Gauss-Markov) process in terms of a diagonal matrix given as

$$\mathbf{S}^g = \begin{bmatrix} s_x^g & 0 & 0 \\ 0 & s_y^g & 0 \\ 0 & 0 & s_z^g \end{bmatrix}$$

Similarly, the gyro bias errors are modeled as as first-order vector Gauss-Markov processes as

$$\mathbf{b}^g = \begin{bmatrix} b_x^g \\ b_y^g \\ b_z^g \end{bmatrix}$$

Finally, the gyro noise is represented by  $\epsilon_g$ . Hence

$$\omega_m^{\mathcal{C}} = (\mathbf{I}_3 + \mathbf{\Gamma} + \mathbf{S}^g) (\omega^{\mathcal{C}} + \mathbf{b}_g + \epsilon_g) \quad (6.17)$$

where  $\mathbf{I}_3$  is a  $3 \times 3$  identity matrix, the superscript  $\mathcal{C}$  indicates that this is an inertial measurement at the ‘box-level’ expressed in *case-frame* co-ordinates, and  $\omega^{\mathcal{C}}$  is the ‘true’ angular velocity in the case frame. If we let  $\mathbf{\Delta}^g \triangleq \mathbf{\Gamma} + \mathbf{S}^g$  and  $(\mathbf{I} + \mathbf{\Delta}^g)^{-1} \approx \mathbf{I} - \mathbf{\Delta}^g$ <sup>‡</sup>, we can express the actual angular velocity in terms of the measured angular velocity as

$$\omega^{\mathcal{C}} = (\mathbf{I}_3 - \mathbf{\Delta}^g) \omega_m^{\mathcal{C}} - \mathbf{b}_g - \epsilon_g \quad (6.18)$$

**6.6.2. The Accumulated  $\Delta\theta$**  In order to find the accumulated angle (not as a function of the measurement, but purely as a function of the true angular velocity), we define  $\Delta\theta$  as

$$\left( \Delta\theta_{\mathcal{C}_{k-1}}^{\mathcal{C}_k} \right)_m \triangleq \int_{t_{k-1}}^{t_k} \omega_m^{\mathcal{C}}(\tau) + \frac{1}{2} \phi_{\mathcal{C}_{ref}}^{\mathcal{C}} \times \omega_m^{\mathcal{C}}(\tau) d\tau \quad (6.19)$$

$$= \int_{t_{k-1}}^{t_k} \omega_m^{\mathcal{C}}(\tau) + \frac{1}{2} \left[ \int_{t_{k-1}}^{\tau} \dot{\phi}_{\mathcal{C}_{ref}}^{\mathcal{C}}(\chi) d\chi \right] \times \omega_m^{\mathcal{C}}(\tau) d\tau \quad (6.20)$$

$$= \int_{t_{k-1}}^{t_k} \omega_m^{\mathcal{C}}(\tau) + \frac{1}{2} \left[ \int_{t_{k-1}}^{\tau} \left( \omega_m^{\mathcal{C}}(\chi) + \frac{1}{2} \phi_{\mathcal{C}_{ref}}^{\mathcal{C}} \times \omega_m^{\mathcal{C}}(\chi) \right) d\chi \right] \times \omega_m^{\mathcal{C}}(\tau) d\tau$$

Ignoring second-order terms, we get

$$\left( \Delta\theta_{\mathcal{C}_{k-1}}^{\mathcal{C}_k} \right)_m = \int_{t_{k-1}}^{t_k} \left[ \omega_m^{\mathcal{C}}(\tau) + \frac{1}{2} \int_{t_{k-1}}^{\tau} \omega_m^{\mathcal{C}}(\chi) d\chi \times \omega_m^{\mathcal{C}}(\tau) \right] d\tau \quad (6.21)$$

With this expression, we find that, by analogy, we can express  $\left( \Delta\theta_{\mathcal{C}_{k-1}}^{\mathcal{C}_k} \right)$  as

$$\left( \Delta\theta_{\mathcal{C}_{k-1}}^{\mathcal{C}_k} \right) = \int_{t_{k-1}}^{t_k} \left[ \omega^{\mathcal{C}}(\tau) + \frac{1}{2} \int_{t_{k-1}}^{\tau} \omega^{\mathcal{C}}(\chi) d\chi \times \omega^{\mathcal{C}}(\tau) \right] d\tau \quad (6.22)$$

**6.6.3. The Accelerometer Model** The accelerometer package will likely be mis-aligned relative to the IMU reference frame. This is due to the fact that the three accelerometers (contained in the accelerometer package) are not mounted orthogonal to each

<sup>‡</sup>In order to evaluate  $(\mathbf{I} + \mathbf{\Delta})^{-1}$  we recall the Woodbury matrix identity

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1}$$

Using this, we obtain the following relation (with  $\mathbf{A} = \mathbf{I}$ ,  $\mathbf{U} = \mathbf{\Delta}$ ,  $\mathbf{C} = \mathbf{I}$  and  $\mathbf{V} = \mathbf{I}$ ),

$$(\mathbf{I} + \mathbf{\Delta})^{-1} = \mathbf{I} - \mathbf{\Delta}(\mathbf{I} + \mathbf{\Delta})^{-1}$$

which, worked recursively, yields the following approximation

$$(\mathbf{I} + \mathbf{\Delta})^{-1} = \mathbf{I} - \mathbf{\Delta} + \mathbf{\Delta}^2 - \mathbf{\Delta}^3 + \mathbf{\Delta}^4 - \mathbf{\Delta}^5 + \dots$$

Therefore, to first-order (neglecting second-order and higher terms in the above equation), we get

$$(\mathbf{I} + \mathbf{\Delta})^{-1} \approx \mathbf{I} - \mathbf{\Delta}$$

other and these errors are expressed in terms of six different small angles as:

$$\mathbf{\Xi}^a = \begin{bmatrix} 0 & \xi_{xy}^a & \xi_{xz}^a \\ \xi_{yx}^a & 0 & \xi_{yz}^a \\ \xi_{zx}^a & \xi_{zy}^a & 0 \end{bmatrix}$$

Similar to the gyros, the accelerometer scale factor represents the error in conversion from raw sensor outputs (accelerometer digitizer pulses) to useful units. In general we model the scale-factor error as a first-order (Gauss-) Markov process in terms of a diagonal matrix given as

$$\mathbf{S}^a = \begin{bmatrix} s_x^a & 0 & 0 \\ 0 & s_y^a & 0 \\ 0 & 0 & s_z^a \end{bmatrix}$$

Similarly, the bias errors are modeled as as first-order Gauss-Markov processes as

$$\mathbf{b}^a = \begin{bmatrix} b_x^a \\ b_y^a \\ b_z^a \end{bmatrix}$$

So, the accelerometer measurements,  $\mathbf{a}_m^{\mathcal{C}}$  are modeled as:

$$\mathbf{a}_m^{\mathcal{C}} = (\mathbf{I}_3 + \mathbf{\Xi}^a) (\mathbf{I}_3 + \mathbf{S}^a) (\mathbf{a}^{\mathcal{C}} + \mathbf{b}^a + \mathbf{v}_a) \quad (6.23)$$

where  $\mathbf{I}_3$  is a  $3 \times 3$  identity matrix, the superscript  $\mathcal{C}$  indicates that this is an inertial measurement at the ‘box-level’ expressed in *case-frame* co-ordinates, and  $\mathbf{a}^{\mathcal{C}}$  is the ‘true’ non-gravitational acceleration in the case frame. The quantity  $\mathbf{v}_a$  is the velocity random walk, a zero-mean white sequence on acceleration that integrates into a velocity random walk, which is the ‘noise’ on the accelerometer output. If we assume that the errors are small, then to first-order

$$(\mathbf{I}_3 + \mathbf{\Xi}^a) (\mathbf{I}_3 + \mathbf{S}^a) \approx \mathbf{I} + \mathbf{\Xi}^a + \mathbf{S}^a$$

So, the linear accelerometer measurements (in the case frame) are:

$$\mathbf{a}_m^{\mathcal{C}} = (\mathbf{I}_3 + \mathbf{\Xi}^a + \mathbf{S}^a) (\mathbf{a}^{\mathcal{C}} + \mathbf{b}^a + \mathbf{v}_a) \quad (6.24)$$

**6.6.4. Accumulated  $\Delta \mathbf{v}$**  We note that the measured  $\Delta \mathbf{v}$  in the case frame,  $\Delta \mathbf{v}_m^{\mathcal{C}}$ , is mapped to the end of its corresponding time interval by the sculling algorithm within the IMU firmware, so that we can write

$$(\Delta \mathbf{v}_m^{\mathcal{C}})_k = \int_{t_{k-1}}^{t_k} \mathbf{T}_{\mathcal{C}(t)}^{\mathcal{C}_k} \mathbf{a}_m^{\mathcal{C}(t)} dt \quad (6.25)$$

where  $(\Delta \mathbf{v}_m^{\mathcal{C}})_k$  covers the time interval from  $t_{k-1}$  to  $t_k$  ( $t_k > t_{k-1}$ ) and  $\mathcal{C}(t)$  is the instantaneous case frame<sup>§</sup>. We recall that a transformation matrix can be written in terms of the

<sup>§</sup>Or equivalently,

$$(\Delta \mathbf{v}_m^B)_k = \int_{t_{k-1}}^{t_k} \mathbf{T}_{B(t)}^{B_k} \mathbf{a}_m^{B(t)} dt \quad (6.26)$$

But since  $\mathbf{T}_{B(t)}^{B_k} \approx \mathbf{I}_3 - [\phi_{B(t)}^{B_k} \times]$ , we find

$$(\Delta \mathbf{v}_m^B)_k = \int_{t_{k-1}}^{t_k} [\mathbf{I}_3 - [\phi_{B(t)}^{B_k} \times]] \mathbf{a}_m^{B(t)} dt \quad (6.27)$$

Euler axis/angle as

$$T(\boldsymbol{\phi}) = \cos(\phi)\mathbf{I} - \frac{\sin \phi}{\phi} [\boldsymbol{\phi} \times] + \frac{1 - \cos \phi}{\phi^2} \boldsymbol{\phi} \boldsymbol{\phi}^\top \quad (6.28)$$

$$= \mathbf{I} - \frac{\sin \phi}{\phi} [\boldsymbol{\phi} \times] + \frac{1 - \cos \phi}{\phi^2} [\boldsymbol{\phi} \times] [\boldsymbol{\phi} \times] \quad (6.29)$$

which, for  $\boldsymbol{\phi} \sim \mathbf{0}$  can be approximated as

$$T(\boldsymbol{\phi}) = \mathbf{I} - [\boldsymbol{\phi} \times] \quad (6.30)$$

With this in mind,  $\mathbf{T}_{\mathcal{C}(t)}^{\mathcal{C}_k} = \mathbf{I}_3 - [\boldsymbol{\theta}_{\mathcal{C}(t)}^{\mathcal{C}_k} \times]$ ,  $(\Delta \mathbf{v}_m^B)_k$  using Eq. (6.24), becomes

$$(\Delta \mathbf{v}_m^{\mathcal{C}})_k = \int_{t_{k-1}}^{t_k} \left[ \mathbf{I}_3 - [\boldsymbol{\theta}_{\mathcal{C}(t)}^{\mathcal{C}_k} \times] \right] [(\mathbf{I}_3 + \boldsymbol{\Delta}^a) \mathbf{a}^{\mathcal{C}} + \mathbf{b}^a + \mathbf{v}_a] dt \quad (6.31)$$

We can expand this equation, neglecting terms of second-order, as follows

$$\begin{aligned} (\Delta \mathbf{v}_m^{\mathcal{C}})_k &= \int_{t_{k-1}}^{t_k} \left[ \mathbf{I}_3 - [\boldsymbol{\theta}_{\mathcal{C}(t)}^{\mathcal{C}_k} \times] \right] \mathbf{a}^{\mathcal{C}} dt + \int_{t_{k-1}}^{t_k} (\mathbf{b}^a + \mathbf{v}_a) dt \\ &\quad + \int_{t_{k-1}}^{t_k} \boldsymbol{\Delta}^a \mathbf{a}^{\mathcal{C}} dt \end{aligned} \quad (6.32)$$

The first term in the above equation (Eq. (6.32)) becomes

$$\int_{t_{k-1}}^{t_k} \left[ \mathbf{I}_3 - [\boldsymbol{\theta}_{\mathcal{C}(t)}^{\mathcal{C}_k} \times] \right] \mathbf{a}^{\mathcal{C}} dt = (\Delta \mathbf{v}^{\mathcal{C}})_k \quad (6.33)$$

and the third term becomes

$$\int_{t_{k-1}}^{t_k} \boldsymbol{\Delta}^a \mathbf{a}^{\mathcal{C}} dt = \boldsymbol{\Delta}^a \int_{t_{k-1}}^{t_k} \mathbf{a}^{\mathcal{C}} dt \approx \boldsymbol{\Delta}^a (\Delta \mathbf{v}^{\mathcal{C}})_k \quad (6.34)$$

Finally, the accelerometer noise, which is zero-mean process with spectral density  $\mathbf{S}_a$  becomes

$$\int_{t_k}^{t_{k+1}} \mathbf{v}_a dt = \mathbf{u}_a \quad (6.35)$$

where  $\mathbf{u}_a$  is a random vector with covariance  $\mathbf{S}_a(t_k - t_{k-1})$ . So, Eq. (6.32) becomes

$$(\Delta \mathbf{v}_m^{\mathcal{C}})_k = [\mathbf{I}_3 + \boldsymbol{\Delta}^a] (\Delta \mathbf{v}^{\mathcal{C}})_k + \mathbf{b}^a \Delta t + \mathbf{v}_a \Delta t \quad (6.36)$$

Since we have established that  $[\mathbf{I}_3 + \boldsymbol{\Delta}^a]^{-1} \approx [\mathbf{I}_3 - \boldsymbol{\Delta}^a]$ , and neglecting terms of second-order,

$$(\Delta \mathbf{v}^{\mathcal{C}})_k = [\mathbf{I}_3 - \boldsymbol{\Delta}^a] (\Delta \mathbf{v}_m^{\mathcal{C}})_k - \mathbf{b}^a \Delta t - \mathbf{v}_a \Delta t \quad (6.37)$$

The average acceleration in the case frame is

$$\mathbf{a}_{ave}^{\mathcal{C}} = \frac{(\Delta \mathbf{v}^{\mathcal{C}})_k}{\Delta t} \quad (6.38)$$

and the average measured acceleration in the case frame is

$$(\mathbf{a}_m^{\mathcal{C}})_{ave} = \frac{(\Delta \mathbf{v}_m^{\mathcal{C}})_k}{\Delta t} \quad (6.39)$$

so we find that

$$\mathbf{a}_{ave}^C = [\mathbf{I}_3 - \Delta^a] (\mathbf{a}_m^C)_{ave} - \mathbf{b}^a - \mathbf{v}_a \quad (6.40)$$

Recalling that the IMU measures accelerations except for gravity, total acceleration is

$$\mathbf{a}^I = \mathbf{g}^I(\mathbf{r}) + \left( \mathbf{T}_{Bref}^I \right)_k \mathbf{T}_B^{Bref} \mathbf{T}_C^B \mathbf{a}_{ave}^C \quad (6.41)$$

**6.6.5. The Gravity Call** One of the more expensive computations involving the propagation of the trajectory with IMU data is the gravity calculation. This is particularly acute when the gravity field used is of high order. The gravity gradient matrix requires even more computation. Hence it goes without saying that if a way is found to minimize the gravity calls, that would make the navigation software more tractable. Taking advantage of the fact that propagation of the trajectory using IMU data occurs at a high rate (usually at 40 Hz or higher), we expand the gravity vector in terms of a Taylor series about  $\mathbf{r}^*$  as

$$\mathbf{g}(\mathbf{r}) = \mathbf{g}(\mathbf{r}^*) + \left. \frac{\partial \mathbf{g}}{\partial \mathbf{r}} \right|_{\mathbf{r}=\mathbf{r}^*} [\mathbf{r} - \mathbf{r}^*] + \frac{1}{2} [\mathbf{r} - \mathbf{r}^*]^\top \left. \frac{\partial^2 \mathbf{g}}{\partial \mathbf{r}^2} \right|_{\mathbf{r}=\mathbf{r}^*} [\mathbf{r} - \mathbf{r}^*] + \dots \quad (6.42)$$

Knowing that the gravity gradient matrix,  $\mathbf{G}$ , is

$$\mathbf{G}(\mathbf{r}^*) = \left. \frac{\partial \mathbf{g}}{\partial \mathbf{r}} \right|_{\mathbf{r}=\mathbf{r}^*} \quad (6.43)$$

and truncating after the first-order in  $\mathbf{r}$ , we find that

$$\mathbf{g}(\mathbf{r}) \approx \mathbf{g}(\mathbf{r}^*) + \mathbf{G}(\mathbf{r}^*) [\mathbf{r} - \mathbf{r}^*] \quad (6.44)$$

$$\approx \mathbf{G}(\mathbf{r}^*) \mathbf{r} + [\mathbf{g}(\mathbf{r}^*) - \mathbf{G}(\mathbf{r}^*) \mathbf{r}^*] \quad (6.45)$$

where now the gravity vector and the gravity gradient matrix need only to be evaluated at the beginning of the major cycle.

## CHAPTER 7

# Factorization Methods

Contributed by Christopher D'Souza

Of the various covariance factorization methods\*, the  $UDU$  covariance factorization technique is among the most commonly used covariance matrix factorization methodologies used in practice. It is implemented in GEONS and flew on MMS and is the heart of the Orion Absolute Navigation System. This chapter is intended to present the  $UDU$  triangular factorization method and the rationale for its use.

Above all, we demonstrate that the  $UDU$  factorization results in a significant reduction in the arithmetic operations (specifically adds and multiplies) compared with the usual  $\Phi\bar{P}\Phi^T + Q$  time update and the Joseph measurement update.

In the next section, we present some notational and preliminary operations for the matrix factors  $\mathbf{U}$  and  $\mathbf{D}$ . In the section that follows, we will derive the time update equations for the aforementioned covariance matrix factors. Next, we will derive the measurement update equations for the covariance matrix factors. Finally, we will present some concluding comments.

### 7.1. Why Use the $UDU$ Factorization?

The usual Kalman filter equations work well for rather simple problems. But once the state-space becomes large, the condition numbers of the covariance matrix becomes large and nonlinear effects begin to affect the numerical characteristics, problems such as filter divergence and non-positive definiteness of the covariance matrix occur. These issues began to be observed almost as soon as Kalman filters began to be used in real problems. Matrix factorization techniques were introduced to solve (at least) some of these issues. The earliest was the Potter Square Root Factorization, which was used in the on-board Apollo navigation filters.

In fact Bierman and Thornton, in a 1976 JPL Report, rather cheekily compare those who insist on using the conventional Kalman filtering and batch least-squares algorithms (*contra* the matrix factorization algorithms) to unrepentant smokers by describing “*an attitude often encountered among estimation practitioners [is] that they will switch to the more accurate and stable algorithms if and when numerical problems occur. An analogy comes to mind of a smoker who promises to stop when cancer or heart ailment symptoms are detected. To expand on the analogy, one may note the following:*

- *Most smokers do not get cancer or heart disease. (Most applications of the Kalman algorithms work.)*

---

\*Other options include the Square Root Covariance Factorization and the Square Root Information Filter (SRIF).

- *Even when catastrophic illness does not occur, there is diminished health. (Even when algorithms work, performance may be degraded.)*
- *Smokers can take precautions to lessen the danger, such as smoking low tar or filtered cigarettes. (Engineers can scale their variables to reduce the dynamic range or use double-precision arithmetic.)*
- *Lung cancer may not be diagnosed until it is too advanced for treatment. (Numerical problems may not be detected in time to be remedied.)* ” [71]

In addition, a little advertised, but incredibly useful feature of the UDU factorization is the ability to interrogate for the positive definiteness of the covariance matrix for ‘free’, since definiteness of the  $\mathbf{D}$  matrix may be trivially evaluated.

This sets the stage for the need for the matrix factorization techniques and the UDU technique in particular.

## 7.2. Preliminaries

Let us factor a covariance matrix,  $\mathbf{P}$ , into the following form

$$\mathbf{P} = \mathbf{U}\mathbf{D}\mathbf{U}^T \quad (7.1)$$

where  $\mathbf{U}$  is an upper triangular matrix with 1’s on the diagonals and 0’s on the lower portion of the matrix,  $\mathbf{D}$  is a diagonal matrix. We can write  $\mathbf{U}$  and  $\mathbf{D}$  compactly as

$$\mathbf{U} = \{u_{ij}\}, \quad i < j \quad (7.2)$$

$$\mathbf{D} = \{d_{ii}\} \quad (7.3)$$

as well

$$u_{ii} = 1 \quad (7.4)$$

It should be noted that Eq. (7.2) gives the upper triangular portion of the covariance; the lower triangular matrix can be obtained by reflection or by evaluation of Eq. (7.2), with  $u_{lm} = 0$  for  $l > m$ .

Equally valid is

$$p_{ij} = \sum_{k=i}^n u_{ik}d_{kk}u_{jk}, \quad j < i \quad (7.5)$$

and for the diagonals we find,

$$p_{ii} = \sum_{k=i}^n u_{ik}^2 d_{kk} \quad (7.6)$$

So, given an  $n \times n$  symmetric, positive semi-definite matrix  $\mathbf{P}$ , the unit upper triangular factor  $\mathbf{U}$  and the diagonal factor  $\mathbf{D}$  (such that  $\mathbf{P} = \mathbf{U}\mathbf{D}\mathbf{U}^T$ ) is obtained using the following equations. We begin with the  $(n, n)$  element and work upwards (along the columns).

The algorithm is as follows:

```

dn,n = pn,n
un,n = 1.0
for j = n - 1 : -1 : 1 do
    uj,n = pj,n/dn,n
end for
for j = n : -1 : 2 do
    dj,j = pj,j
    for k = j + 1 : -n do

```

```

     $d_{j,j} = d_{j,j} + d_{k,k}u_{j,k}^2$ 
end for
 $u_{j,j} = 1.0$ 
for  $i = j - 1 : -1 : 1$  do
     $u_{i,j} = p_{i,j}$ 
    for  $k = j + 1 : n$  do
         $u_{i,j} = u_{i,j} + d_{k,k}u_{j,k}u_{i,k}$ 
    end for
     $u_{i,j} = u_{i,j}/d_{j,j}$ 
end for
end for

```

A word of caution: Maybeck [50] on page 392, rather unusually and uncharacteristically, calls the matrix  $\mathbf{U}$  a unitary matrix, ostensibly because there are 1's on the diagonals. Strictly speaking a if  $\mathbf{U}$  were a unitary matrix,  $\mathbf{U}\mathbf{U}^\top = \mathbf{U}^\top\mathbf{U} = \mathbf{I}$ , which is clearly not the case for the  $\mathbf{U}$  in the UDU factorization.

In practice, particularly when storage limitations are driving the design, the  $n \times n$  matrix  $\mathbf{D}$ , which is a diagonal matrix, can be stored as a  $n$ -vector. Likewise, the matrix  $n \times n$  matrix  $\mathbf{U}$  which is upper triangular with 1's on the diagonal, can be stored as a  $n(n-1)/2$  vector. The storage savings can be particularly significant as  $n$  increases. Of course the algorithms need to be designed to ensure that the entries of  $\mathbf{U}$  above the diagonal are the only ones used in the computations.

To complicate matters further, if one uses the Matlab `qr` function to triangularize a matrix  $\mathbf{A}$ , it outputs two matrices,  $\mathbf{Q}$  and  $\mathbf{R}$ , of which  $\mathbf{Q}$  is a unitary matrix in the 'classic' mathematical definition (*i.e.*  $\mathbf{Q}\mathbf{Q}^\top = \mathbf{Q}^\top\mathbf{Q} = \mathbf{I}$ ) and  $\mathbf{R}$  is an upper triangular matrix so that  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ .

### 7.3. The Time Update of the Covariance

As is necessary in Kalman Filtering, we wish to propagate the UDU factorization of the covariance matrix. We loosely follow Maybeck [50] in this development. First, we will pose the more general problem and then we will specialize it for navigation problems with a large number of sensor biases.

We begin by expressing the equations for the general time update problem. Next, we specialize the general problem to the case where a subset of states, which we will call 'parameters', whose dynamics are uncorrelated with any other state other than themselves. Finally, we present the arithmetic operation (numbers of adds, multiplies, and divides) of the time update of the covariance matrix.

**7.3.1. The General Time Update Problem** Given a state,  $\mathbf{x}$ , that evolves according to

$$\mathbf{x}(t_k) = \Phi(t_k, t_{k-1})\mathbf{x}(t_{k-1}) + \mathbf{G}_k\mathbf{w}_k$$

where  $\mathbf{w}_k$  is the process noise at time  $t_k$ , where  $\mathbf{x}$  is an  $n \times 1$  vector, and  $\mathbf{w}$  is an  $m \times 1$  vector. With this in hand, the general problem is as follows [70]: we wish to propagate the covariance matrix defined by

$$\bar{\mathbf{P}}(t_k) = \Phi(t_k, t_{k-1})\mathbf{P}(t_{k-1})\Phi^\top(t_k, t_{k-1}) + \mathbf{G}_k\mathbf{Q}_k\mathbf{G}_k^\top \quad (7.7)$$

where  $\bar{\mathbf{P}}$  is the propagated covariance (the overbar indicates a propagated quantity) and  $\mathbf{P}$  is the updated covariance at the prior time step,  $\mathbf{Q}_k$  is the diagonal process noise covariance

matrix and  $\mathbf{G}_k$  is the mapping of the noise to the state. To save memory, since  $\mathbf{Q}_k$  and  $\mathbf{D}_k$  are diagonal matrices, in the implementation, we pass  $\mathbf{Q}_k$  and  $\mathbf{D}_k$  as vectors.

We want to find the propagated factors  $\bar{\mathbf{U}}_k$  and  $\bar{\mathbf{D}}_k$ , such that  $\bar{\mathbf{P}}_k \triangleq \bar{\mathbf{U}}_k \bar{\mathbf{D}}_k \bar{\mathbf{U}}_k^\top$ . For compactness, we now drop the time subscripts. Given the  $UDU$  factorization of covariance matrices, we can write Eq. (7.7) as

$$\bar{\mathbf{U}}_k \bar{\mathbf{D}}_k \bar{\mathbf{U}}_k^\top = \Phi_k \mathbf{U}_{k-1} \mathbf{D}_{k-1} \mathbf{U}_{k-1}^\top \Phi_k^\top + \mathbf{G}_k \mathbf{Q}_k \mathbf{G}_k^\top \quad (7.8)$$

$$= \begin{bmatrix} \Phi_k \mathbf{U}_{k-1} & \mathbf{G}_k \end{bmatrix} \begin{bmatrix} \mathbf{D}_{k-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_k \end{bmatrix} \begin{bmatrix} \Phi_k \mathbf{U}_{k-1} & \mathbf{G}_k \end{bmatrix}^\top \quad (7.9)$$

Since  $\mathbf{x}$  is an  $n \times 1$  vector, and  $\mathbf{w}$  is an  $m \times 1$  vector,  $\Phi_k$  is an  $n \times n$  matrix,  $\mathbf{G}_k$  is an  $n \times m$  matrix, and  $\mathbf{Q}_k$  is an  $m \times m$  matrix.

We recall that both  $\bar{\mathbf{U}}_k$  and  $\mathbf{U}_{k-1}$  are  $n \times n$  upper triangular matrices, with 1's on the diagonal and  $\bar{\mathbf{D}}$  and  $\mathbf{D}$  are purely diagonal matrices. So, we have some work to do on Eq. (7.9) because  $\begin{bmatrix} \Phi_k \mathbf{U}_{k-1} & \mathbf{G}_k \end{bmatrix}$  is an  $n \times (n+m)$  matrix and  $\begin{bmatrix} \mathbf{D}_{k-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_k \end{bmatrix}$  is an  $(n+m) \times (n+m)$  diagonal matrix.

Defining the first matrix ( $\begin{bmatrix} \Phi_k \mathbf{U}_{k-1} & \mathbf{G}_k \end{bmatrix}$  which is an  $n \times (n+m)$  matrix) on the right hand side of Eq.(7.9) as

$$\mathbf{Y} \triangleq \begin{bmatrix} \Phi_k \mathbf{U}_{k-1} & \mathbf{G}_k \end{bmatrix} \quad (7.10)$$

we seek a matrix  $\mathbf{T}_k$  that transforms  $\mathbf{Y}_k$  such that

$$\mathbf{Y}_k \mathbf{T}_k^{-\top} = \begin{bmatrix} \bar{\mathbf{U}}_k & \mathbf{0}_{n \times m} \end{bmatrix} \quad (7.11)$$

where  $\bar{\mathbf{U}}_k$  is an  $n \times n$  upper triangular matrix with 1's on the diagonal. In order to find the desired matrix  $\mathbf{T}_k$  we perform a Gram-Schmidt orthogonalization<sup>‡</sup>. We define the diagonal matrix  $\tilde{\mathbf{D}}_k$  as

$$\tilde{\mathbf{D}}_k \triangleq \begin{bmatrix} \mathbf{D}_{k-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_k \end{bmatrix} \quad (7.12)$$

which will be important as we define the weighted inner product in the Gram-Schmidt Orthogonalization. Eq. (7.9) which was

$$\bar{\mathbf{U}}_k \bar{\mathbf{D}}_k \bar{\mathbf{U}}_k^\top = \mathbf{Y} \tilde{\mathbf{D}}_k \mathbf{Y}^\top \quad (7.13)$$

can now be rewritten as

$$\bar{\mathbf{U}}_k \bar{\mathbf{D}}_k \bar{\mathbf{U}}_k^\top = \mathbf{Y}_k \mathbf{T}_k^{-\top} \left[ \mathbf{T}_k^\top \tilde{\mathbf{D}}_k \mathbf{T}_k \right] \mathbf{T}_k^{-1} \mathbf{Y}_k^\top \quad (7.14)$$

We note that the matrix  $\left[ \mathbf{T}_k^\top \tilde{\mathbf{D}}_k \mathbf{T}_k \right]$  is a diagonal matrix.

“All” that remains is to find  $\mathbf{T}_k$ . This is where we harness the power of the modified Gram-Schmidt orthogonalization process which provides us with what we were after:  $\bar{\mathbf{U}}_k$  and  $\bar{\mathbf{D}}_k$ .

---

<sup>‡</sup>Whereas we can use a **qr** factorization, we specifically perform a Gram-Schmidt factorization specialized to the  $UDU$  factorization on the correlated states.

7.3.1.1. *The Modified Gram-Schmidt Algorithm* Given the  $(n + m) \times n$  matrix  $\mathbf{Y}_k = [\Phi_k \mathbf{U}_{k-1} \quad \mathbf{G}_k]$ ,  $\Phi_k \mathbf{U}_{k-1}$  can be constructed taking advantage of the structure of  $\mathbf{U}_{k-1}$  which is an upper triangular matrix with 1's on the diagonal, with  $\frac{1}{2}(n^3 - n^2)$  adds and  $\frac{1}{2}(n^3 - n^2)$  multiplies. We recall that  $\mathbf{Y}_k$  and  $\mathbf{b}_k$  are  $(n + m) \times 1$  vectors. In the following algorithm, the number of adds, multiplies and divides as a consequence of each operation is expressed in terms of ‘[adds, multiplies, divides]’. We only go to  $j = 2$  because  $\bar{\mathbf{U}}_{11} = 1$ . The MGS algorithm can be expressed as:

```

for  $k = n, \dots, 2$  do
     $\mathbf{b}_k = \mathbf{Y}_k$ 

end for

for  $j = n, \dots, 2$  do
     $\mathbf{f}_j = \tilde{\mathbf{D}}\mathbf{b}_j$                                  $[0, n(n + m), 0]$ 
     $\bar{\mathbf{D}}_{jj} = \mathbf{b}_j^\top \mathbf{f}_j$                          $[n(n + m), n(n + m), 0]$ 
     $\mathbf{f}_j = \mathbf{f}_j / \bar{\mathbf{D}}_{jj}$                            $[0, 0, (n + m)(n - 1)]$ 

    for  $i = 1, \dots, j - 1$  do
         $\bar{\mathbf{U}}_{ij} = \mathbf{b}_i^\top \mathbf{f}_j$                      $[(n + m) \frac{(n^2 - n)}{2}, (n + m) \frac{(n^2 - n)}{2}, 0]$ 
         $\mathbf{b}_i = \mathbf{b}_i - \bar{\mathbf{U}}_{ij} \mathbf{b}_j$                  $[(n + m) \frac{(n^2 - n)}{2}, (n + m) \frac{(n^2 - n)}{2}, 0]$ 
    end for

     $\bar{\mathbf{U}}_{11} = 1$ 
     $\mathbf{f}_1 = \tilde{\mathbf{D}}\mathbf{b}_1$ 
     $\bar{\mathbf{D}}_{11} = \mathbf{b}_1^\top \mathbf{f}_1$ 
end for

```

Thus, the algorithm not only provides the orthogonal basis vectors,  $\mathbf{b}_j, j = 1, \dots, n_{\mathbf{x}}$ , but it also provides the triangular matrix factors  $\bar{\mathbf{U}}$  and  $\bar{\mathbf{D}}$ .

Since we are also interested in the arithmetic operations, we find that there are  $[n_{\mathbf{x}}^3 + n_{\mathbf{x}}^2 m_{\mathbf{x}}]$  adds,  $[n_{\mathbf{x}}(n_{\mathbf{x}} + 1)(n_{\mathbf{x}} + m_{\mathbf{x}})]$  multiplies and  $[(n_{\mathbf{x}} + m_{\mathbf{x}})(n_{\mathbf{x}} - 1)]$  divides. For the case when  $m_{\mathbf{x}} = 0$ , i.e. no process noise, we have  $n_{\mathbf{x}}^3$  adds and  $[n_{\mathbf{x}}^3 + n_{\mathbf{x}}^2]$  multiplies and  $[n_{\mathbf{x}}^2 - n_{\mathbf{x}}]$  divides.

Finally, the entire covariance update algorithm, including the computation of  $\mathbf{Y}_k$  uses  $[1.5n_{\mathbf{x}}^3 + 0.5n_{\mathbf{x}}^2(2m_{\mathbf{x}} - 1)]$ , and  $[0.5n_{\mathbf{x}}^2(3n_{\mathbf{x}} + 1) + n_{\mathbf{x}}m_{\mathbf{x}}(n_{\mathbf{x}} + 1)]$  multiplies and  $[(n_{\mathbf{x}} + m_{\mathbf{x}})(n_{\mathbf{x}} - 1)]$  divides.

The Modified Gram-Schmidt orthogonalization process makes no assumptions regarding the structure of  $\Phi$ . For a large number of states (say,  $n_{\mathbf{x}} = 35$  with process noise inputs,  $m_{\mathbf{x}} = 35$ ), most of which might be biases (or Gauss Markov processes), much of this is wasted considering the sparseness of  $\Phi$ . In Appendix B.1, we show that we use  $[n_{\mathbf{x}}^3 + n_{\mathbf{x}}^2 m_{\mathbf{x}}]$  adds and  $[n_{\mathbf{x}}(n_{\mathbf{x}} + 1)(n_{\mathbf{x}} + m_{\mathbf{x}})]$  multiplies to obtain  $\bar{\mathbf{U}}_k$  and  $\bar{\mathbf{D}}_k$ . For  $n_{\mathbf{x}} = 35$  and  $m_{\mathbf{x}} = 35$ , we require 85,750 adds and 88,200 multiplies – quite a large number of computations. We can stop here and all will be well – if we are willing to pay the heavy computational price.

But we can do better! We can vastly improve (reduce) on the number of computations by partitioning the original state vector into ‘states’ and ‘parameters’, where the parameters

will be modeled as first-order Gauss-Markov processes. Unlike the parameters, the states can vary in any manner. This motivates the next section.

**7.3.2. An Improvement for the Case of Parameters** As stated earlier, for a large number of states (say,  $n_{\mathbf{x}} = 35$ ), the  $UDU$  time update for the full covariance matrix (a la Gram-Schmidt orthogonalization) is computationally expensive, requiring  $2n_{\mathbf{x}}^3 + 2n_{\mathbf{x}}^2$  multiplies and additions. This is not competitive with the ‘standard’  $\Phi\mathbf{P}\Phi + \mathbf{Q}$  formulation (which uses  $2n_{\mathbf{x}}^3$  multiplies). However, one might guess that an improvement can be made. This is particularly significant because, normally, most of the states are biases or ECRVs (Exponentially Correlated Random Variables) or first-order Gauss-Markov processes. In order to generalize the development, we assume ECRVs for the ‘parameter’ (or bias) states.

For most space-borne navigation applications, we can usually partition the states into position, velocity, attitude (if applicable) and clock states, all of which we group together and denote as  $\mathbf{x}$ , and parameter states which usually comprise the sensor biases, scale factors, *etc.*, which we denote as  $\mathbf{p}$ . This means that the full state space is

$$\mathcal{X} = \begin{bmatrix} \mathbf{x} \\ \mathbf{p} \end{bmatrix} \quad (7.15)$$

The ‘states’ partition must comprise all those quantities whose time evolution cannot be described as purely self-auto correlated processes. With this in hand, we partition  $\mathbf{U}$  and  $\mathbf{D}$  as

$$\mathbf{U} = \begin{bmatrix} \mathbf{U}_{\mathbf{xx}} & \mathbf{U}_{\mathbf{xp}} \\ \mathbf{0} & \mathbf{U}_{\mathbf{pp}} \end{bmatrix} \quad \mathbf{D} = \begin{bmatrix} \mathbf{D}_{\mathbf{xx}} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_{\mathbf{pp}} \end{bmatrix} \quad (7.16)$$

Also, partition  $\Phi$  according to

$$\Phi = \begin{bmatrix} \Phi_{\mathbf{xx}} & \Phi_{\mathbf{xp}} \\ \mathbf{0} & \mathbf{M} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{M} \end{bmatrix} \begin{bmatrix} \Phi_{\mathbf{xx}} & \Phi_{\mathbf{xp}} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} = \Phi_2 \Phi_1 \quad (7.17)$$

where  $\mathbf{M}$  is a diagonal matrix, representing an ECRV whose propagation for  $p_k$  is

$$\bar{p}_k = e^{-\Delta t/\tau} \bar{p}_{k-1} \quad (7.18)$$

so that

$$\mathbf{M}(i, i) = m_i = e^{-\Delta t/\tau_i} \quad (7.19)$$

where  $\tau_i$  is the time constant of the  $i^{\text{th}}$  ECRV state and  $\mathbf{Q}$  is partitioned according to

$$\mathbf{Q} = \begin{bmatrix} \mathbf{Q}_{\mathbf{xx}} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_{\mathbf{pp}} \end{bmatrix} = \begin{bmatrix} \mathbf{Q}_{\mathbf{xx}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_{\mathbf{pp}} \end{bmatrix} = \mathbf{Q}_1 + \mathbf{Q}_2 \quad (7.20)$$

Recall that the original propagation equation was

$$\overline{\mathbf{U}}\mathbf{D}\overline{\mathbf{U}}^{\text{T}} = \Phi\mathbf{U}\mathbf{D}\mathbf{U}^{\text{T}}\Phi^{\text{T}} + \mathbf{Q} \quad (7.21)$$

Harnessing the development in Appendix B.1,  $\overline{\mathbf{U}}\overline{\mathbf{D}}\overline{\mathbf{U}}^{\text{T}}$  becomes

$$\overline{\mathbf{U}}\overline{\mathbf{D}}\overline{\mathbf{U}}^{\text{T}} = \Phi_2 [\Phi_1 \mathbf{U}\mathbf{D}\mathbf{U}^{\text{T}}\Phi_1^{\text{T}} + \mathbf{Q}_1] \Phi_2^{\text{T}} + \mathbf{Q}_1 \quad (7.22)$$

This suggests the following two-step process:

1) Find  $\tilde{\mathbf{U}}$  and  $\tilde{\mathbf{D}}$  from

$$\tilde{\mathbf{U}}\tilde{\mathbf{D}}\tilde{\mathbf{U}}^{\text{T}} = \Phi_1 \mathbf{U}\mathbf{D}\mathbf{U}^{\text{T}}\Phi_1^{\text{T}} + \mathbf{Q}_1 \quad (7.23)$$

2) Find  $\overline{\mathbf{U}}$  and  $\overline{\mathbf{D}}$  from

$$\overline{\mathbf{U}}\overline{\mathbf{D}}\overline{\mathbf{U}}^{\text{T}} = \Phi_2 \tilde{\mathbf{U}}\tilde{\mathbf{D}}\tilde{\mathbf{U}}^{\text{T}}\Phi_2^{\text{T}} + \mathbf{Q}_2 \quad (7.24)$$

The following sub-sections will describe each of these steps.

7.3.2.1. *The First Sub-Problem* ( $\Phi_1 \mathbf{U} \mathbf{D} \mathbf{U}^\top \Phi_1^\top + \mathbf{Q}_1$ ) Lets look at 1). The left hand side of Eq. (7.23) is

$$\tilde{\mathbf{U}} \tilde{\mathbf{D}} \tilde{\mathbf{U}}^\top = \begin{bmatrix} \tilde{\mathbf{U}}_{xx} \tilde{\mathbf{D}}_{xx} \tilde{\mathbf{U}}_{xx}^\top + \tilde{\mathbf{U}}_{xp} \tilde{\mathbf{D}}_{pp} \tilde{\mathbf{U}}_{xp}^\top & \tilde{\mathbf{U}}_{xp} \tilde{\mathbf{D}}_{pp} \tilde{\mathbf{U}}_{pp}^\top \\ \tilde{\mathbf{U}}_{pp} \tilde{\mathbf{D}}_{pp} \tilde{\mathbf{U}}_{xp}^\top & \tilde{\mathbf{U}}_{pp} \tilde{\mathbf{D}}_{pp} \tilde{\mathbf{U}}_{pp}^\top \end{bmatrix} \quad (7.25)$$

The right hand side of Eq. (7.23) is

$$\tilde{\mathbf{U}} \tilde{\mathbf{D}} \tilde{\mathbf{U}}^\top = \left[ \begin{array}{l|l} \Phi_{xx} (\mathbf{U}_{xx} \mathbf{D}_{xx} \mathbf{U}_{xx}^\top + \mathbf{U}_{xp} \mathbf{D}_{pp} \mathbf{U}_{xp}^\top) \Phi_{xx}^\top & \Phi_{xx} \mathbf{U}_{xp} \mathbf{D}_{pp} \mathbf{U}_{pp}^\top \\ + \Phi_{xx} \mathbf{U}_{xp} \mathbf{D}_{pp} \mathbf{U}_{pp}^\top \Phi_{xp}^\top & + \Phi_{xp} \mathbf{U}_{pp} \mathbf{D}_{pp} \mathbf{U}_{pp}^\top \\ + \Phi_{xp} \mathbf{U}_{pp} \mathbf{D}_{pp} \mathbf{U}_{xp}^\top \Phi_{xx}^\top & \\ + \Phi_{xp} \mathbf{U}_{pp} \mathbf{D}_{pp} \mathbf{U}_{pp}^\top \Phi_{xp}^\top + \mathbf{Q}_{xx} & \\ \hline \mathbf{U}_{pp} \mathbf{D}_{pp} \mathbf{U}_{xp}^\top \Phi_{xx}^\top + \mathbf{U}_{pp} \mathbf{D}_{pp} \mathbf{U}_{pp}^\top \Phi_{xp}^\top & \mathbf{U}_{pp} \mathbf{D}_{pp} \mathbf{U}_{pp}^\top \end{array} \right] \quad (7.26)$$

Equating each component of the matrix in Eqs (7.25) and (7.26), we find that the (2, 2) component yields

$$\tilde{\mathbf{U}}_{pp} = \mathbf{U}_{pp} \quad (7.27)$$

$$\tilde{\mathbf{D}}_{pp} = \mathbf{D}_{pp} \quad (7.28)$$

Equating the (1, 2) (or (2, 1)) component yields

$$\tilde{\mathbf{U}}_{xp} = \Phi_{xx} \mathbf{U}_{xp} + \Phi_{xp} \mathbf{U}_{pp} \quad (7.29)$$

Finally, equating the (1, 1) components of Eqs.(7.25) and (7.26), and using Eqs. (7.27), (7.28) and (7.29), we find that

$$\tilde{\mathbf{U}}_{xx} \tilde{\mathbf{D}}_{xx} \tilde{\mathbf{U}}_{xx}^\top = \Phi_{xx} \mathbf{U}_{xx} \mathbf{D}_{xx} \mathbf{U}_{xx}^\top \Phi_{xx}^\top + \mathbf{Q}_{xx} \quad (7.30)$$

Since we have partitioned the states such that  $\mathbf{x}$  comprises the position, velocity and attitude and is a 9-vector, we use the *modified Gram-Schmidt algorithm* to update  $\tilde{\mathbf{U}}_{xx}$  and  $\tilde{\mathbf{D}}_{xx}$ . And then we compute  $\tilde{\mathbf{U}}_{pp}$ ,  $\tilde{\mathbf{D}}_{pp}$  and  $\tilde{\mathbf{U}}_{xp}$  according to Eqs.(7.27) - (7.29).

Thus, given  $n_x$  states with  $m_x$  process noise parameters associated with those states, the number of computations associated with the the first sub-problem is:  $[1.5n_x^3 + 0.5n_x^2(2m_x - 1)]$  adds,  $[0.5n_x^2(3n_x + 1) + n_x m_x (n_x + 1)]$  multiplies, and  $[(n_x + m_x)(n_x - 1)]$  divides.

7.3.2.2. *The Second Sub-Problem* ( $\Phi_2 \tilde{\mathbf{U}} \tilde{\mathbf{D}} \tilde{\mathbf{U}}^\top \Phi_2^\top + \mathbf{Q}_2$ ) Now we look at 2). We now partition  $\tilde{\mathbf{U}}$  and  $\tilde{\mathbf{D}}$  as

$$\tilde{\mathbf{U}} = \begin{bmatrix} \tilde{\mathbf{U}}_{aa} & \tilde{\mathbf{U}}_{ab} & \tilde{\mathbf{U}}_{ac} \\ \mathbf{0} & 1 & \tilde{\mathbf{U}}_{bc} \\ \mathbf{0} & \mathbf{0} & \tilde{\mathbf{U}}_{cc} \end{bmatrix} \begin{array}{l} \} n_a \\ \} 1 \\ \} n_c \end{array} \quad \text{and} \quad \tilde{\mathbf{D}} = \begin{bmatrix} \tilde{\mathbf{D}}_{aa} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \tilde{d}_b & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \tilde{\mathbf{D}}_{cc} \end{bmatrix} \begin{array}{l} \} n_a \\ \} 1 \\ \} n_c \end{array} \quad (7.31)$$

in order to isolate a parameter. Correspondingly,

$$\Phi_2 = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{M}_c \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & m_b & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I} \end{bmatrix} = \Phi_c \Phi_b \quad (7.32)$$

and

$$\mathbf{Q}_2 = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & q_b & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{Q}_c \end{bmatrix} = \mathbf{Q}_b + \mathbf{Q}_c \quad (7.33)$$

As in the previous section, we note that  $\Phi_c^{-1} \mathbf{Q}_b \Phi_c^{-\top} = \mathbf{Q}_b$ . So, now Eq. (7.24) becomes

$$\overline{\mathbf{U}} \overline{\mathbf{D}} \overline{\mathbf{U}}^\top = \Phi_c \left[ \Phi_b \tilde{\mathbf{U}} \tilde{\mathbf{D}} \tilde{\mathbf{U}}^\top \Phi_b^\top + \mathbf{Q}_b \right] \Phi_c^\top + \mathbf{Q}_c \quad (7.34)$$

The term in the square bracket in Eq. (7.34) is

$$\check{\mathbf{U}} \check{\mathbf{D}} \check{\mathbf{U}}^\top = \Phi_b \tilde{\mathbf{U}} \tilde{\mathbf{D}} \tilde{\mathbf{U}}^\top \Phi_b^\top + \mathbf{Q}_b \quad (7.35)$$

In Appendix B the above equation is expanded until we find that

$$\check{\mathbf{U}}_{\mathbf{ac}} = \tilde{\mathbf{U}}_{\mathbf{ac}}, \quad \check{\mathbf{D}}_{\mathbf{cc}} = \tilde{\mathbf{D}}_{\mathbf{cc}}, \quad \check{\mathbf{U}}_{\mathbf{cc}}^\top = \tilde{\mathbf{U}}_{\mathbf{cc}}^\top \quad (7.36)$$

Additionally,

$$\check{\mathbf{U}}_{\mathbf{bc}} = m_b \tilde{\mathbf{U}}_{\mathbf{bc}} \quad (7.37)$$

and

$$\check{d}_b = m_b^2 \tilde{d}_b + q_b \quad (7.38)$$

and

$$\check{\mathbf{U}}_{\mathbf{ab}} = m_b \frac{\tilde{d}_b}{\check{d}_b} \tilde{\mathbf{U}}_{\mathbf{ab}} \quad (7.39)$$

Finally, we find that

$$\check{\mathbf{U}}_{\mathbf{aa}} \check{\mathbf{D}}_{\mathbf{aa}} \check{\mathbf{U}}_{\mathbf{aa}}^\top = \tilde{\mathbf{U}}_{\mathbf{aa}} \tilde{\mathbf{D}}_{\mathbf{aa}} \tilde{\mathbf{U}}_{\mathbf{aa}}^\top + \left( \frac{\tilde{d}_b q_b}{\check{d}_b} \right) \tilde{\mathbf{U}}_{\mathbf{ab}} \tilde{\mathbf{U}}_{\mathbf{ab}}^\top \quad (7.40)$$

We note that  $\tilde{\mathbf{U}}_{\mathbf{ab}}$  is a column vector so Eq.(7.40), and hence is of rank 1, constitutes a ‘rank one’ update. Since  $\check{d}_b$ ,  $\tilde{d}_b$  and  $q_b$  are all positive (assuming  $m_b$  is a positive quantity), we can use the Agee-Turner Rank One update [1]. It should be pointed out that as the algorithm proceeds down the ‘list’ of parameters, the size of the states  $\mathbf{a}$  increases by one (and consequently the size of the parameters  $\mathbf{c}$  decreases by one. Hence  $\check{\mathbf{U}}_{\mathbf{aa}}$  and  $\check{\mathbf{D}}_{\mathbf{aa}}$  begins with a dimension of  $n_{\mathbf{x}}$  and concludes with dimension  $n_{\mathbf{x}} + n_{\mathbf{p}} - 1$ .

Therefore, this is done recursively for all the (sensor) parameters  $\mathbf{p}$  which are of size  $n_{\mathbf{p}}$ .

### The Algorithm for $\Phi_2 \tilde{\mathbf{U}} \tilde{\mathbf{D}} \tilde{\mathbf{U}}^\top \Phi_2^\top + \mathbf{Q}_2$

The algorithm can be expressed as follows (with the arithmetic operations (adds, multiplies, divides) in square brackets per  $k$ ):

```

for  $k = 1, \dots, n_{\mathbf{p}}$  do
   $\check{\mathbf{D}}(n_{\mathbf{x}} + k, n_{\mathbf{x}} + k) = \mathbf{M}(k, k)^2 \tilde{\mathbf{D}}(n_{\mathbf{x}} + k, n_{\mathbf{x}} + k) + \mathbf{Q}_{\mathbf{pp}}(k, k)$  (Eq. (7.38)) [1, 2, 0]
   $\alpha = \mathbf{M}(k, k) \frac{\tilde{\mathbf{D}}(n_{\mathbf{x}} + k, n_{\mathbf{x}} + k)}{\tilde{\mathbf{D}}(n_{\mathbf{x}} + k, n_{\mathbf{x}} + k)}$  [0, 1, 1]
  for  $i = 1, \dots, (n_{\mathbf{x}} + k - 1)$  do
     $\check{\mathbf{U}}(i, n_{\mathbf{x}} + k) = \alpha \tilde{\mathbf{U}}(i, n_{\mathbf{x}} + k)$  (Eq. (7.37)) [0,  $n_{\mathbf{x}} + k - 1$ , 0]
  end for

```

**for**  $j = n_{\mathbf{x}} + k + 1, \dots, (n_{\mathbf{x}} + n_{\mathbf{p}})$  **do**  
 $\check{\mathbf{U}}(n_{\mathbf{x}} + k, j) = \mathbf{M}(k, k)\check{\mathbf{U}}(n_{\mathbf{x}} + k, j)$  (Eq. (7.39))  $[0, n_{\mathbf{p}} - k, 0]$   
**end for**  
 Solve for  $\check{\mathbf{U}}_{\mathbf{xx}}^{(k)}, \check{\mathbf{D}}_{\mathbf{xx}}^{(k)}$  using<sup>1</sup> the Rank-One update  $[(n_{\mathbf{x}} + k)^2, (n_{\mathbf{x}} + k)^2 + 3(n_{\mathbf{x}} + k) + 2, 0]$   
**end for**

Thus, the arithmetic operations are as follows:

Adds:

$$\sum_{k=1}^{n_{\mathbf{p}}} ((n_{\mathbf{x}} + k)^2 + 1) = n_{\mathbf{x}}^2 n_{\mathbf{p}} + n_{\mathbf{p}} + n_{\mathbf{x}}(n_{\mathbf{p}} + 1)n_{\mathbf{p}} + \sum_{k=1}^{n_{\mathbf{p}}} k^2 \quad (7.41)$$

Multiplies:

$$\begin{aligned} \sum_{k=1}^{n_{\mathbf{p}}} (3 + (n_{\mathbf{x}} + k - 1) + n_{\mathbf{p}} - k + (n_{\mathbf{x}} + k)^2 + 3(n_{\mathbf{x}} + k) + 2) &= (5.5 + 5n_{\mathbf{x}} + n_{\mathbf{x}}^2)n_{\mathbf{p}} \\ &+ \frac{1}{2}(2n_{\mathbf{x}} + 5)n_{\mathbf{p}}^2 + \sum_{k=1}^{n_{\mathbf{p}}} k^2 \quad (7.42) \end{aligned}$$

and  $n_{\mathbf{p}}$  divides.

### The Agee-Turner Rank One Update Algorithm

Appendix B.3 contains the development of the Agee-Turner Rank-One update which is the key to reducing the numerical operations on the UDU Time update. Given  $\mathbf{U}$  and  $\mathbf{D}$ , along with  $c$ , and the vector  $\mathbf{x}$ , we are interested in obtaining  $\tilde{\mathbf{U}}$  and  $\tilde{\mathbf{D}}$  along the lines of

$$\tilde{\mathbf{U}}\tilde{\mathbf{U}}^{\mathbf{T}} = \mathbf{U}\mathbf{D}\mathbf{U}^{\mathbf{T}} + c\mathbf{x}\mathbf{x}^{\mathbf{T}} \quad (7.43)$$

The algorithm to compute  $\tilde{U}_{ij}$  and  $\tilde{D}_{ii}$  ss:

$\mathcal{C}^n = c$   
**for**  $j = n, \dots, 2$  **do**  
 $\tilde{D}_{jj} = D_{jj} + \mathcal{C}^j x_j^2$   $[n - 1, 2(n - 1), 0]$   
 $\tilde{U}_{jj} = 1$   
 $\beta_j = \mathcal{C}^j / \tilde{D}_{jj}$   $[0, 0, (n - 1)]$   
 $v_j = \beta_j x_j$   $[0, (n - 1), 0]$   
**for**  $i = 1, \dots, j - 1$  **do**  
 $x_i := x_i - U_{ij} x_j$   $[\frac{1}{2}(n^2 - n), \frac{1}{2}(n^2 - n), 0]$   
 $\tilde{U}_{ij} = U_{ij} + x_i v_j$   $[\frac{1}{2}(n^2 - n), \frac{1}{2}(n^2 - n), 0]$   
**end for**  
 $\mathcal{C}^{j-1} = \beta_j D_{jj}$   $[0, (n - 1), 0]$   
**end for**  
 $\tilde{D}_{11} = D_{11} + \mathcal{C}^1 x_1^2$   $[1, 2, 0]$

This algorithm has  $n^2$  adds,  $(n^2 + 3n + 2)$  multiplies and  $n - 1$  divides.

<sup>1</sup>We are using the nomenclature  $\mathbf{U}^{(k)}$  and  $\mathbf{D}^{(k)}$  to denote the upper left  $n_{\mathbf{x}} + k - 1$  rows and columns of the  $\mathbf{U}$  and  $\mathbf{D}$  matrices

**7.3.3. Arithmetic Operations for Time Update** For the time update of the covariance matrix, we will have (from Section XX and XX), we have the following arithmetic operations:

$$\text{Adds : } 1.5n_{\mathbf{x}}^3 + n_{\mathbf{x}}^2 m_{\mathbf{x}} + n_{\mathbf{x}}^2 n_{\mathbf{p}} - 0.5n_{\mathbf{x}}^2 + n_{\mathbf{p}} + n_{\mathbf{x}}(n_{\mathbf{p}} + 1)n_{\mathbf{p}} + \sum_{k=1}^{n_{\mathbf{p}}} k^2$$

$$\text{Multiplies : } 0.5n_{\mathbf{x}}^2(3n_{\mathbf{x}} + 1) + n_{\mathbf{x}}m_{\mathbf{x}}(n_{\mathbf{x}} + 1) + (5.5 + 5n_{\mathbf{x}} + n_{\mathbf{x}}^2)n_{\mathbf{p}} + \frac{1}{2}(2n_{\mathbf{x}} + 5)n_{\mathbf{p}}^2 + \sum_{k=1}^{n_{\mathbf{p}}} k^2$$

$$\text{Divides : } (n_{\mathbf{x}} + m_{\mathbf{x}})(n_{\mathbf{x}} - 1) + n_{\mathbf{p}}$$

For  $n_{\mathbf{x}} = 9$ ,  $m_{\mathbf{x}} = 9$ ,  $n_{\mathbf{p}} = 26$ , we will utilize 16,407 adds, 19,338 multiplies, and 170 divides. In contrast, if we did the MGS on all 35 states ( $n_{\mathbf{x}} = 35$ ,  $m_{\mathbf{x}} = 35$  and  $n_{\mathbf{p}} = 0$ ), we would use 85,750 adds, 88,200 multiplies, and 34 divides. Finally, if the covariance were updated (without any consideration given to the structure of  $\Phi$  from  $\Phi\bar{\mathbf{P}}\Phi^T$ ) in the conventional manner, with  $n_{\mathbf{x}} = 35$ ,  $m_{\mathbf{x}} = 35$ , it would cost 84,525 adds, 85,750 multiplies and no divides<sup>†</sup>. Thus, a very strong case is made for using the *UDU* factorization and harnessing the benefit of updating the sensor parameters using the Agee-Turner Rank-One update. Thus, the *UDU* time update taking advantage of the fact that most of the states are sensor parameters results in nearly **five times fewer** adds and multiplies and 170 more divides that if we operated on the full covariance matrix.

#### 7.4. The Measurement Update

The *UDU* factorization requires that we process the measurements one-at-a-time [2]. This should not be construed as a weakness of the formulation. If the measurements are correlated, they can be ‘decorrelated’ as in Appendix B.4

So, the covariance update equations are

$$\mathbf{P} = \bar{\mathbf{P}} - \mathbf{K}\mathbf{H}\bar{\mathbf{P}} \quad (7.45)$$

where  $\mathbf{K}$  is the Kalman Gain matrix,  $\mathbf{H}$  is the measurement partial,  $\bar{\mathbf{P}}$  is the *a priori* covariance, and  $\mathbf{P}$  is the *a posteriori* covariance matrix. Using the covariance factors  $\mathbf{U}$  and  $\mathbf{D}$ , we rewrite the above equation as

$$\mathbf{U}\mathbf{D}\mathbf{U}^T = \bar{\mathbf{U}}\bar{\mathbf{D}}\bar{\mathbf{U}}^T - \mathbf{K}\mathbf{H}\bar{\mathbf{U}}\bar{\mathbf{D}}\bar{\mathbf{U}}^T \quad (7.46)$$

We note that  $\mathbf{U}$  and  $\bar{\mathbf{U}}$  and  $\mathbf{D}$  and  $\bar{\mathbf{D}}$  are  $n \times n$  matrices, and because we are using the paradigm of processing the measurements one at a time,  $\mathbf{H}$  is an  $1 \times n$  vector and  $\mathbf{K}$  is an  $n \times 1$  vector. Recalling that  $\mathbf{K}$  is defined as

$$\mathbf{K} = \bar{\mathbf{P}}\mathbf{H}^T (\mathbf{H}\bar{\mathbf{P}}\mathbf{H}^T + \mathbf{R})^{-1} = \frac{\bar{\mathbf{U}}\bar{\mathbf{D}}\bar{\mathbf{U}}^T\mathbf{H}^T}{\mathbf{H}\bar{\mathbf{U}}\bar{\mathbf{D}}\bar{\mathbf{U}}^T\mathbf{H}^T + \mathbf{R}} = \frac{\bar{\mathbf{U}}\bar{\mathbf{D}}\bar{\mathbf{U}}^T\mathbf{H}^T}{\alpha} \quad (7.47)$$

where the scalar  $\alpha$  is defined to be

$$\alpha \triangleq \mathbf{H}\bar{\mathbf{U}}\bar{\mathbf{D}}\bar{\mathbf{U}}^T\mathbf{H}^T + \mathbf{R} \quad (7.48)$$

<sup>†</sup>For matrices,  $\mathbf{A}$  and  $\mathbf{B}$  of dimension  $n \times m$  and  $m \times p$ , respectively, the product

$$\mathbf{C} = \mathbf{A}\mathbf{B} \quad (7.44)$$

results in  $n(m-1)p$  adds,  $nmp$  multiplies and no divides.

We find that Eq. (7.46) becomes

$$\mathbf{U}\mathbf{D}\mathbf{U}^\top = \bar{\mathbf{U}} \left[ \bar{\mathbf{D}} - \frac{\bar{\mathbf{D}}\bar{\mathbf{U}}^\top\mathbf{H}^\top}{\alpha}\mathbf{H}\bar{\mathbf{U}}\bar{\mathbf{D}} \right] \bar{\mathbf{U}}^\top \quad (7.49)$$

If we define the  $n \times 1$  vector  $\mathbf{v}$  as

$$\mathbf{v} \triangleq \bar{\mathbf{D}}\bar{\mathbf{U}}^\top\mathbf{H}^\top \quad (7.50)$$

Eq. (7.49) becomes

$$\mathbf{U}\mathbf{D}\mathbf{U}^\top = \bar{\mathbf{U}} \left[ \bar{\mathbf{D}} - \frac{1}{\alpha}\mathbf{v}\mathbf{v}^\top \right] \bar{\mathbf{U}}^\top \quad (7.51)$$

We now analyze the bracketed term in Eq. (7.51) and find that we can define

$$\tilde{\mathbf{U}}\tilde{\mathbf{D}}\tilde{\mathbf{U}}^\top \triangleq \bar{\mathbf{D}} - \frac{1}{\alpha}\mathbf{v}\mathbf{v}^\top \quad (7.52)$$

Therefore,

$$\mathbf{U} = \bar{\mathbf{U}}\tilde{\mathbf{U}} \quad \text{and} \quad \mathbf{D} = \tilde{\mathbf{D}} \quad (7.53)$$

So, how do we proceed? This has all the marks of a rank-one update, for after all  $\mathbf{v}$  is of rank one. We *can* proceed by using the Agee-Turner rank-one update. Except for one thing – that pesky minus sign in Eq. (7.52). That minus sign portends all sorts of numerical issues because there is a strong possibility that we can lose numerical precision if the Agee-Turner update is used blindly. It turns out that we can have ‘our cake and eat it too’, for Neil Carlson developed a rank-one update to remedy precisely our issue. The mathematical development of this algorithm is detailed in Appendix B.5.

The algorithm is as follows: Given  $\bar{\mathbf{U}}$ ,  $\bar{\mathbf{D}}$ ,  $\mathbf{R}$ , and  $\mathbf{H}$

$$\begin{aligned} \mathbf{f} &= \bar{\mathbf{U}}^\top\mathbf{H}^\top \quad \text{where} \quad \mathbf{f} = [f_1 \ f_2 \ \cdots \ f_n]^\top && [\frac{1}{2}(n_{\mathbf{x}}^2 - n_{\mathbf{x}}), \frac{1}{2}(n_{\mathbf{x}}^2 - n_{\mathbf{x}}), 0] \\ \mathbf{v} &= \bar{\mathbf{D}}^\top\mathbf{f} \quad \text{where} \quad \mathbf{v} = [\bar{d}_1f_1 \ \bar{d}_2f_2 \ \cdots \ \bar{d}_nf_n]^\top && [0, n_{\mathbf{x}}, 0] \\ \bar{\mathbf{K}}_1 &= [v_1 \ 0 \ \cdots \ 0]^\top && \\ \alpha_1 &= \mathbf{R} + v_1f_1 && [1, 1, 0] \\ d_1 &= \left(\frac{\mathbf{R}}{\alpha_1}\right)\bar{d}_1 && [0, 1, 1] \\ \mathbf{for} \ j &= 2, \dots, n \ \mathbf{do} && \\ \alpha_j &= \alpha_{j-1} + v_jf_j && [n_{\mathbf{x}} - 1, n_{\mathbf{x}} - 1, 0] \\ d_j &= \left(\frac{\alpha_{j-1}}{\alpha_j}\right)\bar{d}_j && [0, n_{\mathbf{x}} - 1, n_{\mathbf{x}} - 1] \\ \lambda_j &= -(f_j/\alpha_{j-1}) && [0, 0, n_{\mathbf{x}} - 1] \\ \mathbf{U}_j &= \bar{\mathbf{U}}_j + \lambda_j\bar{\mathbf{K}}_{j-1} && [\frac{1}{2}(n_{\mathbf{x}}^2 - n_{\mathbf{x}}), \frac{1}{2}(n_{\mathbf{x}}^2 - n_{\mathbf{x}}), 0] \\ \bar{\mathbf{K}}_j &= \bar{\mathbf{K}}_{j-1} + v_j\bar{\mathbf{U}}_j && [\frac{1}{2}(n_{\mathbf{x}}^2 - n_{\mathbf{x}}), \frac{1}{2}(n_{\mathbf{x}}^2 - n_{\mathbf{x}}), 0] \\ \mathbf{end} \ \mathbf{for} &&& \\ \mathbf{K} &= \bar{\mathbf{K}}_n/\alpha && [0, 0, n_{\mathbf{x}}] \end{aligned}$$

Thus, taking advantage of the triangularity of the  $\bar{\mathbf{U}}$  matrix (and the fact that  $\bar{\mathbf{U}}^\top\mathbf{H}^\top$  and  $\lambda_j\bar{\mathbf{K}}_{j-1}$  and  $v_j\bar{\mathbf{U}}_j$  use  $n_{\mathbf{x}}(n_{\mathbf{x}} - 1)/2$  multiplies and adds), for each measurement processed, the covariance update results in  $1.5n_{\mathbf{x}}^2 - 0.5n_{\mathbf{x}}$  adds,  $1.5n_{\mathbf{x}}^2 + 1.5n_{\mathbf{x}}$  multiplies and  $3n_{\mathbf{x}} - 1$  divides.

For the normal, Joseph Kalman filter update, for a scalar measurement, we find that if we use efficient methods of calculating and storing quantities [2], we use  $4.5n_x^2 + 3.5n_x$  adds,  $4n_x^2 + 4.5n_x$  multiplies and 1 divide.

For the ‘‘Conventional’’ Kalman filter update ( $\mathbf{P} = \bar{\mathbf{P}} - \mathbf{KHP}$  in Eq.(61)), for a scalar measurement, we find that [2] we use  $1.5n_x^2 + 1.5n_x$  adds,  $1.5n_x^2 + 0.5n_x$  multiplies and 1 divide.

Thus, for  $n_x = 35$ , the covariance update due to measurement processing with the  $UDU$  factorization uses 1820 adds, 1890 divides and 104 divides compared with 5635 adds, 5058 multiplies and 1 divide for the efficient Joseph update. The ‘‘Conventional’’ Kalman update uses 1890 adds, 1855 multiplies, and 1 divide.

Hence there almost a factor of 2.5 improvement in the adds and multiplies using the triangular ( $UDU$ ) update compared with the Joseph update. This rivals the efficiency of the ‘‘conventional’’ Kalman Filter update.

### 7.5. Consider Covariance and Its Implementation in the UDU Filter

‘Consider’ Analysis was first introduced by S. F. Schmidt of NASA Ames in the mid 1960s as a means to account for errors in both the dynamic and measurement models due to uncertain parameters [64]. The Consider Kalman Filter, also called the Schmidt-Kalman Filter, resulted from this body of work. The consider approach is especially useful when parameters have low observability.

We partition the state-vector,  $\mathbf{x}$ , into the  $n_s$  ‘‘estimated states’’,  $\mathbf{s}$ , and the  $n_p$  ‘‘consider’’ parameters,  $\mathbf{p}$ , as

$$\mathbf{x} \triangleq \begin{bmatrix} \mathbf{s} \\ \mathbf{p} \end{bmatrix} \quad (7.54)$$

so that

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_{ss} & \mathbf{P}_{sp} \\ \mathbf{P}_{ps} & \mathbf{P}_{pp} \end{bmatrix}, \mathbf{H} = \begin{bmatrix} \mathbf{H}_s & \mathbf{H}_p \end{bmatrix}, \mathbf{K}_{opt} = \begin{bmatrix} \mathbf{K}_{s,opt} \\ \mathbf{K}_{p,opt} \end{bmatrix} = \begin{bmatrix} \mathbf{P}_{ss}^- \mathbf{H}_s^T + \mathbf{P}_{sp}^- \mathbf{H}_p^T \\ \mathbf{P}_{ps}^- \mathbf{H}_s^T + \mathbf{P}_{pp}^- \mathbf{H}_p^T \end{bmatrix} \mathbf{W}^{-1}$$

where  $\mathbf{K}_{opt}$  is the optimal Kalman gain computed for the full state,  $\mathbf{x}$ . Therefore, if we now choose the  $\mathbf{K}_s$  and  $\mathbf{K}_p$  carefully such that the  $\mathbf{K}_s = \mathbf{K}_{s,opt}$ , the *a posteriori* covariance matrix is

$$\mathbf{P}^+ = \begin{bmatrix} \mathbf{P}_{ss}^- - \mathbf{K}_s \mathbf{W} \mathbf{K}_s^T & \mathbf{P}_{sp}^- - \mathbf{K}_s \mathbf{H} \begin{bmatrix} \mathbf{P}_{sp}^- \\ \mathbf{P}_{pp}^- \end{bmatrix} \\ \mathbf{P}_{ps}^- - \begin{bmatrix} \mathbf{P}_{sp}^- \\ \mathbf{P}_{pp}^- \end{bmatrix}^T \mathbf{H}^T \mathbf{K}_s^T & \mathbf{P}_{pp}^- - \mathbf{K}_p \mathbf{W} \mathbf{K}_p^T \end{bmatrix} \quad (7.55)$$

This equation is valid for any value of  $\mathbf{K}_p$ . Notice that there is no  $\mathbf{K}_p$  in the correlation terms of the covariance matrix. **Therefore, what is remarkable about this equation is that once the optimal  $\mathbf{K}_s$  is chosen, the correlation between  $\mathbf{s}$  and  $\mathbf{p}$  is independent of the choice of  $\mathbf{K}_p$ .**

In its essence, the consider parameters are not updated; therefore, the Kalman gain associated with the consider parameters,  $\mathbf{p}$ , is zero, *i.e.*  $\mathbf{K}_p = \mathbf{0}$ . However, several comments are in order:

- (1) When using the Schmidt-Kalman filter, the *a priori* and *a posteriori* covariance of the parameters ( $\mathbf{P}_{pp}$ ) are the same.
- (2) The *a posteriori* covariance matrix of the states and the correlation between the states and the parameters are the same regardless of whether one uses the Schmidt-Kalman filter or the optimal Kalman update

Therefore, the *consider* covariance,  $\mathbf{P}_{con}^+$  is

$$\mathbf{P}_{con}^+ = \begin{bmatrix} \mathbf{P}_{ss}^- - \mathbf{K}_s \mathbf{W} \mathbf{K}_s^\top & \mathbf{P}_{sp}^- - \mathbf{K}_s \mathbf{H} \begin{bmatrix} \mathbf{P}_{sp}^- \\ \mathbf{P}_{pp}^- \end{bmatrix} \\ \mathbf{P}_{ps}^- - \begin{bmatrix} \mathbf{P}_{sp}^- \\ \mathbf{P}_{pp}^- \end{bmatrix}^\top \mathbf{H}^\top \mathbf{K}_s^\top & \mathbf{P}_{pp}^- \end{bmatrix} \quad (7.56)$$

Of course, the “full” optimal covariance matrix update is

$$\mathbf{P}_{opt}^+ = \begin{bmatrix} \mathbf{P}_{ss}^- - \mathbf{K}_{s,opt} \mathbf{W} \mathbf{K}_{s,opt}^\top & \mathbf{P}_{sp}^- - \mathbf{K}_{s,opt} \mathbf{H} \begin{bmatrix} \mathbf{P}_{sp}^- \\ \mathbf{P}_{pp}^- \end{bmatrix} \\ \mathbf{P}_{ps}^- - \begin{bmatrix} \mathbf{P}_{sp}^- \\ \mathbf{P}_{pp}^- \end{bmatrix}^\top \mathbf{H}^\top \mathbf{K}_{s,opt}^\top & \mathbf{P}_{pp}^- - \mathbf{K}_{p,opt} \mathbf{W} \mathbf{K}_{p,opt}^\top \end{bmatrix} \quad (7.57)$$

The UDU formulation, while numerically stable and tight, is quite inflexible to making any changes in the framework. At first blush, it would seem that the consider analysis would not fit into the framework. However, all is not in vain. With some clever rearrangements, we can allow for a rank-one update to include consider states in the measurement update. The measurement update, expressed in terms of the consider covariance [79], is

$$\mathbf{P}_{opt}^+ = \mathbf{P}_{con}^+ - W (\mathbf{S} \mathbf{K}_{opt}) (\mathbf{S} \mathbf{K}_{opt})^\top \quad (7.58)$$

where  $\mathbf{S}$  is an  $n_x \times n_x$  matrix (defining  $n_x \triangleq n_s + n_p$ , where  $n_x$  is the total number of states,  $n_p$  is the number of consider states, and  $n_s$  is the number of “non-consider” states) defined as

$$\mathbf{S} \triangleq \begin{bmatrix} \mathbf{0}_{n_s \times n_s} & \mathbf{0}_{n_s \times n_p} \\ \mathbf{0}_{n_p \times n_s} & \mathbf{I}_{n_p \times n_p} \end{bmatrix} \quad (7.59)$$

Since we are processing scalar measurements, we note that  $W = \frac{1}{\alpha}$  is a scalar and  $\mathbf{K}_{opt}$  is an  $n_x \times 1$  vector. Therefore  $\mathbf{S} \mathbf{K}_{opt}$  is an  $n_x \times 1$  vector. Therefore, solving for the consider covariance,

$$\mathbf{P}_{con}^+ = \mathbf{P}_{opt}^+ + W (\mathbf{S} \mathbf{K}_{opt}) (\mathbf{S} \mathbf{K}_{opt})^\top \quad (7.60)$$

Eq. (7.58) has the same form as the original rank-one update *i.e.*  $\mathbf{P}^+ = \mathbf{P}^- + \mathbf{c} \mathbf{a} \mathbf{a}^\top$ . With this in mind, we can use the (un-modified) rank-one update which is a backward-recursive update [1]. If, for example, all the consider parameters are in the lower part of the state-space, we can effectively reduce the computations by ending the update when the covariance of the state of the last consider parameter is updated.

Therefore, the procedure is as follows: first perform a complete rank-one measurement update with the optimal Kalman Gain ( $\mathbf{K}_{opt}$ ) according to the *modified* rank-one update – on the full covariance matrix. Second, perform another rank-one update with  $\mathbf{a} = \mathbf{S} \mathbf{K}_{opt}$  and  $c = W$ , according to the (un-modified) rank-one update.

Therefore, since there is an additional rank-one update associated with the consider states and if no rearrangement of the consider states are performed, then there will be an additional  $n_x^2$  adds, and  $n_x^2 + 3n_x + 2$  multiplies, and  $n_x - 1$  divides per measurement.

The use of the ‘consider state’ option, if it is exercised, is likely to be used in ‘consider’ing the attitude states, particularly during entry. The rationale for this is that in certain degenerate cases, when GPS satellites are reacquired after entry blackout, the attitude could be adversely affected. So, to protect for this, the ‘consider’ option may be exercised with respect to the attitude states.

### 7.6. Conclusions

Matrix factorization methods, particularly the UDU factorization, are very useful – indeed essential – for onboard navigation algorithms. They are numerically stable and computationally efficient, competitive with the classic Kalman filter implementation. In addition, they allow the navigation designer to investigate the positive definiteness of the covariance matrix for ‘free’, via the entries of the diagonal matrix  $\mathbf{D}$ .

## CHAPTER 8

# Attitude Estimation

Contributed by F. Landis Markley

The particular complications of attitude estimation arise from a fundamental difference between rotational kinematics and translational kinematics. The translational state of motion can be completely specified in a nonsingular way by the cartesian components of the position vector  $\mathbf{r}(t)$  and the velocity vector  $\mathbf{v}(t)$ . The integral of any reasonable function  $\mathbf{v}(t)$  between two times gives the translational displacement of  $\mathbf{r}(t)$  between these two times. Other parameterizations of the translational state may be singular; the classical Keplerian orbit elements are singular for zero inclination or zero eccentricity, for example. It is the case, however, that globally nonsingular six-parameter representations exist.

Rotations in three-dimensional space have three degrees of freedom, just like translations in three dimensions, and the angular velocity vector  $\boldsymbol{\omega}(t)$  is the rotational analog of the velocity vector. However, two different time histories  $\boldsymbol{\omega}(t)$  that have the same integral over a time interval can result in different rotational displacements over the interval. This is because the order in which rotations are performed is significant, unlike the order in which translations are performed. Thus integration of  $\boldsymbol{\omega}(t)$  does not result in a three-vector rotational analog of the position vector. In fact, it can be proven that no global three-component parameterization of rotations without singular points exists [68]. Rotational analysis is forced to deal with either higher-dimensional representations of rotations or with three-dimensional representations possessing singularities or discontinuities. The following seven sections will briefly present a nine-parameter representation, two four-parameter representations, and five three-parameter representations. Fuller discussions can be found in Refs. [49, 67]. The discussion of attitude representations is followed by two sections on extended Kalman filters for attitude estimation.

### 8.1. Attitude Matrix Representation

Attitude representations are the methods of representing the orientation of an orthonormal triad of basis vectors in one reference frame with respect to an orthonormal triad in some other reference frame. The attitude matrix, in particular, represents the orientation of a vehicle's body frame with respect to a frame that is often, but not always, an inertial frame. The attitude determination of earth-pointing spacecraft, for example, typically employs a reference frame in which one basis vector is pointed from the spacecraft toward the center of the earth and another points opposite to the orbital angular velocity. The body frame of a rigid vehicle is simply defined as a frame fixed in the vehicle. No vehicle is completely rigid, though, so it is quite common to define the body frame operationally as the orientation of some *navigation base*, a sufficiently rigid subsystem of the spacecraft including the most critical attitude sensors and payload instruments.

For actual applications, representations are  $3 \times 3$  matrices that transform the representations of vectors in one frame, i.e. their components along the basis vectors in that frame, to their representations in a different frame. Thus attitude representations describe a fixed physical vector in a rotated frame rather than a rotated vector. This is the *passive* interpretation of a transformation, also known as the *alias* sense (from the Latin word for “otherwise,” in the sense of “otherwise known as”) [67]. The alternative *active* interpretation (also known as the *alibi* sense from the Latin word for “elsewhere”) considers the representation in a fixed reference frame of a rotated physical vector. It is crucial to keep this distinction in mind, because an active rotation in one direction corresponds to a passive rotation in the opposite direction. Overlooking this point has led to errors in flight software.<sup>1</sup>

Now consider transforming the representation of a vector  $\vec{x}$  in a frame  $\mathcal{F}$  to its representation in a frame  $\mathcal{G}$  and then from frame  $\mathcal{G}$  to frame  $\mathcal{H}$  or directly from frame  $\mathcal{F}$  to frame  $\mathcal{H}$ , so

$$\mathbf{x}_{\mathcal{H}} = \mathbf{A}_{\mathcal{H}\mathcal{G}}\mathbf{x}_{\mathcal{G}} = \mathbf{A}_{\mathcal{H}\mathcal{G}}(\mathbf{A}_{\mathcal{G}\mathcal{F}}\mathbf{x}_{\mathcal{F}}) = \mathbf{A}_{\mathcal{H}\mathcal{F}}\mathbf{x}_{\mathcal{F}} \quad (8.1)$$

These transformations must be equivalent for any vector  $\mathbf{x}_{\mathcal{F}}$ , so successive transformations are accomplished by simple matrix multiplication:

$$\mathbf{A}_{\mathcal{H}\mathcal{F}} = \mathbf{A}_{\mathcal{H}\mathcal{G}}\mathbf{A}_{\mathcal{G}\mathcal{F}} \quad (8.2)$$

This may appear to be an obvious result, but only one other attitude representation has such a simple composition rule. Matrix multiplication is associative, meaning that  $\mathbf{A}_{\mathcal{H}\mathcal{G}}(\mathbf{A}_{\mathcal{G}\mathcal{F}}\mathbf{A}_{\mathcal{F}\mathcal{E}}) = (\mathbf{A}_{\mathcal{H}\mathcal{G}}\mathbf{A}_{\mathcal{G}\mathcal{F}})\mathbf{A}_{\mathcal{F}\mathcal{E}}$ . Matrix multiplication is not commutative, however, which means that  $\mathbf{A}_{\mathcal{H}\mathcal{G}}\mathbf{A}_{\mathcal{G}\mathcal{F}} \neq \mathbf{A}_{\mathcal{G}\mathcal{F}}\mathbf{A}_{\mathcal{H}\mathcal{G}}$  in general. The non-commutativity of matrix multiplication is at the heart of the problem of finding a suitable attitude representation.

Transforming from frame  $\mathcal{F}$  to frame  $\mathcal{G}$  and back to frame  $\mathcal{F}$  is effected by the matrix  $\mathbf{A}_{\mathcal{F}\mathcal{F}} = \mathbf{A}_{\mathcal{F}\mathcal{G}}\mathbf{A}_{\mathcal{G}\mathcal{F}}$ , which must be the identity matrix. Rotations must also preserve inner products and norms of vectors, so

$$\mathbf{x}_{\mathcal{F}} \cdot \mathbf{y}_{\mathcal{F}} = \mathbf{x}_{\mathcal{G}} \cdot \mathbf{y}_{\mathcal{G}} = (\mathbf{A}_{\mathcal{G}\mathcal{F}}\mathbf{x}_{\mathcal{F}})^{\top} \mathbf{A}_{\mathcal{G}\mathcal{F}}\mathbf{y}_{\mathcal{F}} = \mathbf{x}_{\mathcal{F}}^{\top} \mathbf{A}_{\mathcal{G}\mathcal{F}}^{\top} \mathbf{A}_{\mathcal{G}\mathcal{F}}\mathbf{y}_{\mathcal{F}} \quad (8.3)$$

These two observations mean that

$$\mathbf{A}_{\mathcal{G}\mathcal{F}}^{\top} = \mathbf{A}_{\mathcal{G}\mathcal{F}}^{-1} = \mathbf{A}_{\mathcal{F}\mathcal{G}} \quad (8.4)$$

A matrix with its transpose is equal to its inverse is called an *orthogonal* matrix, and its determinant must equal  $\pm 1$ . The attitude matrix must be a *proper* orthogonal matrix, i.e. have determinant  $+1$ , in order to transform a right-handed coordinate frame to a right-handed coordinate frame.

The nine-component attitude matrix is in some ways the ideal representation of a vehicle’s attitude. It has a 1:1 correspondence with physical attitudes, it varies smoothly as the physical attitude varies smoothly, its elements all have magnitudes less than or equal to one, and it follows a simple rule for combining successive rotations. It is not an efficient representation, though; only three of its nine parameters are independent because the orthogonality constraint is equivalent to six independent scalar constraints. This provides the opportunity to specify an attitude or an attitude matrix using only three parameters, but not, as was pointed out above, in a globally continuous and nonsingular fashion.

---

<sup>1</sup>One example is an incorrect sign for the velocity aberration correction for star tracker measurements on the WMAP spacecraft, which fortunately was easily corrected.

## 8.2. Euler Axis/Angle Representation

The *Euler axis/angle* representation of a rotation matrix is based on *Euler's Theorem*, which states that the general displacement of a rigid body with one point fixed is a rotation about a fixed axis [22]. Specify the axis by a unit vector  $\vec{e}$  and the rotation angle by  $\vartheta$ , and denote the matrix that maps the representations of vectors from frame  $\mathcal{F}$  to frame  $\mathcal{G}$  by  $\mathbf{A}_{\mathcal{G}\mathcal{F}}(\mathbf{e}, \vartheta)$ . The rotation axis is fixed, so  $\mathbf{e}$  can be its representation in either frame  $\mathcal{F}$  or frame  $\mathcal{G}$ , which are identical. Consider the mapping of a vector  $\vec{x}$  whose representation in frame  $\mathcal{F}$  is

$$\mathbf{x}_{\mathcal{F}} = (\mathbf{e}\mathbf{e}^T)\mathbf{x}_{\mathcal{F}} + (\mathbf{I}_3 - \mathbf{e}\mathbf{e}^T)\mathbf{x}_{\mathcal{F}} \quad (8.5)$$

where  $\mathbf{I}_n$  denotes the  $n \times n$  identity matrix. The first term on the right side is parallel and the second is perpendicular to the rotation axis. The rotation does not affect the parallel component and rotates the perpendicular component through an angle  $\vartheta$  around the rotation axis out of the plane defined by that axis and  $\mathbf{x}_{\mathcal{F}}$ , so

$$\mathbf{x}_{\mathcal{G}} = \mathbf{A}_{\mathcal{G}\mathcal{F}}(\mathbf{e}, \vartheta)\mathbf{x}_{\mathcal{F}} = (\mathbf{e}\mathbf{e}^T)\mathbf{x}_{\mathcal{F}} + \cos \vartheta (\mathbf{I}_3 - \mathbf{e}\mathbf{e}^T)\mathbf{x}_{\mathcal{F}} - \sin \vartheta (\mathbf{e} \times \mathbf{x}_{\mathcal{F}}) \quad (8.6)$$

This formula preserves the norm of  $\mathbf{x}_{\mathcal{F}}$  because  $\mathbf{e} \times \mathbf{x}_{\mathcal{F}}$  and  $(\mathbf{I}_3 - \mathbf{e}\mathbf{e}^T)\mathbf{x}_{\mathcal{F}}$  are orthogonal and have equal magnitude. Since  $\mathbf{x}_{\mathcal{F}}$  is an arbitrary vector, Eq. (8.6) means that

$$\mathbf{A}(\mathbf{e}, \vartheta) = (\cos \vartheta) \mathbf{I}_3 - \sin \vartheta [\mathbf{e} \times] + (1 - \cos \vartheta) \mathbf{e}\mathbf{e}^T \quad (8.7)$$

where  $[\mathbf{e} \times]$  is the cross-product matrix:

$$[\mathbf{e} \times] \equiv \begin{bmatrix} 0 & -e_3 & e_2 \\ e_3 & 0 & -e_1 \\ -e_2 & e_1 & 0 \end{bmatrix} \quad (8.8)$$

The cross-product matrix is defined so that  $[\mathbf{x} \times] \mathbf{y} = \mathbf{x} \times \mathbf{y}$ . Equation (8.7) is the *Euler axis/angle* parameterization of an attitude matrix, with explicit frame indices omitted. It requires four parameters, but only three are independent because of the unit vector constraint  $\|\mathbf{e}\| = 1$ .

The  $\sin \vartheta$  terms in Eqs. (8.6) and (8.7) are negative because the rotation in Euler's theorem is an active rotation of the frame  $\mathcal{G}$  relative to the frame  $\mathcal{F}$ , while the rotation matrix  $\mathbf{A}(\mathbf{e}, \vartheta)$  specifies the passive mapping of the representation of  $\vec{x}$  from frame  $\mathcal{F}$  to frame  $\mathcal{G}$ .

The Euler axis/angle representation can be used to find the time dependence of the rotation matrix. The fundamental definition of a derivative gives

$$\begin{aligned} \dot{\mathbf{A}}_{\mathcal{G}\mathcal{F}}(t) &\equiv \lim_{\Delta t \rightarrow 0} \frac{\mathbf{A}_{\mathcal{G}\mathcal{F}}(t + \Delta t) - \mathbf{A}_{\mathcal{G}\mathcal{F}}(t)}{\Delta t} \\ &= \left( \lim_{\Delta t \rightarrow 0} \frac{\mathbf{A}_{\mathcal{G}\mathcal{F}}(t + \Delta t) \mathbf{A}_{\mathcal{F}\mathcal{G}}(t) - \mathbf{I}_3}{\Delta t} \right) \mathbf{A}_{\mathcal{G}\mathcal{F}}(t) \end{aligned} \quad (8.9)$$

because  $\mathbf{A}_{\mathcal{F}\mathcal{G}}(t) \mathbf{A}_{\mathcal{G}\mathcal{F}}(t)$  is equal to the identity matrix. For small  $\Delta t$ , the matrix product  $\mathbf{A}_{\mathcal{G}\mathcal{F}}(t + \Delta t) \mathbf{A}_{\mathcal{F}\mathcal{G}}(t)$  differs from the identity matrix by a small rotation, so it can be represented by a small angle approximation of Eq. (8.7):

$$\mathbf{A}_{\mathcal{G}\mathcal{F}}(t + \Delta t) \mathbf{A}_{\mathcal{F}\mathcal{G}}(t) \approx \mathbf{I}_3 - \vartheta [\mathbf{e} \times] \quad (8.10)$$

Inserting this into Eq. (8.9), taking the limit of as  $\Delta t$  goes to zero, and omitting time arguments gives

$$\dot{\mathbf{A}}_{\mathcal{G}\mathcal{F}} = -[\boldsymbol{\omega}_{\mathcal{G}}^{\mathcal{G}\mathcal{F}} \times] \mathbf{A}_{\mathcal{G}\mathcal{F}} \quad (8.11)$$

where

$$\boldsymbol{\omega}_{\mathcal{G}}^{\mathcal{G}\mathcal{F}} \equiv \lim_{\Delta t \rightarrow 0} \frac{\vartheta \mathbf{e}}{\Delta t} \quad (8.12)$$

is the vector representation in frame  $\mathcal{G}$  of the angular velocity of frame  $\mathcal{G}$  with respect to frame  $\mathcal{F}$ . The angular velocity is known to be represented in frame  $\mathcal{G}$  because the product  $\mathbf{A}_{\mathcal{G}\mathcal{F}}(t + \Delta t)\mathbf{A}_{\mathcal{F}\mathcal{G}}(t)$  is a rotation from frame  $\mathcal{G}$  at one time to frame  $\mathcal{G}$  at a different time, and these two frames coincide in the limit that  $\Delta t$  goes to zero. The angular velocity is usually written simply as  $\boldsymbol{\omega}$ , with the frames understood. Its units are rad/sec, assuming that time is measured in seconds, because radian measure was assumed in taking the small angle limit of  $\sin \vartheta$ .

Equation (8.11) is the fundamental equation of attitude kinematics. It does not distinguish between the situations where frame  $\mathcal{F}$  or frame  $\mathcal{G}$  or both frames are rotating in an absolute sense; it only cares about the relative rotation between the two frames. This equation can also be written as

$$\dot{\mathbf{A}}_{\mathcal{G}\mathcal{F}} = -\mathbf{A}_{\mathcal{G}\mathcal{F}}\mathbf{A}_{\mathcal{F}\mathcal{G}}[\boldsymbol{\omega}_{\mathcal{G}}^{\mathcal{G}\mathcal{F}} \times] \mathbf{A}_{\mathcal{F}\mathcal{G}}^{\top} = -\mathbf{A}_{\mathcal{G}\mathcal{F}}[\mathbf{A}_{\mathcal{F}\mathcal{G}}\boldsymbol{\omega}_{\mathcal{G}}^{\mathcal{G}\mathcal{F}} \times] = -\mathbf{A}_{\mathcal{G}\mathcal{F}}[\boldsymbol{\omega}_{\mathcal{F}}^{\mathcal{G}\mathcal{F}} \times] \quad (8.13)$$

which expresses the kinematics in terms of the representation in frame  $\mathcal{F}$  of  $\boldsymbol{\omega}^{\mathcal{G}\mathcal{F}}$ . The second equality uses an identity that holds for any proper orthogonal matrix. These kinematic equations, if integrated exactly, preserve the orthogonality of the attitude matrix.

The Euler axis/angle representation is fundamental for analysis, as demonstrated above, but it has been entirely superseded for practical applications by a superior four-parameter representation described in the next section.

### 8.3. Quaternion Representation

Applying trigonometric half-angle identities to Eq. (8.7) gives

$$\mathbf{A}(\mathbf{q}) = (q_4^2 - \|\mathbf{q}_{1:3}\|^2) \mathbf{I}_3 - 2q_4[\mathbf{q}_{1:3} \times] + 2\mathbf{q}_{1:3} \mathbf{q}_{1:3}^{\top} \quad (8.14)$$

where the three-component vector  $\mathbf{q}_{1:3}$  and the scalar  $q_4$  are defined by

$$\mathbf{q}_{1:3} = \mathbf{e} \sin(\vartheta/2) \quad (8.15a)$$

$$q_4 = \cos(\vartheta/2) \quad (8.15b)$$

This representation has the advantage over the Euler axis/angle representation of requiring no trigonometric function evaluations, and its four components are more economical than the nine-component attitude matrix.

The four parameters of this representation were first considered by Euler but their full significance was revealed by Rodrigues, so they are often referred to as the *Euler symmetric parameters* or *Euler-Rodrigues parameters*. This representation is called the quaternion representation and denoted  $\mathbf{A}(\mathbf{q})$  because the four parameters can be regarded as the components of a quaternion

$$\mathbf{q} = \begin{bmatrix} \mathbf{q}_{1:3} \\ q_4 \end{bmatrix} \quad (8.16)$$

with vector part  $\mathbf{q}_{1:3}$  and scalar  $q_4$ . A quaternion is basically four-component vector with some additional operations defined for it.<sup>2</sup> The attitude quaternion

$$\mathbf{q}(\mathbf{e}, \vartheta) = \begin{bmatrix} \mathbf{e} \sin(\vartheta/2) \\ \cos(\vartheta/2) \end{bmatrix} \quad (8.17)$$

is a unit quaternion, obeying the norm constraint  $\|\mathbf{q}\| = 1$ , but not all quaternions are unit quaternions. It is clear from Eq. (8.14) that  $\mathbf{q}$  and  $-\mathbf{q}$  give the same attitude matrix. This 2:1 mapping of quaternions to rotations is a minor annoyance that cannot be removed without introducing discontinuities in the representation.

The most important added quaternion operations are two different products of two quaternions  $\mathbf{q}$  and  $\bar{\mathbf{q}}$ . They can be implemented in matrix form similar to the matrix form of the vector cross product:<sup>3</sup>

$$\mathbf{q} \otimes \bar{\mathbf{q}} = [\Psi(\mathbf{q}) \ \mathbf{q}] \bar{\mathbf{q}} \quad (8.18a)$$

$$\mathbf{q} \odot \bar{\mathbf{q}} = [\Xi(\mathbf{q}) \ \mathbf{q}] \bar{\mathbf{q}} \quad (8.18b)$$

where  $\Psi(\mathbf{q})$  and  $\Xi(\mathbf{q})$  are the  $4 \times 3$  matrices

$$\Psi(\mathbf{q}) \equiv \begin{bmatrix} q_4 \mathbf{I}_3 - [\mathbf{q}_{1:3} \times] \\ -\mathbf{q}_{1:3}^T \end{bmatrix} \quad (8.19a)$$

$$\Xi(\mathbf{q}) \equiv \begin{bmatrix} q_4 \mathbf{I}_3 + [\mathbf{q}_{1:3} \times] \\ -\mathbf{q}_{1:3}^T \end{bmatrix} \quad (8.19b)$$

Unlike the vector cross product, though, the norm of either product of two quaternions is equal to the product of their norms.

Both quaternion products are associative but not commutative in general, in parallel with matrix products. The two product definitions differ only in the signs of the cross product matrices in Eqs. (8.19a) and (8.19b), from which it follows that

$$\mathbf{q} \otimes \bar{\mathbf{q}} = \bar{\mathbf{q}} \odot \mathbf{q} \quad (8.20)$$

The identity quaternion

$$\mathbf{I}_q \equiv [0 \ 0 \ 0 \ 1]^T \quad (8.21)$$

acts in quaternion multiplication like the identity matrix in matrix multiplication. The conjugate  $\mathbf{q}^*$  of a quaternion is obtained by changing the sign of its vector part:

$$\mathbf{q}^* = \begin{bmatrix} \mathbf{q}_{1:3} \\ q_4 \end{bmatrix}^* \equiv \begin{bmatrix} -\mathbf{q}_{1:3} \\ q_4 \end{bmatrix} \quad (8.22)$$

Either product of a quaternion with its conjugate is equal to the square of its norm times the identity quaternion.

The inverse of any quaternion having nonzero norm is defined by

$$\mathbf{q}^{-1} \equiv \mathbf{q}^* / \|\mathbf{q}\|^2 \quad (8.23)$$

<sup>2</sup>This is *conceptually* different from the quaternion introduced by Hamilton in 1844, before the introduction of vector notation, as a hypercomplex extension  $q = q_0 + iq_1 + jq_2 + kq_3$  of a complex number  $z = x + iy$ . The scalar part of a quaternion is often labeled  $q_0$  and put at the top of the column vector. Care must be taken to thoroughly understand the conventions embodied in any quaternion equation that one chooses to reference.

<sup>3</sup>The notation  $\mathbf{q} \otimes \bar{\mathbf{q}}$  was introduced in Ref. [43], and  $\mathbf{q} \odot \bar{\mathbf{q}}$  is a modification of notation introduced in Ref. [57]. Hamilton's product  $\bar{q}q$  corresponds to  $\mathbf{q} \odot \bar{\mathbf{q}}$ , but  $\mathbf{q} \otimes \bar{\mathbf{q}}$  has proven to be more useful in attitude analysis. The order of the quaternion products in Eqs. (8.27) and (8.28) would be reversed with the classical definition of quaternion multiplication.

A unit quaternion, such as the attitude quaternion, always has an inverse, which is identical with its conjugate. The conjugate of the product of two quaternions is equal to the product of their conjugates in the opposite order:  $(\bar{\mathbf{q}} \otimes \mathbf{q})^* = \mathbf{q}^* \otimes \bar{\mathbf{q}}^*$ . The same relationship holds for the other product definition and with conjugates replaced by inverses.

Equation (8.14) can be compactly written as

$$\mathbf{A}(\mathbf{q}) = \Xi^T(\mathbf{q})\Psi(\mathbf{q}) \quad (8.24)$$

Now consider, for a unit quaternion  $\mathbf{q}$ , the product

$$\begin{aligned} \mathbf{q} \otimes \begin{bmatrix} \mathbf{x} \\ 0 \end{bmatrix} \otimes \mathbf{q}^* &= \mathbf{q}^* \odot \left( \mathbf{q} \otimes \begin{bmatrix} \mathbf{x} \\ 0 \end{bmatrix} \right) = [\Xi(\mathbf{q}^*) \quad \mathbf{q}^*] [\Psi(\mathbf{q}) \quad \mathbf{q}] \begin{bmatrix} \mathbf{x} \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} \Xi^T(\mathbf{q}) \\ \mathbf{q}^T \end{bmatrix} \Psi(\mathbf{q}) \mathbf{x} = \begin{bmatrix} \mathbf{A}(\mathbf{q}) \mathbf{x} \\ 0 \end{bmatrix} \end{aligned} \quad (8.25)$$

Applying a transformation by a second quaternion  $\bar{\mathbf{q}}$  gives

$$\begin{aligned} \bar{\mathbf{q}} \otimes \left( \mathbf{q} \otimes \begin{bmatrix} \mathbf{x} \\ 0 \end{bmatrix} \otimes \mathbf{q}^* \right) \otimes \bar{\mathbf{q}}^* &= \bar{\mathbf{q}} \otimes \begin{bmatrix} \mathbf{A}(\mathbf{q}) \mathbf{x} \\ 0 \end{bmatrix} \otimes \bar{\mathbf{q}}^* = \begin{bmatrix} \mathbf{A}(\bar{\mathbf{q}}) \mathbf{A}(\mathbf{q}) \mathbf{x} \\ 0 \end{bmatrix} \\ &= (\bar{\mathbf{q}} \otimes \mathbf{q}) \otimes \begin{bmatrix} \mathbf{x} \\ 0 \end{bmatrix} \otimes (\bar{\mathbf{q}} \otimes \mathbf{q})^* = \begin{bmatrix} \mathbf{A}(\bar{\mathbf{q}} \otimes \mathbf{q}) \mathbf{x} \\ 0 \end{bmatrix} \end{aligned} \quad (8.26)$$

Because this must hold for any  $\mathbf{x}$ , it shows that

$$\mathbf{A}(\bar{\mathbf{q}} \otimes \mathbf{q}) = \mathbf{A}(\bar{\mathbf{q}}) \mathbf{A}(\mathbf{q}) \quad (8.27)$$

This and Eq. (8.2) mean that the quaternion corresponding to successive rotations is just the product

$$\mathbf{q}_{\mathcal{H}\mathcal{F}} = \mathbf{q}_{\mathcal{H}\mathcal{G}} \otimes \mathbf{q}_{\mathcal{G}\mathcal{F}} \quad (8.28)$$

A simple bilinear composition rule of this type holds only for the attitude matrix and quaternion representations, a major reason for the popularity of quaternions.

Representing the rotation between times  $t$  and  $t + \Delta t$  by an Euler axis and angle, Eqs. (8.28), (8.17), (8.20), and (8.18b) give

$$\begin{aligned} \mathbf{q}(t + \Delta t) &= \begin{bmatrix} \mathbf{e} \sin(\vartheta/2) \\ \cos(\vartheta/2) \end{bmatrix} \otimes \mathbf{q}(t) = \cos(\vartheta/2) \mathbf{q}(t) + \sin(\vartheta/2) \begin{bmatrix} \mathbf{e} \\ 0 \end{bmatrix} \otimes \mathbf{q}(t) \\ &= \cos(\vartheta/2) \mathbf{q}(t) + \sin(\vartheta/2) \Xi(\mathbf{q}(t)) \mathbf{e} \end{aligned} \quad (8.29)$$

This quaternion propagation equation has proven to be very useful. It preserves the unity norm of the attitude quaternion. If the angular velocity is well approximated as constant over the time interval, then  $\vartheta = \|\boldsymbol{\omega}\| \Delta t$  and  $\mathbf{e} = \boldsymbol{\omega} / \|\boldsymbol{\omega}\|$ . Alternatively, and particularly for onboard applications,  $\vartheta \mathbf{e}$  can be computed by differencing the outputs of rate-integrating gyroscopes.

Using small angle approximations for the sine and cosine leads to the kinematic equation for the quaternion:

$$\dot{\mathbf{q}} \equiv \lim_{\Delta t \rightarrow 0} \frac{\mathbf{q}(t + \Delta t) - \mathbf{q}(t)}{\Delta t} = \frac{1}{2} \begin{bmatrix} \boldsymbol{\omega} \\ 0 \end{bmatrix} \otimes \mathbf{q} = \frac{1}{2} \Xi(\mathbf{q}) \boldsymbol{\omega} \quad (8.30)$$

where  $\boldsymbol{\omega}$  is defined by Eq. (8.12) and several time arguments have been omitted. Exact integration of this equation preserves the unit norm of the quaternion. The inverse of Eq. (8.30) is often useful:

$$\boldsymbol{\omega} = 2\Xi^T(\mathbf{q})\dot{\mathbf{q}} \quad (8.31)$$

The quaternion representation of a given attitude matrix  $\mathbf{A}$  can be found by normalizing any one of the four vectors [46]

$$\begin{aligned} \begin{bmatrix} 1 + 2A_{11} - \text{tr}\mathbf{A} \\ A_{12} + A_{21} \\ A_{13} + A_{31} \\ A_{23} - A_{32} \end{bmatrix} &= 4q_1\mathbf{q}, & \begin{bmatrix} A_{21} + A_{12} \\ 1 + 2A_{22} - \text{tr}\mathbf{A} \\ A_{23} + A_{32} \\ A_{31} - A_{13} \end{bmatrix} &= 4q_2\mathbf{q} \\ \begin{bmatrix} A_{31} + A_{13} \\ A_{32} + A_{23} \\ 1 + 2A_{33} - \text{tr}\mathbf{A} \\ A_{12} - A_{21} \end{bmatrix} &= 4q_3\mathbf{q}, & \begin{bmatrix} A_{23} - A_{32} \\ A_{31} - A_{13} \\ A_{12} - A_{21} \\ 1 + \text{tr}\mathbf{A} \end{bmatrix} &= 4q_4\mathbf{q} \end{aligned} \quad (8.32)$$

Numerical errors are minimized by choosing the vector with the greatest norm, which can be found by the following procedure. Find the largest of  $\text{tr}\mathbf{A}$  and the three diagonal elements of  $\mathbf{A}$ . If  $\text{tr}\mathbf{A}$  is the largest, then  $|q_4|$  is the largest of the  $|q_i|$ , otherwise the largest value of  $|q_i|$  is the one with the same index as the largest diagonal element of  $\mathbf{A}$ . The overall sign of the normalized vector is not determined, reflecting the twofold ambiguity of the quaternion representation.

Extracting many of the other parameterizations from an attitude matrix is most easily accomplished by first extracting a quaternion and then converting to the desired representation. The kinematic equations of these other representations can also be readily derived through the intermediary of the quaternion representation.

A three-parameter representation can be obtained by using only three components of the quaternion, say the  $i, j, k$  components, with the fourth component given by

$$q_\ell = \pm \sqrt{1 - q_i^2 - q_j^2 - q_k^2} \quad (8.33)$$

The sign is not arbitrary, but is determined by the four-component quaternion being represented by three of its components. Once so determined, the sign will not change if the quaternion varies smoothly unless  $q_\ell$  passes through zero. To avoid a sign error if  $|q_\ell|$  becomes small, the representation is switched to make  $q_\ell$  one of three components in the three-parameter representation, replacing one of the other components, which is then given by the square root with the correct sign. The unit norm constraint on the attitude quaternion means that at least one of its components must have magnitude 1/2 or greater. To minimize switching, it should be done when  $|q_\ell|$  is significantly less than 1/2, and the component replaced by  $q_\ell$  in the three-component representation should have magnitude no smaller than 1/2. To first order in the errors, the error in the fourth component of the quaternion is

$$\Delta q_\ell = -q_\ell^{-1}(q_i\Delta q_i + q_j\Delta q_j + q_k\Delta q_k) \quad (8.34)$$

This approaches the indeterminate quantity 0/0 as  $|q_\ell| \rightarrow 0$ , providing another reason for switching.

#### 8.4. Rodrigues Parameter Representation

The three *Rodrigues parameters* appeared in 1840 [61], but were first arranged in a vector by Gibbs, who invented modern vector notation. For this reason, the vector of Rodrigues parameters is often called the *Gibbs vector* and denoted by  $\mathbf{g}$ . The Rodrigues parameters are related to the Euler axis/angle and the quaternion by

$$\mathbf{g} = \mathbf{e} \tan(\vartheta/2) = \frac{\mathbf{q}_{1:3}}{q_4} \quad (8.35)$$

The quaternion as a function of the Gibbs vector is

$$\mathbf{q} = \frac{\pm 1}{\sqrt{1 + \|\mathbf{g}\|^2}} \begin{bmatrix} \mathbf{g} \\ 1 \end{bmatrix} \quad (8.36)$$

It is clear from Eq. (8.35) that  $\mathbf{q}$  and  $-\mathbf{q}$  map to the same Gibbs vector, so the Rodrigues parameters provide a 1:1 mapping of rotations. The price paid for this is that the Gibbs vector is infinite for a  $180^\circ$  rotation. Thus this parameterization is not recommended as a global attitude representation, but it provides an excellent representation of small rotations.

The Rodrigues parameter representation of the attitude matrix is

$$\mathbf{A}(\mathbf{g}) = \mathbf{I}_3 + 2 \frac{[\mathbf{g} \times]^2 - [\mathbf{g} \times]}{1 + \|\mathbf{g}\|^2} \quad (8.37)$$

This resembles the quaternion representation in requiring no transcendental function evaluations, but it is a rational function rather than a simple polynomial.

The composition rule for the Rodrigues parameters corresponding to the quaternion product  $\bar{\mathbf{q}} \boxtimes \mathbf{q}$  is

$$\bar{\mathbf{g}} \boxtimes \mathbf{g} \equiv \frac{\bar{\mathbf{g}} + \mathbf{g} - \bar{\mathbf{g}} \times \mathbf{g}}{1 - \bar{\mathbf{g}} \cdot \mathbf{g}} \quad (8.38)$$

This composition law is associative but not commutative in general, in parallel with matrix and quaternion products. Because it is not a bilinear function of the constituent Gibbs vectors, it cannot be represented as a matrix product like quaternion composition.

The kinematic equation for the Rodrigues parameters is

$$\dot{\mathbf{g}} = (1/2) (\mathbf{I}_3 + [\mathbf{g} \times] + \mathbf{g}\mathbf{g}^\top) \boldsymbol{\omega} \quad (8.39)$$

with the inverse

$$\boldsymbol{\omega} = 2 (1 + \|\mathbf{g}\|^2)^{-1} (\dot{\mathbf{g}} - \mathbf{g} \times \dot{\mathbf{g}}) \quad (8.40)$$

### 8.5. Modified Rodrigues Parameters

The modified Rodrigues parameters (MRPs) were invented by T. F. Wiener in 1962 [75], rediscovered by Marandi and Modi in 1987 [45], and have been championed by Junkins and Schaub [62]. The vector of MRPs is related to the Euler axis/angle and the quaternion by

$$\mathbf{p} = \mathbf{e} \tan(\vartheta/4) = \frac{\mathbf{q}_{1:3}}{1 + q_4} \quad (8.41)$$

The quaternion as a function of the MRPs is

$$\mathbf{q} = \frac{1}{1 + \|\mathbf{p}\|^2} \begin{bmatrix} 2\mathbf{p} \\ 1 - \|\mathbf{p}\|^2 \end{bmatrix} \quad (8.42)$$

and the attitude matrix is given by

$$\mathbf{A}(\mathbf{p}) = \mathbf{I}_3 + \frac{8 [\mathbf{p} \times]^2 - 4 (1 - \|\mathbf{p}\|^2) [\mathbf{p} \times]}{(1 + \|\mathbf{p}\|^2)^2} \quad (8.43)$$

Every vector of MRPs has a *shadow*

$$\mathbf{p}^S \equiv -\frac{\mathbf{p}}{\|\mathbf{p}\|^2} = \frac{-\mathbf{q}_{1:3}}{1 - q_4} \quad (8.44)$$

An MRP vector and its shadow represent the same attitude because  $\mathbf{q}$  and  $-\mathbf{q}$  represent the same attitude, so the MRPs are a 2:1 mapping of the rotations just as the quaternions are. It is clear from Eq. (8.44) that  $\|\mathbf{p}^S\| \|\mathbf{p}\| = 1$ , so every attitude can be represented

by either an MRP vector with  $\|\mathbf{p}\| \leq 1$  or an equivalent MRP vector in the *shadow set* of MRPs with  $\|\mathbf{p}\| \geq 1$ .

The MRP vector corresponding to the quaternion product  $\bar{\mathbf{q}} \otimes \mathbf{q}$  follows the composition rule

$$\bar{\mathbf{p}} \square \mathbf{p} \equiv \frac{(1 - \|\mathbf{p}\|^2) \bar{\mathbf{p}} + (1 - \|\bar{\mathbf{p}}\|^2) \mathbf{p} - 2 \bar{\mathbf{p}} \times \mathbf{p}}{1 + \|\mathbf{p}\|^2 \|\bar{\mathbf{p}}\|^2 - 2 \bar{\mathbf{p}} \cdot \mathbf{p}} \quad (8.45)$$

This composition law is associative but not commutative in general, in parallel with matrix and quaternion products. It cannot be represented as a matrix product.

The kinematic equation for the MRPs is

$$\dot{\mathbf{p}} = \frac{1 + \|\mathbf{p}\|^2}{4} \left( \mathbf{I}_3 + 2 \frac{[\mathbf{p} \times]^2 + [\mathbf{p} \times]}{1 + \|\mathbf{p}\|^2} \right) \boldsymbol{\omega} \quad (8.46)$$

The matrix in parentheses is orthogonal, so the inverse of Eq. (8.46) is

$$\boldsymbol{\omega} = \frac{4}{1 + \|\mathbf{p}\|^2} \left( \mathbf{I}_3 + 2 \frac{[\mathbf{p} \times]^2 - [\mathbf{p} \times]}{1 + \|\mathbf{p}\|^2} \right) \dot{\mathbf{p}} \quad (8.47)$$

The norm of an MRP vector can grow without limit during dynamic propagation or attitude estimation; Eq. (8.41) shows that the norm is infinite for  $\vartheta = 2\pi$ . This singularity can be avoided by switching from the MRP vector to its shadow. The norm can be restricted to be less than or equal to unity in theory, but in practice it is best to allow the norm to exceed unity by some amount before switching in order to avoid “chattering” between the MRP and its shadow. An error  $\Delta \mathbf{p}$  in an MRP vector corresponds to an error

$$\Delta \mathbf{p}^S = \frac{\partial \mathbf{p}^S}{\partial \mathbf{p}} \Delta \mathbf{p} = \frac{2\mathbf{p}\mathbf{p}^\top - \|\mathbf{p}\|^2 \mathbf{I}_3}{\|\mathbf{p}\|^4} \Delta \mathbf{p} = [2\mathbf{p}^S (\mathbf{p}^S)^\top - \|\mathbf{p}^S\|^2 \mathbf{I}_3] \Delta \mathbf{p} \quad (8.48)$$

in its shadow. This relation is useful for mapping covariance matrices into and out of the shadow set. The 2:1 four-component quaternion representation does not have these complications because the two quaternions representing the same attitude both have unit norm, so there is no need to switch between them.

### 8.6. Rotation Vector Representation

It is convenient for analysis, but not for computations, to combine the Euler axis and angle into the three-component *rotation vector*

$$\boldsymbol{\vartheta} \equiv \vartheta \mathbf{e} = 2(\cos^{-1} q_4) \frac{\mathbf{q}_{1:3}}{\|\mathbf{q}_{1:3}\|} \quad (8.49)$$

This leads to the very elegant expression

$$\mathbf{A}(\boldsymbol{\vartheta}) = \exp(-[\boldsymbol{\vartheta} \times]) \quad (8.50)$$

where  $\exp(\cdot)$  is the matrix exponential. This equation can be verified by expansion of it and Eq. (8.7) as Taylor series in  $\vartheta$  and repeated applications of the identity  $[\mathbf{e} \times]^2 = \mathbf{e}\mathbf{e}^\top - \mathbf{I}_3$ .

The kinematic equation for the rotation vector is

$$\dot{\boldsymbol{\vartheta}} = \boldsymbol{\omega} + \frac{1}{2} \boldsymbol{\vartheta} \times \boldsymbol{\omega} + \frac{1}{\vartheta^2} \left( 1 - \frac{\vartheta}{2} \cot \frac{\vartheta}{2} \right) \boldsymbol{\vartheta} \times (\boldsymbol{\vartheta} \times \boldsymbol{\omega}) \quad (8.51)$$

The coefficient of  $\boldsymbol{\vartheta} \times (\boldsymbol{\vartheta} \times \boldsymbol{\omega})$  goes to 1/12 as  $\vartheta$  goes to zero, but it is singular for  $\vartheta$  equal to any nonzero multiple of  $2\pi$ . The inverse of Eq. (8.51) is

$$\boldsymbol{\omega} = \dot{\boldsymbol{\vartheta}} - \frac{1 - \cos \vartheta}{\vartheta^2} \boldsymbol{\vartheta} \times \dot{\boldsymbol{\vartheta}} + \frac{\vartheta - \sin \vartheta}{\vartheta^3} \boldsymbol{\vartheta} \times (\boldsymbol{\vartheta} \times \dot{\boldsymbol{\vartheta}}) \quad (8.52)$$

The rotation vector is useful for the analysis of small rotations, but not for large rotations, because of both the computational cost of evaluating the matrix exponential and the kinematic singularity for  $\|\boldsymbol{\vartheta}\| = 2\pi$ . This singularity can be avoided, as for the MRPs, by switching from the rotation vector to its shadow

$$\boldsymbol{\vartheta}^S \equiv (1 - 2\pi\|\boldsymbol{\vartheta}\|^{-1})\boldsymbol{\vartheta} \quad (8.53)$$

which represents the same attitude. This can restrict the norm of the rotation vector to  $\pi$  or less in theory, but in practice it is best to allow the norm to exceed  $\pi$  by some amount before switching in order to avoid “chattering” between the rotation vector and its shadow.

The properties of the rotation vector are very similar to those of the MRPs, and it has no obvious advantages over the MRPs. It has the disadvantage of requiring transcendental function evaluations to compute the attitude matrix, so it is rarely used in practical applications.

### 8.7. Euler Angles

An Euler angle representation parameterizes a rotation by the product of three rotations about coordinate frame unit vectors:

$$\mathbf{A}_{ijk}(\phi, \theta, \psi) = \mathbf{A}(\mathbf{e}_k, \psi)\mathbf{A}(\mathbf{e}_j, \theta)\mathbf{A}(\mathbf{e}_i, \phi) \quad (8.54)$$

where  $\mathbf{e}_j$ ,  $\mathbf{e}_j$ , and  $\mathbf{e}_j$  are selected from the set

$$\mathbf{e}_1 = [1 \ 0 \ 0]^\top, \quad \mathbf{e}_2 = [0 \ 1 \ 0]^\top, \quad \mathbf{e}_3 = [0 \ 0 \ 1]^\top \quad (8.55)$$

The possible choice of axes is constrained by the requirements  $i \neq j$  and  $j \neq k$ , leaving six symmetric sets with  $ijk$  equal to 121, 131, 232, 212, 313, and 323 and six asymmetric sets with  $ijk$  equal to 123, 132, 231, 213, 312, and 321. Symmetric Euler angle sets are used in classical studies of rigid body motion [22, 28, 35, 49, 74].

The asymmetric sets of angles are called the *Tait-Bryan* angles or *roll*, *pitch*, and *yaw* angles. The latter terminology originally described the motions of ships and then was carried over into aircraft and spacecraft. Roll is a rotation about the vehicle body axis that is closest to the vehicle’s usual direction of motion, and hence would be perceived as a screwing motion. The roll axis is conventionally assigned index 1. Yaw is a rotation about the vehicle body axis that is usually closest to the direction of local gravity, and hence would be perceived as a motion that points the vehicle left or right. The yaw axis is conventionally assigned index 3. Pitch is a rotation about the remaining vehicle body axis, and hence would be perceived as a motion that points the vehicle up or down. The pitch axis is conventionally assigned index 2. Note that Eq. (8.54) assigns the variables  $\phi$ ,  $\theta$ , and  $\psi$  based on the order of rotations in the sequence, making no definite association between these variables and either the axis labels 1, 2, and 3 or the names roll, pitch and yaw. Many authors follow a different convention, denoting roll by  $\phi$ , pitch by  $\theta$ , and yaw by  $\psi$ . As always, the reader consulting any source should be careful to understand the conventions that it follows.

Using the product rule and Eq. (8.11) to compute the time derivative of Eq. (8.54) gives

$$\begin{aligned} -[\boldsymbol{\omega} \times] \mathbf{A}_{ijk}(\phi, \theta, \psi) = & \left\{ -[(\dot{\psi} \mathbf{e}_k) \times] - \mathbf{A}(\mathbf{e}_k, \psi)[(\dot{\theta} \mathbf{e}_j) \times] \mathbf{A}^\top(\mathbf{e}_k, \psi) \right. \\ & \left. - \mathbf{A}(\mathbf{e}_k, \psi) \mathbf{A}(\mathbf{e}_j, \theta)[(\dot{\phi} \mathbf{e}_i) \times] \mathbf{A}^\top(\mathbf{e}_j, \theta) \mathbf{A}^\top(\mathbf{e}_k, \psi) \right\} \mathbf{A}_{ijk}(\phi, \theta, \psi) \end{aligned} \quad (8.56)$$

The identity  $\mathbf{A}[\mathbf{x} \times] \mathbf{A}^\top = [(\mathbf{A}\mathbf{x}) \times]$ , which holds for any proper orthogonal  $\mathbf{A}$ , gives

$$\boldsymbol{\omega} = \dot{\psi} \mathbf{e}_k + \dot{\theta} \mathbf{A}(\mathbf{e}_k, \psi) \mathbf{e}_j + \dot{\phi} \mathbf{A}(\mathbf{e}_k, \psi) \mathbf{A}(\mathbf{e}_j, \theta) \mathbf{e}_i = \mathbf{A}(\mathbf{e}_k, \psi) \mathbf{M}[\dot{\phi} \ \dot{\theta} \ \dot{\psi}]^\top \quad (8.57)$$

where

$$\mathbf{M} = [\mathbf{A}(\mathbf{e}_j, \theta)\mathbf{e}_i \quad \mathbf{e}_j \quad \mathbf{e}_k] \quad (8.58)$$

The second equality in Eq. (8.57) makes use of  $\mathbf{A}(\mathbf{e}_k, \psi)\mathbf{e}_k = \mathbf{e}_k$ . The Euler angle rates as functions of the angular velocity are

$$[\dot{\phi} \quad \dot{\theta} \quad \dot{\psi}]^T = \mathbf{M}^{-1}\mathbf{A}^T(\mathbf{e}_k, \psi)\boldsymbol{\omega} \quad (8.59)$$

This kinematic equation is singular if the determinant of  $\mathbf{M}$  is zero. Equation (8.7) and the Euler axis requirement that  $\mathbf{e}_i \cdot \mathbf{e}_j = \mathbf{e}_j \cdot \mathbf{e}_k = 0$  gives

$$\det \mathbf{M} = [\mathbf{A}(\mathbf{e}_j, \theta)\mathbf{e}_i] \cdot (\mathbf{e}_j \times \mathbf{e}_k) = \cos \theta [\mathbf{e}_i \cdot (\mathbf{e}_j \times \mathbf{e}_k)] - \sin \theta (\mathbf{e}_i \cdot \mathbf{e}_k) \quad (8.60)$$

The triple vector product  $\mathbf{e}_i \cdot (\mathbf{e}_j \times \mathbf{e}_k)$  is zero for the symmetric Euler angles, so the kinematic equations of these representations are singular if  $\sin \theta = 0$ . The dot product  $\mathbf{e}_i \cdot \mathbf{e}_k$  is zero for the asymmetric Euler angles, so the kinematics of these representations are singular if  $\cos \theta = 0$ . This singularity is known as *gimbal lock* and is caused by collinearity of the *physical* rotation axis vectors of the first and third rotations in the sequence. Note that the column vector representations of these rotation axes are always parallel for the symmetric Euler angle sequences and always perpendicular for the asymmetric sequences, but this neither causes nor prevents gimbal lock.

Because Euler angles are discussed in many references on rotational motion and because they are not widely used in navigation filters, they will not be discussed further here. Kinematic equations and explicit forms of the attitude matrices for all twelve sets can be found in Refs. [28, 35, 49, 74].

### 8.8. Additive EKF (AEKF)

Three-component representations are the most natural representations for filtering, because only three parameters are needed to represent rotations. As was pointed out at the beginning of this Chapter, though, all three-parameter representations of the rotation group have discontinuities or singularities. Any filter using a three-dimensional attitude representation must provide some guarantee of avoiding these singular points, either by restricting the vehicle's attitude or by switching between different parameter sets if the representation approaches a discontinuity or singularity.

The earliest Kalman filters for spacecraft attitude estimation used a roll, pitch, yaw sequence of Euler or Tait-Bryan angles discussed in Section 8.7 [14, 15]. This is a very useful representation if the middle angle of the sequence, generally the pitch angle, stays well away from  $\pm 90^\circ$ , and has been used for Earth-pointing spacecraft with small pitch angles. One disadvantage of this representation is its trigonometric function evaluations, but this is less of an issue with the computing power now available, especially in onboard computers.

An EKF can estimate three components of a quaternion, with the fourth component being determined by the quaternion unit norm constraint, as discussed at the end of Section 8.3 [43]. If the fourth component becomes small, it must be added to the set of estimated quaternion components, with one of the other components switched out. This switch should be made when the magnitude of the fourth component is not too close to either end of the range from 0 and 1/2 to avoid "chattering" between component sets. The switch must be accompanied by the following covariance matrix transformation. Assuming the three components in the pre-switch representation have indices  $i, j, k$  in ascending order, their

$3 \times 3$  covariance matrix is the symmetric matrix

$$\mathbf{P}_{3 \times 3} = \begin{bmatrix} P_{ii} & P_{ij} & P_{ik} \\ P_{ji} & P_{jj} & P_{jk} \\ P_{ki} & P_{kj} & P_{kk} \end{bmatrix} \quad (8.61)$$

The  $4 \times 4$  covariance matrix of the full quaternion is formed by adding the  $\ell$ th row and column, keeping the indices in ascending order. The needed covariance components, using using Eqs. (8.33) and (8.34), are

$$P_{m\ell} = P_{\ell m} = -q_\ell^{-1}(q_i P_{im} + q_j P_{jm} + q_k P_{km}) \quad \text{for } m = i, j, k \quad (8.62a)$$

$$P_{\ell\ell} = q_\ell^{-2}(q_i^2 P_{ii} + q_j^2 P_{jj} + q_k^2 P_{kk} + 2q_i q_j P_{ji} + 2q_i q_k P_{ik} + 2q_j q_k P_{jk}) \quad (8.62b)$$

Then the row and column corresponding to the quaternion component switched out of the representation are deleted to form the new  $3 \times 3$  covariance matrix.

The modified Rodrigues parameters (MRPs) are non-singular for rotations of less than  $360^\circ$ , and the singularity can be avoided by switching to an MRP in the shadow set, as discussed in Section 8.5. The switch to the shadow set is made at some angle greater than  $180^\circ$  to avoid “chattering” between the two parameter sets. The switch must be accompanied by a covariance matrix transformation using Eq. (8.48)

$$\mathbf{P}_{pp}^S = [2\mathbf{p}^S(\mathbf{p}^S)^\top - \|\mathbf{p}^S\|^2 \mathbf{I}_3] \mathbf{P}_{pp} [2\mathbf{p}^S(\mathbf{p}^S)^\top - \|\mathbf{p}^S\|^2 \mathbf{I}_3] \quad (8.63)$$

where  $\mathbf{P}_{pp}$  is the covariance before the switch, and  $\mathbf{P}_{pp}^S$  is the covariance of  $\mathbf{p}^S$  after the switch [37]. This covariance mapping is simpler than the corresponding mapping for the three-component quaternion representation, and the MRP representation avoids a square root computation. The appearance of  $\vartheta/4$  in Eq. (8.41) as opposed to  $\vartheta/2$  in Eq. (8.17) means that switching will be less frequent when the MRP representation is used. For these reasons the MRP representation has become the preferred three-parameter attitude filter when the attitude is unrestricted.

The Gibbs vector or Rodrigues parameter representation has been used in an EKF [29], but it is not well suited to filtering because of its inability to represent  $180^\circ$  rotations, as discussed in Section 8.4. It provides an excellent representation of small attitude errors, however. The rotation vector, discussed in Section 8.6, is also not recommended for application in an EKF, as it has no clear advantage over the MRPs and has the disadvantage of requiring transcendental function evaluations.

The nine-component attitude matrix and the four-component quaternion represent the entire rotation group without singularities or discontinuities, but the linear measurement update in EKFs employing these representations violates the orthogonality constraint on the attitude matrix or the unit norm constraint on the quaternion. Various methods have been proposed to circumvent this difficulty, but these are not without problems [49]. At the very least, they are inefficient due to the larger dimensionality of the covariance matrix.

### 8.9. Multiplicative EKF (MEKF)

The multiplicative EKF (MEKF) uses the nine-component attitude matrix or the four-component quaternion as the “global” attitude representation and a three-component vector  $\delta\boldsymbol{\vartheta}$  for the “local” representation of attitude errors, so that the true attitude is represented as a product

$$\mathbf{A}^{\text{true}} = \delta\mathbf{A}(\delta\boldsymbol{\vartheta})\hat{\mathbf{A}} \quad (8.64a)$$

$$\mathbf{q}^{\text{true}} = \delta\mathbf{q}(\delta\boldsymbol{\vartheta}) \otimes \hat{\mathbf{q}} \quad (8.64b)$$

The constraints on the representations are satisfied because  $\mathbf{A}^{\text{true}}$ ,  $\delta\mathbf{A}$ , and  $\hat{\mathbf{A}}$  are all proper orthogonal matrices, and  $\mathbf{q}^{\text{true}}$ ,  $\delta\mathbf{q}$ , and  $\hat{\mathbf{q}}$  all have unit norm. The MEKF avoids the attitude restrictions or switching required by additive attitude EKFs because the error vector  $\delta\boldsymbol{\vartheta}$  is assumed to be small enough to completely avoid singularities in the parameterizations  $\delta\mathbf{A}(\delta\boldsymbol{\vartheta})$  or  $\delta\mathbf{q}(\delta\boldsymbol{\vartheta})$ . In some sense, though, Eq. (8.64) incorporates a continuous switching of the attitude reference.

Only the quaternion version of the MEKF, which is much more widely employed, is presented here. Reference [43] reviews its history. Any three-component attitude representation that is related to first order in  $\delta\boldsymbol{\vartheta}$  to the quaternion by

$$\delta\mathbf{q} \approx \begin{bmatrix} \delta\boldsymbol{\vartheta}/2 \\ 1 \end{bmatrix} \quad (8.65)$$

can be used as the error vector. Common choices are the rotation vector, as suggested by the notation of Eq. (8.64), two times the vector part of the quaternion, two times the vector of Rodrigues parameters, four times the vector of MRPs, or a vector of suitably indexed roll, pitch, and yaw angles [49].

The order of the factors on the right side of Eq. (8.64) means that the attitude errors are in the body reference frame. This leads to a major advantage of the MEKF, that the covariance of the attitude error angles in the body frame has a transparent physical interpretation. The covariance of estimators using other attitude representations has a less obvious interpretation unless the attitude matrix is close to the identity matrix. It is possible to reverse the order of the factors on the right side of Eq. (8.64) so the attitude errors are in the reference frame [19]. The covariance can be mapped into the body frame if desired.

The MEKF estimates  $\delta\boldsymbol{\vartheta}$  and any other state variables of interest. This discussion addresses only the attitude part, as the equations for the other components of the state vector obey the usual EKF equations. The MEKF proceeds by iteration of three steps: measurement update, state vector reset, and propagation to the next measurement time. The measurement update step updates the local error state vector. The reset moves the updated information from the local error state to the global attitude representation and resets the components of the local error state to zero. The propagation step propagates the global variables to the time of the next measurement. The local error state variables do not need to be propagated because they are identically zero over the propagation step. These three steps will be discussed in more detail.

**8.9.1. Measurement Update** The observation model is given in terms of the true global state

$$\mathbf{y} = \mathbf{h}(\mathbf{q}^{\text{true}}) + \mathbf{v} \quad (8.66)$$

but the measurement sensitivity matrix is the partial derivative with respect to the local error state, so the measurement sensitivity matrix is

$$\mathbf{H} = \frac{\partial \mathbf{h}}{\partial \mathbf{q}} \frac{\partial \mathbf{q}^{\text{true}}}{\partial (\delta\boldsymbol{\vartheta})} \quad (8.67)$$

Equations (8.64), (8.65), (8.18b), and (8.20) give, to first order in the error vector,

$$\mathbf{q}^{\text{true}} \approx \begin{bmatrix} \delta\boldsymbol{\vartheta}/2 \\ 1 \end{bmatrix} \otimes \hat{\mathbf{q}} = \hat{\mathbf{q}} + \frac{1}{2}\Xi(\hat{\mathbf{q}})\delta\boldsymbol{\vartheta} \quad (8.68)$$

Inserting the partial derivative of this with respect to  $\delta\boldsymbol{\vartheta}$  into Eq. (8.67) then gives

$$\mathbf{H} = \frac{1}{2} \frac{\partial \mathbf{h}}{\partial \mathbf{q}} \Xi(\hat{\mathbf{q}}) \quad (8.69)$$

Simplifications are possible in some special cases.

The Kalman gain computation and covariance update have the standard Kalman filter forms. The error state update employs the first-order Taylor expansion

$$E\{\mathbf{h}(\mathbf{q}^{\text{true}})\} \approx \mathbf{h}(\hat{\mathbf{q}}) + \mathbf{H}\delta\hat{\boldsymbol{\vartheta}} \quad (8.70)$$

giving

$$\delta\hat{\boldsymbol{\vartheta}}^+ = \delta\hat{\boldsymbol{\vartheta}}^- + \mathbf{K} [\mathbf{y} - E\{\mathbf{h}(\mathbf{q}^{\text{true}})\}] = (I - KH)\delta\hat{\boldsymbol{\vartheta}}^- + \mathbf{K} [\mathbf{y} - \mathbf{h}(\hat{\mathbf{q}}^-)] \quad (8.71)$$

**8.9.2. Reset** The discrete measurement update assigns a finite post-update value to  $\delta\hat{\boldsymbol{\vartheta}}^+$ , but the global state still retains the value  $\hat{\mathbf{q}}^-$ . A reset procedure is used to move the update information to a post-update estimate global state vector  $\hat{\mathbf{q}}^+$ , while simultaneously resetting  $\delta\hat{\boldsymbol{\vartheta}}$  to  $\mathbf{0}_3$ , the three-component vector with all zero components. The reset does not change the overall estimate, so the reset must obey

$$\hat{\mathbf{q}}^+ = \delta\mathbf{q}(\mathbf{0}_3) \otimes \hat{\mathbf{q}}^+ = \delta\mathbf{q}(\delta\hat{\boldsymbol{\vartheta}}^+) \otimes \hat{\mathbf{q}}^- \quad (8.72)$$

Thus the reset moves information from one part of the estimate to another part.

Every EKF includes an additive reset of the global state vector, but this is usually implicit rather than explicit. The multiplicative quaternion reset is the special feature of the MEKF. This reset has to preserve the quaternion norm, so an exact unit-norm expression for the functional dependence of  $\delta\mathbf{q}$  on  $\delta\boldsymbol{\vartheta}$  must be used, not the linear approximation of Eq. (8.65). Using the Rodrigues parameter vector has the practical advantage that the reset operation for this parameterization is

$$\hat{\mathbf{q}}^+ = \delta\mathbf{q}(\delta\hat{\boldsymbol{\vartheta}}^+) \otimes \hat{\mathbf{q}}^- = \frac{1}{\sqrt{1 + \|\delta\hat{\boldsymbol{\vartheta}}^+/2\|^2}} \begin{bmatrix} \delta\hat{\boldsymbol{\vartheta}}^+/2 \\ 1 \end{bmatrix} \otimes \hat{\mathbf{q}}^- \quad (8.73)$$

Using an argument similar to Eq. (8.68), this can be accomplished in two steps:

$$\mathbf{q}' = \hat{\mathbf{q}}^- + \frac{1}{2} \Xi(\hat{\mathbf{q}}^-) \delta\hat{\boldsymbol{\vartheta}}^+ \quad (8.74)$$

followed by

$$\hat{\mathbf{q}}^+ = \frac{\mathbf{q}'}{\|\mathbf{q}'\|} \quad (8.75)$$

The first step is just the usual linear Kalman update, and the second step is equivalent to a brute force normalization of the updated quaternion. Thus the MEKF using Rodrigues parameters for the error vector provides a theoretical justification for brute force renormalization, with the added advantage of completely avoiding the accumulation of quaternion norm errors after many updates. The Rodrigues parameters also have the conceptual advantage that they map the rotation group into three-dimensional Euclidean space, with the largest possible  $180^\circ$  attitude errors mapped to points at infinity. Thus probability distributions with infinitely long tails, such as Gaussian distributions, make sense in Rodrigues parameter space.

If a measurement update immediately follows a reset or propagation, the  $\delta\hat{\boldsymbol{\vartheta}}^-$  term on the right side of Eq. (8.71) can be omitted because  $\delta\hat{\boldsymbol{\vartheta}}^-$  is zero. The reset is often delayed

for computational efficiency until all the updates for a set of simultaneous measurements have been performed, though, in which case  $\delta\hat{\boldsymbol{\vartheta}}^-$  may have a finite value and all the terms in Eq. (8.71) must be included. It is imperative to perform a reset before beginning the time propagation, however, to avoid the necessity of propagating  $\delta\hat{\boldsymbol{\vartheta}}$  between measurements.

There is some controversy over the question of whether the reset affects the covariance. One argument holds that it doesn't because the covariance depends not on the actual measurements but on their assumed statistics. Measurement errors are assumed to have zero mean, so the mean reset is zero. But the reset is very different from the measurement update in that it changes the reference frame for the attitude covariance, which might be expected to modify the covariance even though it adds no new information. The change in the covariance of  $\delta\boldsymbol{\vartheta}$  resulting from the effect of the actual update, rather than its zero expectation, can be computed to be [47, 60]

$$\mathbf{P}_{\vartheta\vartheta}^{\text{reset}} = \left( \mathbf{I}_3 - [\delta\hat{\boldsymbol{\vartheta}}^+ \times]/2 \right) \mathbf{P}_{\vartheta\vartheta}^+ \left( \mathbf{I}_3 - [\delta\hat{\boldsymbol{\vartheta}}^+ \times]/2 \right)^{\top} \quad (8.76)$$

to first order in  $\delta\hat{\boldsymbol{\vartheta}}^+$ . Comparison with Eq. (8.7) shows that the covariance reset looks to this order like a rotation by  $\delta\hat{\boldsymbol{\vartheta}}^+/2$ , but this equivalence does not hold in higher orders. Most applications omit this covariance reset, but Reynolds has found that it speeds convergence and adds robustness in the presence of large updates, and that omitting it can even lead to filter divergence in some cases [60].

**8.9.3. Propagation** An EKF must propagate the expectation and covariance of the state. The MEKF is unusual in propagating the expectation  $\hat{\mathbf{q}}$  and the covariance of the error-state vector. The propagation of the attitude error is found by differentiating Eq. (8.64):

$$\dot{\mathbf{q}}^{\text{true}} = \delta\dot{\mathbf{q}} \otimes \hat{\mathbf{q}} + \delta\mathbf{q} \otimes \dot{\hat{\mathbf{q}}} \quad (8.77)$$

The true and estimated quaternions satisfy the kinematic equations

$$\dot{\mathbf{q}}^{\text{true}} = \frac{1}{2} \begin{bmatrix} \boldsymbol{\omega}^{\text{true}} \\ 0 \end{bmatrix} \otimes \mathbf{q}^{\text{true}} \quad (8.78a)$$

$$\dot{\hat{\mathbf{q}}} = \frac{1}{2} \begin{bmatrix} \hat{\boldsymbol{\omega}} \\ 0 \end{bmatrix} \otimes \hat{\mathbf{q}} \quad (8.78b)$$

where  $\boldsymbol{\omega}^{\text{true}}$  and  $\hat{\boldsymbol{\omega}}$  are the true and estimated angular rates, respectively. Substituting these equations and Eq. (8.64) into Eq. (8.77), multiplying on the right by  $\hat{\mathbf{q}}^{-1}$ , and rearranging terms gives [43]

$$\delta\dot{\mathbf{q}} = \frac{1}{2} \left( \begin{bmatrix} \boldsymbol{\omega}^{\text{true}} \\ 0 \end{bmatrix} \otimes \delta\mathbf{q} - \delta\mathbf{q} \otimes \begin{bmatrix} \hat{\boldsymbol{\omega}} \\ 0 \end{bmatrix} \right) \quad (8.79)$$

Substituting Eq. (8.65) and  $\boldsymbol{\omega}^{\text{true}} = \hat{\boldsymbol{\omega}} + \delta\boldsymbol{\omega}$ , where  $\delta\boldsymbol{\omega}$  is the angular velocity error, and multiplying by two leads to

$$\begin{bmatrix} \delta\dot{\boldsymbol{\vartheta}} \\ 0 \end{bmatrix} = \begin{bmatrix} \hat{\boldsymbol{\omega}} \\ 0 \end{bmatrix} \otimes \begin{bmatrix} \delta\boldsymbol{\vartheta}/2 \\ 1 \end{bmatrix} - \begin{bmatrix} \delta\boldsymbol{\vartheta}/2 \\ 1 \end{bmatrix} \otimes \begin{bmatrix} \hat{\boldsymbol{\omega}} \\ 0 \end{bmatrix} + \begin{bmatrix} \delta\boldsymbol{\omega} \\ 0 \end{bmatrix} \otimes \begin{bmatrix} \delta\boldsymbol{\vartheta}/2 \\ 1 \end{bmatrix} \quad (8.80)$$

Ignoring products of the small terms  $\delta\boldsymbol{\omega}$  and  $\delta\boldsymbol{\vartheta}$ , in the spirit of the (linearized) EKF, the first three components of Eq. (8.80) are

$$\delta\dot{\boldsymbol{\vartheta}} = -\hat{\boldsymbol{\omega}} \times \delta\boldsymbol{\vartheta} + \delta\boldsymbol{\omega} \quad (8.81)$$

and the fourth component is  $0 = 0$ . Equation (8.81) is the equation needed to propagate the covariance of the attitude error-angle covariance.

The expectation of Eq. (8.81) is

$$\delta \dot{\hat{\boldsymbol{\vartheta}}} = -\hat{\boldsymbol{\omega}} \times \delta \hat{\boldsymbol{\vartheta}} \quad (8.82)$$

because  $\delta \boldsymbol{\omega}$  has zero expectation. This says that if  $\delta \hat{\boldsymbol{\vartheta}}$  is zero at the beginning of a propagation it will remain zero through the propagation, which is equivalent to saying that  $\delta \hat{\mathbf{q}}$  will be equal to the identity quaternion throughout the propagation.

## Usability Considerations

Contributed by J. Russell Carpenter

This chapter is a catch-all to address best practices for things like selective processing of measurements, backup ephemeris, reinitializations and restarts, availability of ground-commandable tuning parameters, etc.

### 9.1. Editing

Let  $\mathbf{r} = \mathbf{y} - \mathbf{h}(\mathbf{x})$ , where  $\mathbf{y}$  is the observed measurement,  $\mathbf{h}(\mathbf{x})$  is the value of the measurement computed from the state  $\mathbf{x}$ ,  $\mathbf{y} = \mathbf{h}(\mathbf{x}) + \mathbf{v}$ , and  $\mathbf{v}$  is the measurement noise,  $E[\mathbf{v}] = \mathbf{0}$ ,  $E[\mathbf{v}\mathbf{v}^T] = \mathbf{R}$ . The quantity  $\mathbf{r}$  is known as the *innovation* or sometimes the *pre-fit residual*. The covariance of  $\mathbf{r}$  is given by

$$\mathbf{W} = \mathbf{H}\mathbf{P}\mathbf{H}^T + \mathbf{R} \quad (9.1)$$

where  $\mathbf{P} = E[\mathbf{e}\mathbf{e}^T]$ ,  $\mathbf{H} = \partial\mathbf{h}(\mathbf{x})/\partial\mathbf{x}$ , and  $\mathbf{e}$  is the error in the estimate of  $\mathbf{x}$ . The squared *Mahalanobis distance* associated with  $\mathbf{r}$ ,

$$m_{\mathbf{r}}^2 = \mathbf{r}^T\mathbf{W}^{-1}\mathbf{r} \quad (9.2)$$

has a  $\chi^2$  distribution with degrees of freedom equal to the number of measurements contained in the vector  $\mathbf{y}$ . The statistic  $m_{\mathbf{r}}^2$ , also known as the squared residual or innovations ratio, may be compared to a  $\chi^2$  statistic with a given probability in order to edit outlying measurements. For a purely linear estimation scheme, such editing is unnecessary, but for an *ad hoc* linearization such as the EKF, editing is essential to prevent large innovations that would violate Taylor series truncations used to develop the EKF approximation from being violated, even in the unlikely scenario in which sensors produced measurements with noise characteristics that perfectly followed their assumed (Gaussian) probability distributions.

Experience has shown that it is beneficial to provide for a command-able capability to selectively apply a three-way editing flag to each measurement type. This flag may be enumerated with the labels “accept,” “inhibit,” and “force,” or similar. The “accept” label denotes use of the measurement, if it is accepted by the aforementioned edit test. The “inhibit” flag indicates that the measurement should be rejected regardless of the status of the edit test. The “force” flag correspondingly indicates that the measurement should be ingested regardless of the edit test result.

### 9.2. Reinitialization, Restarts, and Backup Ephemeris

Since no EKF can be guaranteed to remain converged, it is prudent to provide for features that can ease the process of recovering nominal filter operation. While to a large extent the particulars of each application will guide a designer to select among a variety of filter recovery features, a few general design principles are broadly relevant.

Experience has shown that the fairly drastic step of completely re-initializing the filter may not always be necessary or prudent. In particular, under conditions in which it is reasonably clear that the filter has begun to edit measurements because its covariance matrix has become overly optimistic, but there remains reason to believe that the state estimate has not yet become corrupted, it may be beneficial to reinitialize the covariance while retaining the current state estimate. It may also be desirable to retain flexibility to retain only the position/velocity state components, while reinitializing the various bias components.

If the filter has halted for some reason other than divergence (e.g. a flight computer reset), or if the start of the divergence can be reliably determined, it may be useful to “restart” the filter from a previously saved state and covariance, especially if there would otherwise be a long time required for re-convergence. To enable such a restart capability, the full state and covariance must have been periodically saved, and then they must be propagated to the restart time.

Periodically saving the filter state also enables the capability to maintain a backup ephemeris, which provides an additional comparison source for evaluating filter divergence. For flight phases of limited duration, the backup ephemeris may be propagated inertially without any measurement updates. For extended operations, it will usually be necessary to re-seed the backup with a current filter state at periodic intervals. In some applications, such as GPS filtering, an independent “point solution” may also serve as a useful comparison source.

### 9.3. Ground System Considerations

While it may seem intuitive that all parameters affecting filter performance should be available for re-tuning from the ground via mechanisms such as commands, table uploads, etc., experience has shown that decisions about ground system design often limit the accessibility of key tuning parameters. Adequate bandwidth in telemetry for full insight into filter performance, including access to full covariance data, must be available, if only for limited periods during commissioning and/or troubleshooting activities. The ground system must be able to reproduce the onboard filter’s performance when provided with corresponding input data via telemetry playback. And the ground system must also be able to form “best estimated trajectories” for comparison to onboard filter performance, e.g. through the use of smoothing algorithms.

## CHAPTER 10

# Smoothing

Contributed by Christopher N. D’Souza and J. Russell Carpenter

Since this work is primarily concerned with onboard navigation filters, one might question the need for a chapter on best practices for smoothing. While the addition of a smoother to an onboard navigation system has usually proved unnecessary, smoothing has nonetheless proved to be a useful ancillary capability for trajectory reconstruction by ground-based analysts. Smoothed trajectories form the basis for our best proxies for truth, in the form of “best estimated trajectories,” (BET) and McReynold’s “filter-smoother consistency test,” propagated by Jim Wright [77], has proven to be a useful aid to tuning a filter using flight data. Also, sequential smoothing techniques can provide optimal estimates of the process noise sequence, as Appendix X of Bierman’s text [2] shows. These estimates may prove useful as part of the filter tuning process.

It is also worth mentioning the topic of “smoothability.” As described in for example Gelb [20], only states that are controllable from the process noise will be affected by smoothing. So for example, estimates of random constant measurement biases cannot be improved by smoothing.

In point of fact there are three types of smoothing: fixed-interval smoothing, fixed-lag smoothing, and fixed-point smoothing. The context described above is concerned with fixed-interval smoothing. Maximum likelihood estimation (MLE) of states over a fixed interval has been subject of investigation ever since the advent of the Kalman filter [36] in 1961. In 1962, Bryson and Frazier first approached the problem from a continuous time perspective [4] and obtained the smoother equations as necessary conditions of an optimal control problem<sup>†</sup>. In 1965, Rauch, Tung and Striebel [59] (RTS) continued the development of the MLE filters but from a discrete time perspective. Their smoother, soon called the RTS smoother, was widely used because of its ease of implementation. However, as Bierman [2] and others [51] pointed out, there can sometimes arise numerical difficulties in implementing the RTS smoother. A short time later Fraser and Potter [17, 18] approached the problem a bit differently, looking at smoothing as a optimal combination of two optimal linear filters and obtained different, yet equivalent, equations. Bierman’s Square-Root Information Filter [2] (SRIF) also has an accompanying smoother form, suitable for applications utilizing the SRIF. Since the Fraser-Potter form avoids the numerical issues of the RTS form, and since it can be easily adapted from existing onboard Kalman-type forward filtering algorithms, it is generally to be preferred.

---

<sup>†</sup>The Bryson-Frazier smoother is a continuous time instantiation of the smoother. It won’t be considered here for we are interested in discrete smoothers. [4]

The boundary conditions for the Fraser-Potter (FP) smoother require the backward covariance at the final time to be infinite, and the backward filter's final state to be zero. Fraser and Potter avoided the infinity by maintaining the backward filter covariance in information form, so that both the information matrix and the information vector are zero. As Brown points out [3], the backward filter may be retained in covariance form, and the infinite boundary condition covariance replaced by either a covariance that is many multiples of the forward filter's initial covariance, or by the covariance and state from a short batch least-squares solution using the final few measurements. Many practical smoother implementations used by NASA have followed an approach of this sort.

Thus, a practical covariance form of the FP smoother results from running whatever implementation of Algorithm 1.3 has proved suitable for the application at hand, but in reverse time, and combining the backward filter results with the forward filter results at each measurement time. Given the forward filter state and covariance,  $\hat{\mathbf{X}}_{F_i}^+$  and  $\mathbf{P}_{F_i}^+$ , which include the measurement at  $t_i$ , and the backward filter state and covariance,  $\hat{\mathbf{X}}_{B_i}^-$  and  $\mathbf{P}_{B_i}^-$ , which *do not* include the measurement at  $t_i$ , the optimally smoothed state and covariance at  $t_i$  are given by

$$\hat{\mathbf{X}}_i^S = \mathbf{P}_i^S \left[ (\mathbf{P}_{F_i}^+)^{-1} \hat{\mathbf{X}}_{F_i}^+ + (\mathbf{P}_{B_i}^-)^{-1} \hat{\mathbf{X}}_{B_i}^- \right] \quad (10.1)$$

$$\mathbf{P}_i^S = \left[ (\mathbf{P}_{F_i}^+)^{-1} + (\mathbf{P}_{B_i}^-)^{-1} \right]^{-1} \quad (10.2)$$

If covariance form is to be retained, the tedious number of inverses apparent in Eqs. (10.1) and (10.2) may be avoided as follows. Suppose we define the smoothed state as a linear fusion of the forward and backward filter states:

$$\hat{\mathbf{X}}_i^S = \mathbf{W}_{F_i} \hat{\mathbf{X}}_{F_i}^+ + \mathbf{W}_{B_i} \hat{\mathbf{X}}_{B_i}^- \quad (10.3)$$

For  $\hat{\mathbf{X}}_i^S$  to be unbiased, we must choose either  $\mathbf{W}_{F_i} = \mathbf{I} - \mathbf{W}_{B_i}$  or  $\mathbf{W}_{B_i} = \mathbf{I} - \mathbf{W}_{F_i}$ . Choosing the latter, the smoothed state becomes

$$\hat{\mathbf{X}}_i^S = \mathbf{W}_{F_i} \hat{\mathbf{X}}_{F_i}^+ + (\mathbf{I} - \mathbf{W}_{F_i}) \hat{\mathbf{X}}_{B_i}^- \quad (10.4)$$

Given the enforced lack of correlation between the forward and backward filters, the fused (smoothed) covariance is given by

$$\mathbf{P}_i^S = \mathbf{W}_{F_i} \mathbf{P}_{F_i}^+ \mathbf{W}_{F_i}^\top + \mathbf{W}_{B_i} \mathbf{P}_{B_i}^- \mathbf{W}_{B_i}^\top \quad (10.5)$$

$$= \mathbf{W}_{F_i} \mathbf{P}_{F_i}^+ \mathbf{W}_{F_i}^\top + (\mathbf{I} - \mathbf{W}_{F_i}) \mathbf{P}_{B_i}^- (\mathbf{I} - \mathbf{W}_{F_i})^\top \quad (10.6)$$

Choosing  $\mathbf{W}_{F_i}$  to minimize the trace of  $\mathbf{P}_i^S$  results in

$$\mathbf{W}_{F_i} = \mathbf{P}_{B_i}^- (\mathbf{P}_{F_i}^+ + \mathbf{P}_{B_i}^-)^{-1} \quad (10.7)$$

To see that Eq. (10.6), with Eq. (10.7), is equal to Eq. (10.2), expand Eq. (10.6), substituting Eq. (10.7), and recall Woodbury's identity:

$$\left[ (\mathbf{P}_{F_i}^+)^{-1} + (\mathbf{P}_{B_i}^-)^{-1} \right]^{-1} = \mathbf{P}_{B_i}^- - \mathbf{P}_{B_i}^- (\mathbf{P}_{F_i}^+ + \mathbf{P}_{B_i}^-)^{-1} \mathbf{P}_{B_i}^- \quad (10.8)$$

$$= \mathbf{P}_{F_i}^+ - \mathbf{P}_{F_i}^+ (\mathbf{P}_{F_i}^+ + \mathbf{P}_{B_i}^-)^{-1} \mathbf{P}_{F_i}^+ \quad (10.9)$$

To see that Eq. (10.4), with Eq. (10.7), is equal to Eq. (10.1), use Eq. (10.8) to show that  $\mathbf{P}_i^S (\mathbf{P}_{B_i}^-)^{-1} = \mathbf{W}_{B_i}$  and use Eq. (10.9) to show that  $\mathbf{P}_i^S (\mathbf{P}_{F_i}^+)^{-1} = \mathbf{W}_{F_i}$ .

In a typical application, the forward filter has been running continuously onboard the vehicle, and ground-based analysts will periodically wish to generate a BET over a particular

span of recently downlinked data. If the telemetry system has recorded and downlinked the full state and covariance at each measurement time, along with the measurements, the ground system need only run a “copy” of the forward filter backwards through the measurements and fuse the data according to Eqs. (10.1) and (10.2). Care must be taken that regeneration of the state transition matrices and process noise covariances is consistent with the forward filter’s modeling, and with the negative flow of time.

For various reasons, it may be necessary or desirable to run the forward filter on the ground as well, e.g. with higher-fidelity modeling than the onboard implementation permits. If so, it is efficient to store the state transition matrices and process noise covariances computed in the forward pass for use in the backward filter. In this case, Brown shows that the backward covariance may be propagated using

$$\mathbf{P}_{B_i}^- = \mathbf{\Phi}_{i+1,i}^{-1} \left[ \mathbf{P}_{B_{i+1}}^+ + \mathbf{S}_{i+1} \right] \mathbf{\Phi}_{i+1,i}^{-\top} \quad (10.10)$$



## Advanced Estimation Algorithms

This chapter will describe advanced estimation algorithms that have yet to achieve the status of best practices, but which appear to the contributors to have good potential for someday reaching such status.

### The Sigma-Point Estimator

Contributed by J. Russell Carpenter

Derivative-free state estimation techniques have received increasing attention in recent years. A particular class of such estimators make use of the columns of the factors of the estimators' error covariance matrices, which are scaled to form vectors that have become generally known as "sigma points." The so-called "Unscented Kalman Filter" [33, 34, 42] is a particular example of a sigma-point filter. A more general form is the divided-difference sigma-point filter, which is a sequential estimator that replaces first-order truncations of Taylor series approximations with second-order numerical differencing equations to approximate nonlinear dynamics and measurement models [53, 54]. If the process and measurement noise enter the system additively, Lee and Alfriend recently showed that several simplifications are possible, including a substantial reduction in the number of sigma-points [41].

This section highlights some broad aspects of sigma-point filtering, then briefly reviews how the ADDSPF works. It concludes with some brief comments comparing the ADDSPF to other sequential filters.

**The Sigma-Point Filter** In its most general form, the sigma-point filter performs sequential estimation of the  $n$ -dimensional state,  $\mathbf{x}$ , whose nonlinear dynamics over the time interval  $[t_k, t_{k+1}]$  are given by

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{w}_k) \quad (11.1)$$

The process noise input,  $\mathbf{w}$ , consists of independent increments whose first two moments are  $E[\mathbf{w}_k] = \mathbf{0}$  and  $E[\mathbf{w}_k \mathbf{w}_\ell] = \mathbf{Q}_k \delta_{k\ell}$ , where  $\delta_{k\ell}$  is the Kronecker delta. Although the second moment may be a function of the time index, this estimator assumes that all of the samples of  $\mathbf{w}$  arise from the same type of distribution, and this work further assumes that this distribution is Gaussian, so that higher-order moments may be neglected.

The filter sequentially processes an ordered set of measurements,  $\mathbb{Y}_k = [\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_k]$  of the form

$$\mathbf{y}_k = \mathbf{h}(\mathbf{x}_k, \mathbf{v}_k) \quad (11.2)$$

where the measurement noise input,  $\mathbf{v}$ , consists of independent and identically distributed (again, in this work, Gaussian) increments whose first two moments are  $E[\mathbf{v}_k] = \mathbf{0}$  and

$E[\mathbf{v}_k \mathbf{v}_\ell] = \mathbf{R}_k \delta_{k\ell}$ . By contrast, the ADDSPF utilizes models where the noise sources enter additively:

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k) + \mathbf{g}(\mathbf{x}_k) \mathbf{w}_k \quad (11.3)$$

and

$$\mathbf{y}_k = \mathbf{h}(\mathbf{x}_k) + \mathbf{v}_k. \quad (11.4)$$

All sigma-point filters utilize a linear measurement update equation of the form

$$\hat{\mathbf{x}}_k^+ = \hat{\mathbf{x}}_k^- + \mathbf{K}_k (\mathbf{y}_k - \hat{\mathbf{y}}_k^-) \quad (11.5)$$

where the accented variables in Eq. 11.5 denote conditional expectations, as in the Kalman filter:

$$\hat{\mathbf{x}}_k^+ = E[\mathbf{x}_k | \mathbb{Y}_k] \quad (11.6)$$

$$\hat{\mathbf{x}}_k^- = E[\mathbf{x}_k | \mathbb{Y}_{k-1}] \quad (11.7)$$

$$\hat{\mathbf{y}}_k^- = E[\mathbf{y}_k | \mathbb{Y}_{k-1}] \quad (11.8)$$

The gain matrix,  $\mathbf{K}$ , is based on conditional covariances, as in the Kalman filter:

$$\mathbf{K}_k = \mathbf{P}_{xy_k}^- (\mathbf{P}_{yy_k}^-)^{-1} \quad (11.9)$$

$$\mathbf{P}_k^- = \mathbf{P}_{xx_k}^- = E[(\mathbf{x}_k - \hat{\mathbf{x}}_k^-) (\mathbf{x}_k - \hat{\mathbf{x}}_k^-)^\top | \mathbb{Y}_{k-1}] \quad (11.10)$$

$$\mathbf{P}_{xy_k}^- = E[(\mathbf{x}_k - \hat{\mathbf{x}}_k^-) (\mathbf{y}_k - \hat{\mathbf{y}}_k^-)^\top | \mathbb{Y}_{k-1}] \quad (11.11)$$

$$\mathbf{P}_{yy_k}^- = E[(\mathbf{y}_k - \hat{\mathbf{y}}_k^-) (\mathbf{y}_k - \hat{\mathbf{y}}_k^-)^\top | \mathbb{Y}_{k-1}] \quad (11.12)$$

$$(11.13)$$

and the covariance associated with state estimate  $\hat{\mathbf{x}}_k^+$  is

$$\mathbf{P}_k^+ = \mathbf{P}_{xx_k}^+ = E[(\mathbf{x}_k - \hat{\mathbf{x}}_k^+) (\mathbf{x}_k - \hat{\mathbf{x}}_k^+)^\top | \mathbb{Y}_k] \quad (11.14)$$

Hereafter, equations will suppress the time index if it is the same for all variables in the equation.

Estimators such as the Kalman filter estimate these conditional expectations by approximating the nonlinear functions  $\mathbf{f}$  and  $\mathbf{h}$  with first-order Taylor series truncations, e.g.:

$$\mathbf{f}(\mathbf{x}) \approx \mathbf{f}(\hat{\mathbf{x}}^-) + \mathbf{f}'(\mathbf{x})(\mathbf{x} - \hat{\mathbf{x}}^-) \quad (11.15)$$

where  $\mathbf{f}'$  is an exact gradient. By contrast, the divided difference filter uses a second-order truncation along with numerical differencing formulas for the derivatives:

$$\mathbf{f}(\mathbf{x}) \approx \mathbf{f}(\hat{\mathbf{x}}^-) + \tilde{\mathbf{D}}_{\Delta \mathbf{x}}^{(1)} \mathbf{f}(\hat{\mathbf{x}}^-) + \tilde{\mathbf{D}}_{\Delta \mathbf{x}}^{(2)} \mathbf{f}(\hat{\mathbf{x}}^-) \quad (11.16)$$

where the divided difference operators,  $\tilde{\mathbf{D}}_{\Delta \mathbf{x}}^{(i)} \mathbf{f}(\hat{\mathbf{x}}^-)$ , approximate the coefficients of the multidimensional Taylor series expansion using Stirling interpolations. These interpolators difference perturbations of  $\mathbf{f}(\hat{\mathbf{x}}^-)$  across an interval,  $h$ , over a spanning basis set. Whether they are first-order, such as the unscented filter, or second order, sigma-point filters choose the interval so as to better approximate the moments required for the gain calculation, and choose as the spanning basis a set of sigma points, which are derived from  $\hat{\mathbf{x}}^-$  and the columns of the Cholesky factors of  $\mathbf{P}^-$  as follows.

**The ADDSPF** Let  $\hat{\mathcal{X}}$  denote the array whose columns are a particular ordering of the sigma points derived from  $\hat{\mathbf{x}}$  and its corresponding covariance,  $\mathbf{P}$ . Then

$$\hat{\mathcal{X}} = \left[ \hat{\mathbf{x}}, \hat{\mathbf{x}} + h\sqrt[3]{\mathbf{P}}_{(:,1)}, \hat{\mathbf{x}} + h\sqrt[3]{\mathbf{P}}_{(:,2)}, \dots, \hat{\mathbf{x}} - h\sqrt[3]{\mathbf{P}}_{(:,1)}, \hat{\mathbf{x}} - h\sqrt[3]{\mathbf{P}}_{(:,2)}, \dots \right] \quad (11.17)$$

where the subscript  $(:,i)$  denotes column  $i$  of the corresponding array, and  $\mathbf{P} = \sqrt[3]{\mathbf{P}}\sqrt[3]{\mathbf{P}}^\top$  denotes a Cholesky factorization. In the sequel, the shorthand notation  $\hat{\mathbf{x}} \pm h\sqrt[3]{\mathbf{P}}$  will denote the array on the right-hand side of the equation above. Then, for the ADDSPF, Ref. 41 shows that as each new measurement becomes available, an array of sigma points generated from the prior update should be propagated to the new measurement time:

$$\hat{\mathcal{X}}_k^- = \mathbf{f}(\hat{\mathcal{X}}_{k-1}^+) \quad (11.18)$$

These propagated sigma points are then merged to form the state estimate just prior to incorporating the new measurement as follows:

$$\hat{\mathbf{x}}^- = \mu_h(\hat{\mathcal{X}}^-) = \frac{h^2 - n}{h^2} \hat{\mathcal{X}}_{(:,1)}^- + \frac{1}{2h^2} \sum_{i=2}^{2n+1} \hat{\mathcal{X}}_{(:,i)}^- \quad (11.19)$$

To form an associated covariance, the following divided-differences are next computed:

$$\tilde{\mathbf{D}}_{\Delta\mathbf{x}}^{(1)}\mathbf{f}(\hat{\mathbf{x}}^-)_{(:,i)} = \frac{1}{2h} \left[ \hat{\mathcal{X}}_{(:,i+1)}^- - \hat{\mathcal{X}}_{(:,i+1+n)}^- \right] \quad (11.20)$$

$$\tilde{\mathbf{D}}_{\Delta\mathbf{x}}^{(2)}\mathbf{f}(\hat{\mathbf{x}}^-)_{(:,i)} = \frac{\sqrt{h^2 - 1}}{2h^2} \left[ \hat{\mathcal{X}}_{(:,i+1)}^- + \hat{\mathcal{X}}_{(:,i+1+n)}^- - 2\hat{\mathcal{X}}_{(:,1)}^- \right] \quad (11.21)$$

Ref. 41 shows that the covariance may then be computed from

$$\mathbf{P}^- = \left[ \tilde{\mathbf{D}}_{\Delta\mathbf{x}}^{(1)}\mathbf{f}(\hat{\mathbf{x}}^-), \tilde{\mathbf{D}}_{\Delta\mathbf{x}}^{(2)}\mathbf{f}(\hat{\mathbf{x}}^-), \sqrt[3]{\mathbf{Q}_d} \right] \left[ \tilde{\mathbf{D}}_{\Delta\mathbf{x}}^{(1)}\mathbf{f}(\hat{\mathbf{x}}^-), \tilde{\mathbf{D}}_{\Delta\mathbf{x}}^{(2)}\mathbf{f}(\hat{\mathbf{x}}^-), \sqrt[3]{\mathbf{Q}_d} \right]^\top \quad (11.22)$$

One advantage of sigma-point filters is that the full covariance need not be maintained, but rather only its Cholesky factor. Although the factors in square brackets in Eq. 11.22 are not Cholesky factors, since each is a full  $n \times 3n$  matrix, one may extract an  $n \times n$  triangular factor from it using the so-called ‘‘thin’’ version [23] of the QR decomposition<sup>1</sup>, or alternatively using a Householder factorization [2]. Thus,

$$\mathbf{M} \begin{bmatrix} \sqrt[3]{\mathbf{P}^-} \\ \mathbf{O}_{2n \times n} \end{bmatrix} = \left[ \tilde{\mathbf{D}}_{\Delta\mathbf{x}}^{(1)}\mathbf{f}(\hat{\mathbf{x}}^-), \tilde{\mathbf{D}}_{\Delta\mathbf{x}}^{(2)}\mathbf{f}(\hat{\mathbf{x}}^-), \sqrt[3]{\mathbf{Q}_d} \right]^\top \quad (11.23)$$

where  $\mathbf{M}$  is a full  $3n \times 3n$  orthonormal matrix, and  $\mathbf{O}_{2n \times n}$  is a  $2n \times n$  matrix of zeros.

For the measurement update, a new array of sigma points must be generated from  $\hat{\mathbf{x}}^-$  and  $\mathbf{P}^-$ ; this array is denoted  $\hat{\mathcal{X}}^*$ . These sigma points are used to generate a set of sigma points representing the measurement:

$$\hat{\mathcal{Y}}^- = \mathbf{h}(\hat{\mathcal{X}}^*) \quad (11.24)$$

In similar fashion to the time update, the sigma points of the measurement are then merged to form the estimated measurement:

$$\hat{\mathbf{y}}^- = \mu_h(\hat{\mathcal{Y}}^-) \quad (11.25)$$

<sup>1</sup>For *Matlab* users, this may be accomplished in several ways, e.g. by passing the transpose of this matrix to the `qr` function, then keeping the first  $n$  non-zero rows from the second output, and transposing this result.

the corresponding divided-differences are computed as:

$$\tilde{D}_{\Delta\mathbf{x}}^{(1)}\mathbf{h}(\hat{\mathbf{x}}^-)_{(:,i)} = \frac{1}{2h} \left[ \hat{\mathcal{Y}}_{(:,i+1)}^- - \hat{\mathcal{Y}}_{(:,i+1+n)}^- \right] \quad (11.26)$$

$$\tilde{D}_{\Delta\mathbf{x}}^{(2)}\mathbf{h}(\hat{\mathbf{x}}^-)_{(:,i)} = \frac{\sqrt{h^2-1}}{2h^2} \left[ \hat{\mathcal{Y}}_{(:,i+1)}^- + \hat{\mathcal{Y}}_{(:,i+1+n)}^- - 2\hat{\mathcal{Y}}_{(:,1)}^- \right] \quad (11.27)$$

and the covariances required for the gain calculation may then be computed from

$$\mathbf{P}_{yy}^- = \left[ \tilde{D}_{\Delta\mathbf{x}}^{(1)}\mathbf{h}(\hat{\mathbf{x}}^-), \tilde{D}_{\Delta\mathbf{x}}^{(2)}\mathbf{h}(\hat{\mathbf{x}}^-), \sqrt{{}^c\mathbf{R}} \right] \left[ \tilde{D}_{\Delta\mathbf{x}}^{(1)}\mathbf{h}(\hat{\mathbf{x}}^-), \tilde{D}_{\Delta\mathbf{x}}^{(2)}\mathbf{h}(\hat{\mathbf{x}}^-), \sqrt{{}^c\mathbf{R}} \right]^\top \quad (11.28)$$

$$\mathbf{P}_{xy}^- = \sqrt{{}^c\mathbf{P}^-} \left[ \tilde{D}_{\Delta\mathbf{x}}^{(1)}\mathbf{h}(\hat{\mathbf{x}}^-) \right]^\top \quad (11.29)$$

Note that the second-order divided difference for the measurement function,  $\tilde{D}_{\Delta\mathbf{x}}^{(2)}\mathbf{h}(\hat{\mathbf{x}}^-)$ , does *not* appear in the cross-covariance update. As with the time update, through the use of the thin QR factorization, only triangular factors need be maintained for  $\mathbf{P}_{yy}^-$ <sup>2</sup>:

$$\mathbf{M}_{yy} \begin{bmatrix} \sqrt{{}^c\mathbf{P}_{yy}^-} \\ \mathbf{O}_{2n \times n} \end{bmatrix}^\top = \left[ \tilde{D}_{\Delta\mathbf{x}}^{(1)}\mathbf{h}(\hat{\mathbf{x}}^-), \tilde{D}_{\Delta\mathbf{x}}^{(2)}\mathbf{h}(\hat{\mathbf{x}}^-), \sqrt{{}^c\mathbf{R}} \right]^\top \quad (11.30)$$

Now, all of the terms required for the state update (Eqs. 11.5 and 11.9), are available. Ref. **41** shows that the corresponding Cholesky factor of the covariance is extracted from

$$\mathbf{M}^+ \begin{bmatrix} \sqrt{{}^c\mathbf{P}^+} \\ \mathbf{O}_{2n \times n} \end{bmatrix}^\top = \left[ \sqrt{{}^c\mathbf{P}^-} - \mathbf{K}\tilde{D}_{\Delta\mathbf{x}}^{(1)}\mathbf{h}(\hat{\mathbf{x}}^-), \mathbf{K} \left[ \tilde{D}_{\Delta\mathbf{x}}^{(2)}\mathbf{h}(\hat{\mathbf{x}}^-), \sqrt{{}^c\mathbf{R}} \right] \right] \quad (11.31)$$

**The ADDSPF vs. Other Sequential Filters** To conclude this section, some observations concerning the ADDSPF in comparison to other filters are offered. These observations concern the number of sigma points, the order of approximation, and the existence and method of choice of free parameters in the algorithms.

Although in many problems of practical interest the noise enters the system additively, if this is not the case, then either the original divided difference filter or the unscented filter may provide superior results to the ADDSPF, at the cost of requiring more sigma points. In both of the former algorithms, the nonlinear functions must be perturbed not only over a basis spanning the state space, but also over the discrete process noise and measurement noise spaces. Thus, rather than  $2n + 1$  sigma points, the more general algorithms require  $2n_a + 1$ , where  $n_a = n + n_w + n_v$ , and  $n_w$  and  $n_v$  are the dimensions of the discrete process noise and measurement noise inputs.

The Kalman filter is an exact algorithm for linear stochastic systems driven by Gaussian noise, and nothing is to be gained from the use the sigma-point filters for such purely linear systems. First-order sigma-point filters such as the UKF retain the Kalman filter's first-order truncation, but avoid the need for the designer to supply explicit gradients. The divided difference filter is comparable to a derivative-free version of the modified second-order Gaussian filter [30] in that, for symmetric distributions, it retains some terms as high as order four.

Unlike the Kalman filter, for which all of the parameters in principle can be associated with properties of the underlying stochastic system, all of the sigma point filters involve at

<sup>2</sup>Although it might seem that the full matrix  $\mathbf{P}_{yy}^-$  is required for the gain computation of Eq. 11.9, Ref. **53** points out that, rather than inverting the product of the factors to compute the gain, the gain may be solved from forward and back substitution directly using the Cholesky factor.

least one free parameter. In the unscented filter, the weights for combining the sigma points involve three parameters whose physical interpretation is perhaps less clear than with the single parameter in the divided difference algorithms, where the free parameter  $h$  is clearly associated with the size of the perturbation in the numerical differencing formulae. Ref. **53** shows that  $h$  should be bounded below by  $h > 1$ , and that for symmetric distributions,  $\sqrt{h}$  should be equal to the kurtosis, which for a Gaussian distribution is three<sup>3</sup>.

---

<sup>3</sup>Some authors subtract three from the definition of kurtosis, so that Gaussian distributions have zero kurtosis.



## APPENDIX A

# Models and Realizations of Random Variables

Contributed by J. Russell Carpenter

A continuous random variable is a function that maps the outcomes of random events to the real line. Realizations of random variables are thus real numbers. A vector of  $n$  random variables maps outcomes of random events to  $\mathbb{R}^n$ . For our purposes, random variables will always be associated with a probability density function that indicates the likelihood that a realization occurs within a particular interval of the real line, or within a particular subspace of  $\mathbb{R}^n$  for the vector case. It is common to assume that this density is the normal or Gaussian density. For the vector case, the normal probability density function is

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{1}{\sqrt{|2\pi\mathbf{P}|}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\top}\mathbf{P}^{-1}(\mathbf{x}-\boldsymbol{\mu})} \quad (\text{A.1})$$

where  $\boldsymbol{\mu}$  is a vector of mean values for each component of  $\mathbf{x}$ , and  $\mathbf{P}$  is a matrix that contains the variances of each component of  $\mathbf{x}$  along its diagonal, and the covariances between each component as its off-diagonal components. The covariances indicate the degree of correlation between the random variables composing  $\mathbf{x}$ . The matrix  $\mathbf{P}$  is thus called the variance-covariance matrix, which we will hereafter abbreviate to just “covariance matrix,” or “covariance.” Since the normal density is completely characterized by its mean and covariance, we will use the following notation as a shorthand to describe drawing a realization from a normally-distributed random vector:

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \mathbf{P}) \quad (\text{A.2})$$

Thus, the model for realizations of a measurement noise vector is

$$\mathbf{v} \sim N(0, \mathbf{R}) \quad (\text{A.3})$$

For the scalar case, or for the vector case when the covariance is diagonal, we may directly generate realizations of a normally-distributed random vector from normal random number generators available in most software libraries. If  $\mathbf{P}$  has non-zero off-diagonal elements, we must model the specified correlations when we generate realizations. If  $\mathbf{P}$  is strictly positive definite, we can factor it as follows:

$$\mathbf{P} = \mathcal{C}\sqrt{\mathbf{P}}\mathcal{C}^{\top} \quad (\text{A.4})$$

where  $\mathcal{C}\sqrt{\mathbf{P}}$  is a triangular matrix known as a Cholesky factor; this can be viewed as a “matrix square root.” The Cholesky factorization is available in many linear algebra libraries. We can then use  $\mathcal{C}\sqrt{\mathbf{P}}$  to generate correlated realizations of  $\mathbf{x}$  as follows. Let  $\mathbf{z}$  be a realization of a normally-distributed random vector of the same dimension as  $\mathbf{x}$ , with zero mean and

unit variance, that is

$$\mathbf{z} \sim N(\mathbf{0}, \mathbf{I}) \quad (\text{A.5})$$

Then, with

$$\mathbf{x} = \sqrt{\mathbf{P}}\mathbf{z} \quad (\text{A.6})$$

we can generate properly correlated realizations of  $\mathbf{x}$ . We can also use a Cholesky factorization of the measurement noise covariance  $\mathbf{R}$ , if  $\mathbf{R}$  is non-diagonal, to transform correlated measurements into uncorrelated auxiliary measurements for cases in which the estimator cannot handle correlated measurement data.

If  $\mathbf{P}$  is only non-negative definite, i.e.  $\mathbf{P} \geq 0$  rather than  $\mathbf{P} > 0$  as above, the Cholesky factorization does not exist. In this case, since  $\mathbf{P}$ 's eigenvalues are real and distinct, it has a diagonal factorization:

$$\mathbf{P} = \mathbf{V}\mathbf{D}\mathbf{V}^T \quad (\text{A.7})$$

where  $\mathbf{V}$  is a matrix of eigenvectors and  $\mathbf{D}$  is a diagonal matrix of eigenvalues. Then, with  $\mathbf{z}$  as above,

$$\mathbf{x} = \mathbf{V}\sqrt{\mathbf{D}}\mathbf{z} \quad (\text{A.8})$$

where  $\sqrt{\mathbf{D}}$  implies taking the square roots of each diagonal element.

## APPENDIX B

# The Mathematics Behind the UDU Factorization

Contributed by Chris D'Souza

### B.1. The Partitioning into Two Subproblems

We can find that the update equation is

$$\overline{\mathbf{UDU}}^T = \mathbf{\Phi UDU}^T \mathbf{\Phi}^T + \mathbf{Q} \quad (\text{B.1})$$

$$\begin{aligned} &= \mathbf{\Phi}_2 \mathbf{\Phi}_1 \mathbf{UDU}^T \mathbf{\Phi}_1^T \mathbf{\Phi}_2^T + \mathbf{Q}_1 + \mathbf{Q}_2 \\ &= \mathbf{\Phi}_2 [\mathbf{\Phi}_1 \mathbf{UDU}^T \mathbf{\Phi}_1^T] \mathbf{\Phi}_2^T + \mathbf{Q}_1 + \mathbf{Q}_2 \end{aligned} \quad (\text{B.2})$$

Recalling that  $\mathbf{Q}_1 = \mathbf{\Phi}_2 \mathbf{\Phi}_2^{-1} \mathbf{Q}_1 \mathbf{\Phi}_2^{-T} \mathbf{\Phi}_2^T$  and

$$\mathbf{\Phi}_2^{-1} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{M}^{-1} \end{bmatrix} \quad (\text{B.3})$$

where  $\mathbf{M}^{-1} = \text{diag}(1/m_i)$ ,  $i = 1, 2, 3, \dots, n_p$ . We note that

$$\mathbf{\Phi}_2^{-1} \mathbf{Q}_1 \mathbf{\Phi}_2^{-T} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{M}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{Q}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{M}^{-1} \end{bmatrix} = \begin{bmatrix} \mathbf{Q}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} = \mathbf{Q}_1 \quad (\text{B.4})$$

### B.2. The Mathematics Behind the Second Subproblem

Recall that we partitioned  $\tilde{\mathbf{U}}$  and  $\tilde{\mathbf{D}}$  as

$$\tilde{\mathbf{U}} = \begin{bmatrix} \tilde{\mathbf{U}}_{aa} & \tilde{\mathbf{U}}_{ab} & \tilde{\mathbf{U}}_{ac} \\ \mathbf{0} & 1 & \tilde{\mathbf{U}}_{bc} \\ \mathbf{0} & \mathbf{0} & \tilde{\mathbf{U}}_{cc} \end{bmatrix} \begin{array}{l} \} n_a \\ \} 1 \\ \} n_c \end{array} \quad \text{and} \quad \tilde{\mathbf{D}} = \begin{bmatrix} \tilde{\mathbf{D}}_{aa} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \tilde{d}_b & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \tilde{\mathbf{D}}_{cc} \end{bmatrix} \begin{array}{l} \} n_a \\ \} 1 \\ \} n_c \end{array} \quad (\text{B.5})$$

in order to isolate a parameter. In fact, the state we choose to isolate is one of the Gauss-Markov states (likely associated with a sensor). Let

$$\mathbf{\Phi}_2 = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{M}_c \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & m_b & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I} \end{bmatrix} = \mathbf{\Phi}_c \mathbf{\Phi}_b \quad (\text{B.6})$$

and

$$\mathbf{Q}_2 = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & q_b & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{Q}_c \end{bmatrix} = \mathbf{Q}_b + \mathbf{Q}_c \quad (\text{B.7})$$

As in the previous exercise, we note that  $\mathbf{\Phi}_c^{-1} \mathbf{Q}_b \mathbf{\Phi}_c^{-T} = \mathbf{Q}_b$ . So, now Eq. (7.24) becomes

$$\overline{\mathbf{UD}} \overline{\mathbf{U}}^T = \mathbf{\Phi}_c \left[ \mathbf{\Phi}_b \tilde{\mathbf{U}} \tilde{\mathbf{D}} \tilde{\mathbf{U}}^T \mathbf{\Phi}_b^T + \mathbf{Q}_b \right] \mathbf{\Phi}_c^T + \mathbf{Q}_c \quad (\text{B.8})$$

The term in the square bracket in Eq. (B.8) is

$$\check{\mathbf{U}}\check{\mathbf{D}}\check{\mathbf{U}}^\top = \Phi_b \check{\mathbf{U}}\check{\mathbf{D}}\check{\mathbf{U}}^\top \Phi_b^\top + \mathbf{Q}_b \quad (\text{B.9})$$

The left side of Eq.(B.9) (recalling that  $\check{\mathbf{U}}_{bb} = 1$ ) is

$$\check{\mathbf{U}}\check{\mathbf{D}}\check{\mathbf{U}}^\top = \begin{bmatrix} \check{\mathbf{U}}_{aa}\check{\mathbf{D}}_{aa}\check{\mathbf{U}}_{aa}^\top & \check{\mathbf{U}}_{ac}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{ac}^\top & \check{\mathbf{U}}_{ac}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{cc}^\top \\ +\check{\mathbf{U}}_{ab}\check{d}_b\check{\mathbf{U}}_{ab}^\top & +\check{\mathbf{U}}_{ab}\check{d}_b & \\ \check{\mathbf{U}}_{ac}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{ac}^\top & & \\ \hline \check{d}_b\check{\mathbf{U}}_{ab}^\top & \check{d}_b + \check{\mathbf{U}}_{bc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{bc}^\top & \check{\mathbf{U}}_{bc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{cc}^\top \\ \check{\mathbf{U}}_{bc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{ac}^\top & & \\ \hline \check{\mathbf{U}}_{cc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{ac}^\top & \check{\mathbf{U}}_{cc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{bc}^\top & \check{\mathbf{U}}_{cc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{cc}^\top \end{bmatrix} \quad (\text{B.10})$$

The right side of Eq.(B.9), once again recalling that  $\check{\mathbf{U}}_{bb} = 1$ , is

$$\Phi_b \check{\mathbf{U}}\check{\mathbf{D}}\check{\mathbf{U}}^\top \Phi_b^\top + \mathbf{Q}_b = \begin{bmatrix} \check{\mathbf{U}}_{aa}\check{\mathbf{D}}_{aa}\check{\mathbf{U}}_{aa}^\top & m_b\check{\mathbf{U}}_{ac}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{ac}^\top & \check{\mathbf{U}}_{ac}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{cc}^\top \\ +\check{\mathbf{U}}_{ab}\check{d}_b\check{\mathbf{U}}_{ab}^\top & +m_b\check{\mathbf{U}}_{ab}\check{d}_b & \\ \check{\mathbf{U}}_{ac}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{ac}^\top & & \\ \hline m_b\check{d}_b\check{\mathbf{U}}_{ab}^\top & m_b^2\check{\mathbf{U}}_{bc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{bc}^\top & m_b\check{\mathbf{U}}_{bc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{cc}^\top \\ +m_b\check{\mathbf{U}}_{bc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{ac}^\top & +m_b^2\check{d}_b + q_b & \\ \hline \check{\mathbf{U}}_{cc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{ac}^\top & m_b\check{\mathbf{U}}_{cc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{bc}^\top & \check{\mathbf{U}}_{cc}\check{\mathbf{D}}_{cc}\check{\mathbf{U}}_{cc}^\top \end{bmatrix} \quad (\text{B.11})$$

Equating the components in Eqs. ((B.10) and (B.11), from the (1,3) and (3,3) element, we find

$$\check{\mathbf{U}}_{ac} = \check{\mathbf{U}}_{ac}, \quad \check{\mathbf{D}}_{cc} = \check{\mathbf{D}}_{cc}, \quad \check{\mathbf{U}}_{cc}^\top = \check{\mathbf{U}}_{cc}^\top \quad (\text{B.12})$$

From the (2,3) element we get,

$$\check{\mathbf{U}}_{bc} = m_b\check{\mathbf{U}}_{bc} \quad (\text{B.13})$$

From the (2,2) element and using the results of Eq. (B.13), we find that

$$\check{d}_b = m_b^2\check{d}_b + q_b \quad (\text{B.14})$$

The (2,1) element yields

$$\check{\mathbf{U}}_{ab} = m_b \frac{\check{d}_b}{\check{d}_b} \check{\mathbf{U}}_{ab} \quad (\text{B.15})$$

What finally remains is the (1,1) element and it is on this we focus. Using the relations in the previous equations, we find that

$$\check{\mathbf{U}}_{aa}\check{\mathbf{D}}_{aa}\check{\mathbf{U}}_{aa}^\top = \check{\mathbf{U}}_{aa}\check{\mathbf{D}}_{aa}\check{\mathbf{U}}_{aa}^\top + \left[ \check{d}_b - m_b^2 \frac{\check{d}_b}{\check{d}_b} \right] \check{\mathbf{U}}_{ab}\check{\mathbf{U}}_{ab}^\top \quad (\text{B.16})$$

The term in the bracket can be simplified as

$$\left[ \check{d}_b - m_b^2 \frac{\check{d}_b}{\check{d}_b} \right] = \frac{m_b^2\check{d}_b^2 + q_b\check{d}_b - m_b^2\check{d}_b^2}{m_b^2\check{d}_b + q_b} = \frac{\check{d}_b q_b}{m_b^2\check{d}_b + q_b} = \frac{\check{d}_b q_b}{\check{d}_b} \quad (\text{B.17})$$

so Eq.(B.16) becomes

$$\check{\mathbf{U}}_{\mathbf{aa}}\check{\mathbf{D}}_{\mathbf{aa}}\check{\mathbf{U}}_{\mathbf{aa}}^T = \tilde{\mathbf{U}}_{\mathbf{aa}}\tilde{\mathbf{D}}_{\mathbf{aa}}\tilde{\mathbf{U}}_{\mathbf{aa}}^T + \left(\frac{\tilde{d}_b q_b}{\check{d}_b}\right) \tilde{\mathbf{U}}_{ab}\tilde{\mathbf{U}}_{ab}^T \quad (\text{B.18})$$

We note that  $\tilde{\mathbf{U}}_{ab}$  is a column vector so Eq.(B.18), and hence is of rank 1, constitutes a ‘rank one’ update. Since  $\check{d}_b$ ,  $\tilde{d}_b$  and  $q_b$  are all positive (assuming  $m_b$  is a positive quantity), we can use the Agee-Turner Rank One update [1]. It should be pointed out that as the algorithm proceeds down the ‘list’ of parameters, the size of the states  $\mathbf{a}$  increases by one (and consequently the size of the parameters  $\mathbf{c}$  decreases by one. Hence  $\check{\mathbf{U}}_{\mathbf{aa}}$  and  $\check{\mathbf{D}}_{\mathbf{aa}}$  begins with a dimension of  $n_{\mathbf{x}}$  and concludes with dimension  $n_{\mathbf{x}} + n_{\mathbf{p}} - 1$ .

Therefore, this is done recursively for all the (sensor) parameters  $\mathbf{p}$  which are of size  $n_{\mathbf{p}}$ .

### B.3. The Agee-Turner Rank-One Update

In trying to get an efficient algorithm for performing the time update of the covariance matrix, we were faced with Eq.(B.18), which is of the form

$$\tilde{\mathbf{U}}\tilde{\mathbf{D}}\tilde{\mathbf{U}}^T = \mathbf{U}\mathbf{D}\mathbf{U}^T + \mathbf{c}\mathbf{x}\mathbf{x}^T \quad (\text{B.19})$$

This is called a ‘rank one’ update because we are updating the matrix factors  $\mathbf{U}$  and  $\mathbf{D}$  based upon products of  $\mathbf{x}$  which is of rank 1.

In order to reduce the number of mathematical operations (adds/subtracts, multiplies and divides), for the case of parameter or ECRV/First-order Gauss Markov processes, for sensor parameters, we consider the rank-one update first introduced by Agee and Turner (of the White Sands Missile Range (WSMR)) in 1972.

Consider a covariance matrix update of the form,

$$\tilde{\mathbf{P}} = \mathbf{P} + \mathbf{c}\mathbf{x}\mathbf{x}^T \quad (\text{B.20})$$

or

$$\tilde{\mathbf{U}}\tilde{\mathbf{D}}\tilde{\mathbf{U}}^T = \mathbf{U}\mathbf{D}\mathbf{U}^T + \mathbf{c}\mathbf{x}\mathbf{x}^T \quad (\text{B.21})$$

so,  $\tilde{p}_{ij}$  can be expressed (and defined) as

$$\tilde{p}_{ij} \triangleq \sum_{k=j}^n \tilde{u}_{ik}\tilde{d}_{kk}\tilde{u}_{jk} = \sum_{k=j}^n u_{ik}d_{kk}u_{jk} + cx_i x_j \quad (\text{B.22})$$

and

$$\tilde{p}_{ii} = \sum_{k=i}^n \tilde{u}_{ik}^2 \tilde{d}_{kk} = \sum_{k=i}^n u_{ik}^2 d_{kk} + cx_i^2 \quad (\text{B.23})$$

We recall that  $\tilde{u}_{ii} = u_{ii} = 1$  and thus, for an  $n \times n$  matrix, for  $j = n$  (i.e. the last column),

$$\tilde{d}_{nn} = d_{nn} + cx_n^2 \quad (\text{B.24})$$

$$\tilde{p}_{in} = \tilde{u}_{in}\tilde{d}_{nn}\tilde{u}_{nn} = d_{nn}u_{in}u_{nn} + cx_i x_n \quad (\text{B.25})$$

so that

$$\tilde{u}_{in} = \frac{1}{\tilde{d}_{nn}} (d_{nn}u_{in} + cx_i x_n) \quad (\text{B.26})$$

The second-to-the-last ( $n - 1$ -th) column of  $\mathbf{U}$  can now be can be operated on, by means of the following decomposition of Eq. (B.22), as

$$\sum_{k=j}^{n-1} \tilde{u}_{ik} \tilde{d}_{kk} \tilde{u}_{jk} + \tilde{u}_{in} \tilde{d}_{nn} \tilde{u}_{jn} = \sum_{k=j}^{n-1} u_{ik} d_{kk} u_{jk} + u_{in} d_{nn} u_{jn} + c x_i x_j \quad (\text{B.27})$$

which leads to

$$\sum_{k=j}^{n-1} \tilde{u}_{ik} \tilde{d}_{kk} \tilde{u}_{jk} = \sum_{k=j}^{n-1} u_{ik} d_{kk} u_{jk} + \Upsilon_n \quad (\text{B.28})$$

If we work on the terms outside the two summations, using Eq. (B.26) for  $\tilde{u}_{in}$  and  $\tilde{u}_{jn}$ ,  $\Upsilon_n$  becomes

$$\begin{aligned} \Upsilon_n &= -\tilde{u}_{in} \tilde{d}_{nn} \tilde{u}_{jn} + u_{in} d_{nn} u_{jn} + c x_i x_j \\ &= -\frac{1}{\tilde{d}_{nn}} [d_{nn} u_{in} + c x_i x_n] [d_{nn} u_{jn} + c x_j x_n] \\ &\quad + \frac{d_{nn} + c x_n^2}{\tilde{d}_{nn}} (u_{in} d_{nn} u_{jn} + c x_i x_j) \\ &= \frac{1}{\tilde{d}_{nn}} [-c d_{nn} u_{jn} x_n x_i - c d_{nn} u_{in} x_n x_j + c d_{nn} x_i x_j + c d_{nn} x_n^2 u_{in} u_{jn}] \\ &= \frac{c d_{nn}}{\tilde{d}_{nn}} (x_i - u_{in} x_n) (x_j - u_{jn} x_n) \end{aligned} \quad (\text{B.29})$$

Therefore, Eq. (B.27) can be written as

$$\sum_{k=j}^{n-1} \tilde{u}_{ik} \tilde{d}_{kk} \tilde{u}_{jk} = \sum_{k=j}^{n-1} u_{ik} d_{kk} u_{jk} + \frac{c d_{nn}}{\tilde{d}_{nn}} (x_i - u_{in} x_n) (x_j - u_{jn} x_n) \quad (\text{B.30})$$

Now, if we operate a bit more on the quantity  $\tilde{u}_{in}$ , we find from Eq. (B.26), that we get

$$\tilde{u}_{in} = \frac{d_{nn}}{\tilde{d}_{nn}} u_{in} + \frac{c}{\tilde{d}_{nn}} x_i x_n \quad (\text{B.31})$$

$$= \frac{\tilde{d}_{nn} - c x_n^2}{\tilde{d}_{nn}} u_{in} + \frac{c}{\tilde{d}_{nn}} x_i x_n \quad (\text{B.32})$$

$$= u_{in} + (x_i - u_{in} x_n) \frac{c x_n}{\tilde{d}_{nn}} \quad (\text{B.33})$$

and if we define  $\alpha_i$  and  $v_n$  as

$$\alpha_i \triangleq (x_i - u_{in} x_n) \quad (\text{B.34})$$

$$v_n \triangleq \frac{c x_n}{\tilde{d}_{nn}} \quad (\text{B.35})$$

$\tilde{u}_{in}$  can be rewritten as

$$\tilde{u}_{in} = u_{in} + \alpha_i v_n \quad (\text{B.36})$$

If we want to generalize this, we can write Eq. (B.30) as

$$\sum_{k=j}^{n-1} \tilde{u}_{ik} \tilde{d}_{kk} \tilde{u}_{jk} = \sum_{k=j}^{n-1} u_{ik} d_{kk} u_{jk} + C^n \mathcal{X}_i^n \mathcal{X}_j^n \quad (\text{B.37})$$

where

$$\mathcal{C}^n \triangleq \frac{c d_{nn}}{\tilde{d}_{nn}} \quad (\text{B.38})$$

$$\mathcal{X}_i^n \triangleq x_i - u_{in} x_n \quad (\text{B.39})$$

with

$$\begin{aligned} \tilde{u}_{in} &= u_{in} + \alpha_i^n v_n \\ \alpha_i^n &= x_i - u_{in} x_n \\ v_n &= \frac{c x_n}{\tilde{d}_{nn}} \\ \tilde{d}_{nn} &= d_{nn} + c x_n^2 \end{aligned}$$

Thus, for the third-to-the-last column ( $j = n - 2$ ), we expand Eq.(B.37) as

$$\begin{aligned} \sum_{k=j}^{n-2} \tilde{u}_{ik} \tilde{d}_{kk} \tilde{u}_{jk} + \tilde{u}_{i,n-1} \tilde{d}_{n-1,n-1} \tilde{u}_{j,n-1} &= \sum_{k=j}^{n-2} u_{ik} d_{kk} u_{jk} \\ &\quad + u_{i,n-1} d_{n-1,n-1} u_{j,n-1} \\ &\quad + \mathcal{C}^n \mathcal{X}_i^n \mathcal{X}_j^n \end{aligned} \quad (\text{B.40})$$

which produces

$$\begin{aligned} \sum_{k=j}^{n-2} \tilde{u}_{ik} \tilde{d}_{kk} \tilde{u}_{jk} &= \sum_{k=j}^{n-2} u_{ik} d_{kk} u_{jk} \\ &\quad + \frac{\mathcal{C}^n d_{n-1,n-1}}{\tilde{d}_{n-1,n-1}} [\mathcal{X}_i^n - u_{i,n-1} \mathcal{X}_n^n] [\mathcal{X}_j^n - u_{j,n-1} \mathcal{X}_n^n] \end{aligned} \quad (\text{B.41})$$

so that using the same machinery as above, we get

$$\sum_{k=j}^{n-2} \tilde{u}_{ik} \tilde{d}_{kk} \tilde{u}_{jk} = \sum_{k=j}^{n-2} u_{ik} d_{kk} u_{jk} + \mathcal{C}^{n-1} \mathcal{X}_i^{n-1} \mathcal{X}_j^{n-1} \quad (\text{B.42})$$

$$\mathcal{C}^{n-1} = \frac{\mathcal{C}^n d_{n-1,n-1}}{\tilde{d}_{n-1,n-1}} \quad (\text{B.43})$$

$$\mathcal{X}_i^{n-1} = \alpha_i^{n-1} = \mathcal{X}_i^n - u_{i,n-1} \mathcal{X}_n^n \quad (\text{B.44})$$

$$v_{n-1} = \frac{\mathcal{C}^n \mathcal{X}_i^n}{\tilde{d}_{n-1,n-1}} \quad (\text{B.45})$$

$$\tilde{u}_{i,n-1} = u_{i,n-1} + \mathcal{X}_i^{n-1} v_{n-1} \quad (\text{B.46})$$

$$\tilde{d}_{n-1,n-1} = d_{n-1,n-1} + \mathcal{C}^n x_{n-1}^2 \quad (\text{B.47})$$

We also are reminded that  $\tilde{u}_{i,i} = 1$ .

#### B.4. Decorrelating Measurements

We normalize the (original) measurement equation

$$\mathbf{z}_{orig} = \mathbf{H}_{orig}\mathbf{x} + \boldsymbol{\nu}_{orig} \quad (\text{B.48})$$

where the measurement noise has statistics

$$E[\boldsymbol{\nu}_{orig}] = \mathbf{0} \quad (\text{B.49})$$

$$E[\boldsymbol{\nu}_{orig}\boldsymbol{\nu}_{orig}^T] = \mathbf{R}_{orig} \quad (\text{B.50})$$

where  $\mathbf{R}_{orig}$  is the measurement noise.

We now change variables so that

$$\mathbf{z} \triangleq \mathbf{R}_{orig}^{-1/2} \mathbf{z}_{orig} \quad (\text{B.51})$$

$$\mathbf{H} \triangleq \mathbf{R}_{orig}^{-1/2} \mathbf{H}_{orig} \quad (\text{B.52})$$

$$\boldsymbol{\nu} \triangleq \mathbf{R}_{orig}^{-1/2} \boldsymbol{\nu}_{orig} \quad (\text{B.53})$$

where  $\mathbf{R}_{orig}^{-1/2}$  is the inverse of the Cholesky factor of  $\mathbf{R}_{orig}$  ( $= \mathbf{R}_{orig}^{1/2} \mathbf{R}_{orig}^{T/2}$ ). With this, the new normalized measurement equation is

$$\mathbf{z} = \mathbf{H}\mathbf{x} + \boldsymbol{\nu} \quad (\text{B.54})$$

where  $E[\boldsymbol{\nu}] = \mathbf{0}$  and  $E[\boldsymbol{\nu}\boldsymbol{\nu}^T] = \mathbf{I}$ . This is sometimes referred to as **pre-whitening** or *decorrelation*, because if the original measurements were correlated, the normalized measurements are now uncorrelated (via the Cholesky decomposition of  $\mathbf{R}_{orig}$ ).

#### B.5. The Carlson Rank-One Update

The Carlson rank-one update [5], introduced in 1973, addresses the problem of updating the covariance due to a loss of precision involved in the differencing of two positive quantities which are nearly equal. In particular, the diagonal elements  $d_{ii}$  have the potential of going negative in certain cases if the Agee-Turner rank-one update is blindly used. Thankfully, we resort to the *Carlson rank-one update* to compute the measurement update without losing numerical precision.

We recall that  $\alpha$  and  $\mathbf{v}$  were defined earlier. We also define the  $n \times 1$  vector  $\mathbf{f}$  as

$$\mathbf{f} \triangleq \bar{\mathbf{U}}^T \mathbf{H}^T \quad (\text{B.55})$$

Therefore, since  $\mathbf{f}$  is an  $n \times 1$  vector and  $\mathbf{D}$  is a diagonal matrix, we can express  $\alpha$  as

$$\alpha = \alpha_n = R + \sum_{i=1}^n f_i^2 d_{ii} \quad (\text{B.56})$$

so that

$$\alpha_j = R + \sum_{i=1}^j f_i^2 d_{ii} = \alpha_{j-1} + f_j^2 d_{jj} \quad (\text{B.57})$$

Since  $\mathbf{v}$  can be written as

$$\mathbf{v} = \mathbf{D}\mathbf{f} \quad (\text{B.58})$$

we can also write

$$v_j = d_{jj} f_j \quad (\text{B.59})$$

and we can write  $\alpha_j$  as

$$\alpha_j = \alpha_{j-1} + \frac{v_j^2}{d_{jj}} \quad (\text{B.60})$$

$$\tilde{\mathbf{U}}\tilde{\mathbf{D}}\tilde{\mathbf{U}}^T \triangleq \bar{\mathbf{D}} - \frac{1}{\alpha}\mathbf{v}\mathbf{v}^T \quad (\text{B.61})$$

which can be rewritten as:

$$\tilde{\mathbf{U}}\tilde{\mathbf{D}}\tilde{\mathbf{U}}^T \triangleq \bar{\mathbf{U}}\bar{\mathbf{D}}\bar{\mathbf{U}}^T - \frac{1}{\alpha}\mathbf{v}\mathbf{v}^T \quad (\text{B.62})$$

where  $\bar{\mathbf{U}} = \mathbf{I}$ . Following the reasoning in the description for the Rank-One update earlier in this Appendix,

$$\tilde{p}_{ij} \triangleq \sum_{k=j}^n \tilde{u}_{ik}\tilde{d}_{kk}\tilde{u}_{jk} = -\frac{1}{\alpha}v_iv_j \quad (\text{B.63})$$

and

$$\tilde{p}_{ii} = \sum_{k=i}^n \tilde{u}_{ik}^2\tilde{d}_{kk} = d_{ii} - \frac{1}{\alpha}v_i^2 \quad (\text{B.64})$$

For  $j = n$ , Eq.(B.63) becomes

$$\tilde{u}_{in}\tilde{d}_{nn}\tilde{u}_{nn} = -\frac{1}{\alpha}v_iv_n \quad (\text{B.65})$$

and from Eq. (B.64),

$$\tilde{u}_{nn}^2\tilde{d}_{nn} = d_{nn} - \frac{1}{\alpha}v_n^2 \quad (\text{B.66})$$

Recalling that  $\tilde{u}_{nn} = 1$ , we get

$$\tilde{d}_{nn} = d_{nn} - \frac{1}{\alpha}v_n^2 \quad (\text{B.67})$$

and

$$\tilde{u}_{in} = -\frac{1}{\alpha\tilde{d}_{nn}}v_iv_n \quad (\text{B.68})$$

So, Eq. (B.63) can be written as

$$\sum_{k=j}^{n-1} \tilde{u}_{ik}\tilde{d}_{kk}\tilde{u}_{jk} + \tilde{u}_{in}\tilde{d}_{nn}\tilde{u}_{jn} = -\frac{1}{\alpha}v_iv_j \quad (\text{B.69})$$

Substituting for  $\tilde{u}_{in}$  and  $\tilde{u}_{jn}$  from Eq. (B.68), we find

$$\sum_{k=j}^{n-1} \tilde{u}_{ik}\tilde{d}_{kk}\tilde{u}_{jk} = -\frac{1}{\alpha} \left[ 1 + \frac{1}{\alpha\tilde{d}_{nn}}v_n^2 \right] v_iv_j \quad (\text{B.70})$$

But since

$$\left[ 1 + \frac{1}{\alpha\tilde{d}_{nn}}v_n^2 \right] = \frac{d_{nn}}{\tilde{d}_{nn}} \quad (\text{B.71})$$

Eq. (B.70) becomes

$$\sum_{k=j}^{n-1} \tilde{u}_{ik} \tilde{d}_{kk} \tilde{u}_{jk} = -\frac{1}{\alpha} \frac{d_{nn}}{\tilde{d}_{nn}} v_i v_j \quad (\text{B.72})$$

So, we can expand Eq.(B.72) as

$$\sum_{k=j}^{n-2} \tilde{u}_{ik} \tilde{d}_{kk} \tilde{u}_{jk} + \tilde{u}_{i,n-1} \tilde{d}_{n-1,n-1} \tilde{u}_{j,n-1} = -\frac{1}{\alpha} \frac{d_{nn}}{\tilde{d}_{nn}} v_i v_j \quad (\text{B.73})$$

We need to obtain  $\tilde{u}_{i,n-1}$  and  $\tilde{d}_{n-1,n-1}$ . First we work on  $\tilde{u}_{i,n-1} \tilde{d}_{n-1,n-1}$  from Eq.(B.63) with  $i = n - 1$  as follows:

$$\tilde{u}_{n-1,n-1}^2 \tilde{d}_{n-1,n-1} + \tilde{u}_{n-1,n}^2 \tilde{d}_{n,n} = d_{n-1} - \frac{1}{\alpha} v_{n-1}^2 \quad (\text{B.74})$$

Recalling that  $\tilde{u}_{n-1,n-1} = 1$  and  $\tilde{u}_{n-1,n}$  was obtained (with  $i = n - 1$ ) in Eq. (B.68), we get

$$\tilde{d}_{n-1,n-1} = d_{n-1,n-1} - \frac{1}{\alpha} \left[ 1 + \frac{1}{\alpha \tilde{d}_{n,n}} v_n^2 \right] v_{n-1}^2 \quad (\text{B.75})$$

Knowing that

$$\left[ 1 + \frac{1}{\alpha \tilde{d}_{n,n}} v_n^2 \right] = \frac{d_{n,n}}{\tilde{d}_{n,n}} \quad (\text{B.76})$$

$\tilde{d}_{n-1,n-1}$  becomes

$$\tilde{d}_{n-1,n-1} = d_{n-1,n-1} - \frac{1}{\alpha} \left( \frac{d_{n,n}}{\tilde{d}_{n,n}} \right) v_{n-1}^2 \quad (\text{B.77})$$

Now we work on  $\tilde{u}_{i,n-1}$ . We recall that from Eq. (B.64), with  $j = n - 1$ , we find that

$$\tilde{u}_{i,n-1} \tilde{d}_{n-1,n-1} \tilde{u}_{n-1,n-1} + \tilde{u}_{i,n} \tilde{d}_{n,n} \tilde{u}_{n-1,n} = -\frac{1}{\alpha} v_i v_{n-1} \quad (\text{B.78})$$

We substitute for  $\tilde{d}_{n-1,n-1}$  from Eq. (B.77), for  $\tilde{u}_{i,n}$  and  $\tilde{u}_{n-1,n}$  from Eq.(B.68) and noting that  $\tilde{u}_{n-1,n-1} = 1$ , we get

$$\tilde{u}_{i,n-1} = -\frac{1}{\alpha \tilde{d}_{n-1,n-1}} \left[ 1 + \frac{1}{\alpha \tilde{d}_{n,n}} v_n^2 \right] v_i v_{n-1} \quad (\text{B.79})$$

Using Eq.(B.76),  $\tilde{u}_{i,n-1}$  becomes

$$\tilde{u}_{i,n-1} = -\frac{1}{\alpha \tilde{d}_{n-1,n-1}} \left( \frac{d_{n,n}}{\tilde{d}_{n,n}} \right) v_i v_{n-1} \quad (\text{B.80})$$

With this in mind, are now prepared to work on Eq. (B.73) and we find that

$$\begin{aligned} \sum_{k=j}^{n-2} \tilde{u}_{ik} \tilde{d}_{kk} \tilde{u}_{jk} &= - \left( \frac{d_{n,n}}{\alpha \tilde{d}_{n,n}} \right) \left[ \frac{d_{n,n}}{\alpha \tilde{d}_{n,n}} v_{n-1}^2 - 1 \right] v_i v_j \\ &= - \left( \frac{d_{n,n}}{\alpha \tilde{d}_{n,n}} \right) \left( \frac{d_{n-1,n-1}}{\tilde{d}_{n-1,n-1}} \right) v_i v_j \end{aligned} \quad (\text{B.81})$$

This has the same form as Eq. (B.69), so this suggests a recursion as follows:

With  $\mathcal{C}_n = -1/\alpha$  for  $j = n, \dots, 1$ :

$$\tilde{d}_{jj} = d_{jj} + \mathcal{C}^j v_j^2 \quad (\text{B.82})$$

$$\tilde{u}_{ij} = \mathcal{C}^j v_i v_j / \tilde{d}_{jj}, \quad k = 1, \dots, j-1 \quad (\text{B.83})$$

$$\mathcal{C}^{j-1} = \mathcal{C}^j d_{jj} / \tilde{d}_{jj} \quad (\text{B.84})$$

From Eq. (B.82), we get

$$\tilde{d}_{jj} = d_{jj} + \mathcal{C}^j v_j^2$$

and from Eq. (B.84), we find that

$$\mathcal{C}^{j-1} = \mathcal{C}^j \frac{d_{jj}}{\tilde{d}_{jj}} = \frac{\mathcal{C}^j d_{jj}}{d_{jj} + \mathcal{C}^j v_j^2} = \frac{d_{jj}}{\frac{d_{jj}}{\mathcal{C}^j} + v_j^2} \quad (\text{B.85})$$

This can be written as

$$\frac{1}{\mathcal{C}^{j-1}} = \frac{1}{\mathcal{C}^j} + \frac{v_j^2}{d_{jj}} \quad (\text{B.86})$$

or

$$-\frac{1}{\mathcal{C}^j} = -\frac{1}{\mathcal{C}^{j-1}} + \frac{v_j^2}{d_{jj}} \quad (\text{B.87})$$

Comparing Eqs. (B.60) and Eq. (B.87) we find that

$$\alpha_j = -\frac{1}{\mathcal{C}^j} \quad (\text{B.88})$$

Using this equation, we find that

$$\tilde{d}_{jj} = d_{jj} \left( \frac{\alpha_{j-1}}{\alpha_j} \right) \quad (\text{B.89})$$

From Eq.(B.83) and using Eqs. (B.89) and (B.88), we can express  $\tilde{u}_{ij}$  as

$$\tilde{u}_{ij} = -\frac{v_i v_j}{d_{jj} \alpha_{j-1}} \quad (\text{B.90})$$

Recalling that  $v_j = d_{jj} f_j$ ,

$$\tilde{u}_{ij} = -\frac{v_i f_j}{\alpha_{j-1}} \quad (\text{B.91})$$

If we define  $\lambda_j$  as

$$\lambda_j \triangleq -\frac{f_j}{\alpha_{j-1}} \quad (\text{B.92})$$

$\tilde{u}_{ij}$  becomes

$$\tilde{u}_{ij} = \lambda_j v_i \quad (\text{B.93})$$

$\tilde{U}_{ij}$  has the structure

$$\tilde{\mathbf{U}} = \begin{bmatrix} 1 & \lambda_2 v_1 & \lambda_3 v_1 & \lambda_4 v_1 & \cdots & \lambda_n v_1 \\ 0 & 1 & \lambda_3 v_2 & \lambda_4 v_2 & \cdots & \lambda_n v_2 \\ 0 & 0 & 1 & \lambda_4 v_3 & \cdots & \lambda_n v_3 \\ 0 & 0 & 0 & 1 & \cdots & \lambda_n v_4 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 \end{bmatrix} \quad (\text{B.94})$$

We can rewrite  $\tilde{\mathbf{U}}$  as

$$\tilde{\mathbf{U}} = \mathbf{I}_n + \begin{bmatrix} \mathbf{0}_{n \times 1} & \lambda_2 \mathbf{v}^{(1)} & \lambda_3 \mathbf{v}^{(2)} & \lambda_4 \mathbf{v}^{(3)} & \cdots & \lambda_n \mathbf{v}^{(n-1)} \end{bmatrix} \quad (\text{B.95})$$

where  $\mathbf{v}^{(j)}$  is an  $n \times 1$  vector defined as

$$\mathbf{v}^{(j)} \triangleq [v_1 \ v_2 \ v_3 \ \cdots \ v_j \ 0 \ \cdots \ 0]^T \quad (\text{B.96})$$

We recall that

$$\mathbf{U} \mathbf{D} \mathbf{U}^T = \bar{\mathbf{U}} \left[ \bar{\mathbf{D}} - \frac{1}{\alpha} \mathbf{v} \mathbf{v}^T \right] \bar{\mathbf{U}}^T$$

and

$$\tilde{\mathbf{U}} \tilde{\mathbf{D}} \tilde{\mathbf{U}}^T \triangleq \bar{\mathbf{D}} - \frac{1}{\alpha} \mathbf{v} \mathbf{v}^T$$

and

$$\mathbf{U} = \bar{\mathbf{U}} \tilde{\mathbf{U}} \quad \text{and} \quad \mathbf{D} = \tilde{\mathbf{D}}$$

With this in mind,  $\mathbf{U}$  is

$$\mathbf{U} = \bar{\mathbf{U}} + \bar{\mathbf{U}} \begin{bmatrix} \mathbf{0}_{n \times 1} & \lambda_2 \mathbf{v}^{(1)} & \lambda_3 \mathbf{v}^{(2)} & \lambda_4 \mathbf{v}^{(3)} & \cdots & \lambda_n \mathbf{v}^{(n-1)} \end{bmatrix} \quad (\text{B.97})$$

If we denote  $\mathbf{U}^{(j)}$  and  $\bar{\mathbf{U}}^{(j)}$  as the  $j$ th columns of  $\mathbf{U}$  and  $\bar{\mathbf{U}}$ , respectively, we find that

$$\mathbf{U}^{(j)} = \bar{\mathbf{U}}^{(j)} + \lambda_j \bar{\mathbf{K}}_{j-1} \quad (\text{B.98})$$

where

$$\bar{\mathbf{K}}_j = \bar{\mathbf{U}} \mathbf{v}^{(j)} = \bar{\mathbf{K}}_{j-1} + v_j \bar{\mathbf{U}}^{(j)} \quad \text{with} \quad \bar{\mathbf{K}}_0 = \mathbf{0}_{n \times 1} \quad (\text{B.99})$$

Finally,

$$\mathbf{K} = \frac{1}{\alpha_n} \bar{\mathbf{K}}_n \quad (\text{B.100})$$

## APPENDIX C

# An Analysis of Dual Inertial-Absolute and Inertial-Relative Navigation Filters

Contributed by Chris D'Souza

This appendix describes a dual inertial-absolute state and dual inertial-relative state navigation filter trade study performed for Orion. The formulation of each of these filters is detailed, the advantages and disadvantages of each are discussed, and a recommendation to use the dual-inertial formulation is made. This appendix is reproduced from **CEV Flight Dynamics Technical Brief** Number FltDyn-CEV-07-141, dated December 21, 2007.

### C.1. Introduction

Orion will need an efficient and well formulated relative navigation filter. Among the many possibilities, two of the most promising will be discussed in this report. The two are the dual inertial-absolute state navigation filter and the dual inertial-relative state navigation filter. The dual inertial-absolute state filter includes the absolute inertial state of both vehicles (with respect to the center of mass of the central body). The dual inertial-relative state navigation filter has as its states the absolute inertial state of the chaser (Orion) vehicle and the relative inertial state of the target with respect to the chaser ( $\mathbf{x}_{rel} = \mathbf{x}_T - \mathbf{x}_C$ ).

### C.2. The Filter Dynamics

**C.2.1. The Dual Inertial-Absolute Filter Dynamics** In general, the inertial states of the chaser vehicle can be expressed as

$$\dot{\mathbf{x}}_C = \mathbf{f}_C(\mathbf{x}_C) + \mathbf{w}_C \quad (\text{C.1})$$

where  $\mathbf{f}_C$  are the chaser nonlinear dynamics and  $\mathbf{w}_C$  is the process (plant) noise (with statistics  $E(\mathbf{w}_C(t)) = \mathbf{0}$ <sup>1</sup> and  $E(\mathbf{w}_C(t)\mathbf{w}_C(\tau)) = Q_C\delta(t - \tau)$ ). Similarly, the inertial states of the target vehicle evolve according to

$$\dot{\mathbf{x}}_T = \mathbf{f}_T(\mathbf{x}_T) + \mathbf{w}_T \quad (\text{C.2})$$

where  $\mathbf{f}_T$  are the target nonlinear dynamics and  $\mathbf{w}_T$  is the process (plant) noise (with statistics  $E(\mathbf{w}_T(t)) = \mathbf{0}$  and  $E(\mathbf{w}_T(t)\mathbf{w}_T(\tau)) = Q_T\delta(t - \tau)$ ). The nominal state dynamics

---

<sup>1</sup>The expectation operator  $E(\cdot)$  for continuous random variables is defined as follows

$$E(X) = \int_{-\infty}^{\infty} x p(x) dx$$

where  $p(x)$  is the probability density function.

can be expressed as

$$\dot{\mathbf{x}}_{C_{nom}} = \mathbf{f}_C(\mathbf{x}_{C_{nom}}) \quad (\text{C.3})$$

$$\dot{\mathbf{x}}_{T_{nom}} = \mathbf{f}_T(\mathbf{x}_{T_{nom}}) \quad (\text{C.4})$$

Defining

$$\delta\mathbf{x}_C \triangleq \mathbf{x}_C - \mathbf{x}_{C_{nom}} \quad (\text{C.5})$$

$$\delta\mathbf{x}_T \triangleq \mathbf{x}_T - \mathbf{x}_{T_{nom}} \quad (\text{C.6})$$

and taking derivatives and expanding to first-order, we get

$$\delta\dot{\mathbf{x}}_C = A_C(\mathbf{x}_{C_{nom}}) \delta\mathbf{x}_C + \mathbf{w}_C \quad (\text{C.7})$$

$$\delta\dot{\mathbf{x}}_T = A_T(\mathbf{x}_{T_{nom}}) \delta\mathbf{x}_T + \mathbf{w}_T \quad (\text{C.8})$$

where

$$A_C(\mathbf{x}_{C_{nom}}) \triangleq \left( \frac{\partial \mathbf{f}_C}{\partial \mathbf{x}_C} \right)_{\mathbf{x}_C = \mathbf{x}_{C_{nom}}} \quad \text{and} \quad A_T(\mathbf{x}_{T_{nom}}) \triangleq \left( \frac{\partial \mathbf{f}_T}{\partial \mathbf{x}_T} \right)_{\mathbf{x}_T = \mathbf{x}_{T_{nom}}} \quad (\text{C.9})$$

Equivalently, we can express the filter errors as

$$\delta\hat{\mathbf{x}}_C \triangleq \hat{\mathbf{x}}_C - \mathbf{x}_C \quad (\text{C.10})$$

$$\delta\hat{\mathbf{x}}_T \triangleq \hat{\mathbf{x}}_T - \mathbf{x}_T \quad (\text{C.11})$$

where, with a bit of abuse of notation<sup>2</sup>

$$\hat{\mathbf{x}}_C \triangleq E(\mathbf{x}_C) \quad \text{and} \quad \hat{\mathbf{x}}_T \triangleq E(\mathbf{x}_T) \quad (\text{C.12})$$

so that the filter error dynamics evolve as

$$\delta\dot{\hat{\mathbf{x}}}_C = A_C(\mathbf{x}_C) \delta\hat{\mathbf{x}}_C + \mathbf{w}_C \quad (\text{C.13})$$

$$\delta\dot{\hat{\mathbf{x}}}_T = A_T(\mathbf{x}_T) \delta\hat{\mathbf{x}}_T + \mathbf{w}_T \quad (\text{C.14})$$

We can, therefore, write the *inertial-absolute*<sup>3</sup> filter error dynamics (dropping the functional dependence for compactness) as

$$\begin{bmatrix} \delta\dot{\hat{\mathbf{x}}}_C \\ \delta\dot{\hat{\mathbf{x}}}_T \end{bmatrix} = \begin{bmatrix} A_C & \mathbf{0} \\ \mathbf{0} & A_T \end{bmatrix} \begin{bmatrix} \delta\hat{\mathbf{x}}_C \\ \delta\hat{\mathbf{x}}_T \end{bmatrix} + \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{w}_C \\ \mathbf{w}_T \end{bmatrix} \quad (\text{C.15})$$

Defining

$$P_{IA} \triangleq E \left\{ \begin{bmatrix} \delta\hat{\mathbf{x}}_C \\ \delta\hat{\mathbf{x}}_T \end{bmatrix} \begin{bmatrix} \delta\hat{\mathbf{x}}_C^T & \delta\hat{\mathbf{x}}_T^T \end{bmatrix} \right\} = \begin{bmatrix} E[\delta\hat{\mathbf{x}}_C \delta\hat{\mathbf{x}}_C^T] & E[\delta\hat{\mathbf{x}}_C \delta\hat{\mathbf{x}}_T^T] \\ E[\delta\hat{\mathbf{x}}_T \delta\hat{\mathbf{x}}_C^T] & E[\delta\hat{\mathbf{x}}_T \delta\hat{\mathbf{x}}_T^T] \end{bmatrix} = \begin{bmatrix} P_{C,C} & P_{C,T} \\ P_{T,C} & P_{T,T} \end{bmatrix} \quad (\text{C.16})$$

where the subscript *IA* denotes that this is the covariance associated with the inertial-absolute filter. The differential equation for the covariance (assuming that the plant/process noise for the two vehicles are independent and are independent of the states of the two vehicles) is

$$\dot{P}_{IA} = A_{IA} P_{IA} + P_{IA} A_{IA}^T + G_{IA} Q_{IA} G_{IA}^T \quad (\text{C.17})$$

<sup>2</sup>To be precise,  $\hat{\mathbf{x}}_C$  should be written in terms of the conditional expectation

$$\hat{\mathbf{x}}_{C_k} = E(\mathbf{x}_C | \mathbf{Z}_1, \dots, \mathbf{Z}_k) \quad \text{and} \quad \hat{\mathbf{x}}_{T_k} = E(\mathbf{x}_T | \mathbf{Z}_1, \dots, \mathbf{Z}_k)$$

with measurements  $\mathbf{Z}_1$  through  $\mathbf{Z}_k$ . At the initial time  $\hat{\mathbf{x}}_{C_0} = E(\mathbf{x}_{C_0})$  and  $\hat{\mathbf{x}}_{T_0} = E(\mathbf{x}_{T_0})$

<sup>3</sup>In order to distinguish between the two filters, we call this filter the (dual) *inertial-absolute* filter because both the chaser and the target states are expressed in terms of absolute inertial coordinates.

where

$$A_{IA} = \begin{bmatrix} A_C & \mathbf{0} \\ \mathbf{0} & A_T \end{bmatrix}, \quad G_{IA} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}, \quad Q_{IA} = \begin{bmatrix} Q_C & \mathbf{0} \\ \mathbf{0} & Q_T \end{bmatrix} \quad (\text{C.18})$$

with the initial condition

$$P_{IA} = \begin{bmatrix} P_{C,C_0} & \mathbf{0} \\ \mathbf{0} & P_{T,T_0} \end{bmatrix} \quad (\text{C.19})$$

where  $P_{C,C_0}$  and  $P_{T,T_0}$  are the initial covariances of the chaser and target inertial (absolute) states, respectively. Finally, the covariance of the relative state is

$$P_{rel,rel} = E \left[ (\delta \hat{\mathbf{x}}_C - \delta \hat{\mathbf{x}}_T) (\delta \hat{\mathbf{x}}_C - \delta \hat{\mathbf{x}}_T)^T \right] \quad (\text{C.20})$$

$$= P_{C,C} + P_{T,T} - P_{T,C} - P_{C,T} = P_{C,C} + P_{T,T} - P_{T,C} - P_{T,C}^T \quad (\text{C.21})$$

**C.2.2. The Dual Inertial-Relative Filter Dynamics** Consistent with the earlier definitions, we define the inertial relative state as

$$\mathbf{x}_{rel} \triangleq \mathbf{x}_T - \mathbf{x}_C \quad (\text{C.22})$$

Taking derivatives of this equation and substituting from Eqs.(1) and (2) yields

$$\dot{\mathbf{x}}_{rel} = \mathbf{f}_T(\mathbf{x}_T) - \mathbf{f}_C(\mathbf{x}_C) + \mathbf{w}_T - \mathbf{w}_C \quad (\text{C.23})$$

Expanding this equation to first-order yields

$$\delta \dot{\mathbf{x}}_{rel} = A_T(\mathbf{x}_T) \delta \hat{\mathbf{x}}_T - A_C(\mathbf{x}_C) \delta \hat{\mathbf{x}}_C + \mathbf{w}_T - \mathbf{w}_C \quad (\text{C.24})$$

$$= A_T(\mathbf{x}_T) (\delta \hat{\mathbf{x}}_{rel} + \delta \hat{\mathbf{x}}_C) - A_C(\mathbf{x}_C) \delta \hat{\mathbf{x}}_C + \mathbf{w}_T - \mathbf{w}_C \quad (\text{C.25})$$

$$= (A_T - A_C) \delta \hat{\mathbf{x}}_C + A_T \delta \hat{\mathbf{x}}_{rel} + (\mathbf{w}_T - \mathbf{w}_C) \quad (\text{C.26})$$

Therefore we write the *inertial-relative*<sup>4</sup> filter error dynamics (once again dropping the functional dependence for compactness) as

$$\begin{bmatrix} \delta \dot{\hat{\mathbf{x}}}_C \\ \delta \dot{\hat{\mathbf{x}}}_{rel} \end{bmatrix} = \begin{bmatrix} A_C & \mathbf{0} \\ A_T - A_C & A_T \end{bmatrix} \begin{bmatrix} \delta \hat{\mathbf{x}}_C \\ \delta \hat{\mathbf{x}}_T \end{bmatrix} + \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{I} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{w}_C \\ \mathbf{w}_T \end{bmatrix} \quad (\text{C.27})$$

Defining, as before,

$$P_{IR} \triangleq E \left\{ \begin{bmatrix} \delta \hat{\mathbf{x}}_C \\ \delta \hat{\mathbf{x}}_{rel} \end{bmatrix} \begin{bmatrix} \delta \hat{\mathbf{x}}_C^T & \delta \hat{\mathbf{x}}_{rel}^T \end{bmatrix} \right\} = \begin{bmatrix} E[\delta \hat{\mathbf{x}}_C \delta \hat{\mathbf{x}}_C^T] & E[\delta \hat{\mathbf{x}}_C \delta \hat{\mathbf{x}}_{rel}^T] \\ E[\delta \hat{\mathbf{x}}_{rel} \delta \hat{\mathbf{x}}_C^T] & E[\delta \hat{\mathbf{x}}_{rel} \delta \hat{\mathbf{x}}_{rel}^T] \end{bmatrix} \quad (\text{C.28})$$

$$= \begin{bmatrix} P_{C,C} & P_{C,rel} \\ P_{rel,C} & P_{rel,rel} \end{bmatrix} \quad (\text{C.29})$$

where the subscript *IR* denotes that this is the covariance associated with the inertial-relative filter. The differential equation for the covariance is

$$\dot{P}_{IR} = A_{IR} P_{IR} + P_{IR} A_{IR}^T + G_{IR} Q_{IR} G_{IR}^T \quad (\text{C.30})$$

where

$$A_{IR} = \begin{bmatrix} A_C & \mathbf{0} \\ A_T - A_C & A_T \end{bmatrix}, \quad G_{IR} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{I} & \mathbf{I} \end{bmatrix}, \quad Q_{IR} = \begin{bmatrix} Q_C & \mathbf{0} \\ \mathbf{0} & Q_T \end{bmatrix} \quad (\text{C.31})$$

<sup>4</sup>We refer to this filter as the *inertial-relative* filter to distinguish it from the prior *inertial-absolute* filter. In this formulation, the filter states consist of the inertial absolute chaser state and the inertial relative target state.

where the initial covariance of the inertial-relative state is found to be<sup>5</sup>

$$P_{IR_0} = \begin{bmatrix} P_{C,C_0} & P_{C,rel_0} \\ P_{rel,C_0} & P_{rel,rel_0} \end{bmatrix} = \begin{bmatrix} P_{C,C_0} & -P_{C,C_0} \\ -P_{C,C_0} & P_{C,C_0} + P_{T,T_0} \end{bmatrix} \quad (C.32)$$

in order to be consistent with the inertial-absolute formulation. We can also express Eqs. (30) and (31) as

$$\dot{P}_{IR} = A_{IR}P_{IR} + P_{IR}A_{IR}^T + Q'_{IR} \quad (C.33)$$

where

$$Q'_{IR} = \begin{bmatrix} Q_C & -Q_C \\ -Q_C & Q_{rel} \end{bmatrix} \quad (C.34)$$

where  $Q_{rel} = Q_T + Q_C$ . This may simplify tuning. While the covariance of the relative state is easily determined since it is lower right partition of the covariance matrix ( $P_{rel,rel}$ ), the covariance of the target vehicle (inertial) state is found (after a bit of manipulation) to be

$$P_{T,T} = P_{rel,rel} + P_{C,C} + P_{rel,C} + P_{C,rel} = P_{rel,rel} + P_{C,C} + P_{rel,C} + P_{rel,C}^T \quad (C.35)$$

### C.3. Incorporation of Measurements

Whereas the previous section analyzed the filter error dynamics/propagation as it applies to the inertial-absolute and inertial-relative filter formulation, this section will analyze the difference in measurement processing between the two filters. Obviously, those measurements that are tied only to the chaser absolute inertial state have the same instantiation in both formulations. This section will only address those measurement types which have different formulations in the two filters. In particular, the measurement and the measurement partials will be discussed.

#### C.3.1. The Dual Inertial-Absolute Measurement Formulations

C.3.1.1. *The Target Inertial State Ground Update* There will be instances during which the on-board filter will need to process ground updates of the target vehicle. For the inertial-absolute formulation, the measurement takes the following expression

$$\mathbf{z}_{TGU}^{IA} = \mathbf{x}_T + \boldsymbol{\nu}_{TGU} \quad \text{and} \quad \boldsymbol{\nu}_{TGU} \sim N(\mathbf{0}, R_{TGU}) \quad (C.36)$$

Since the target state is a member of this filter's state-space, the measurement partials associated with this measurement for the inertial-absolute formulation is

$$H_{TGU}^{IA} = [\mathbf{0} \quad \mathbf{I}] \quad (C.37)$$

C.3.1.2. *Range Measurements* For the case of range measurements (either from the RF link or from the Lidar), the measurement equation can be written simply as

$$\mathbf{z}_{range}^{IA} = \sqrt{(\mathbf{r}_T - \mathbf{r}_C) \cdot (\mathbf{r}_T - \mathbf{r}_C)} + b_{range} + \nu_{range} \quad (C.38)$$

with the range (measurement) noise statistics  $\nu_{range} \sim N(0, R_{range})$ . Since the target state is a member of this filter's state-space, the measurement partials associated with this

---

<sup>5</sup>We assume that at the initial time, the chaser and target initial error covariance matrices are uncorrelated.

measurement for the inertial-absolute measurement are

$$\frac{\partial \mathbf{z}_{range}^{IA}}{\partial \mathbf{r}_C} = -\mathbf{u}_{rel}^T \quad (\text{C.39})$$

$$\frac{\partial \mathbf{z}_{range}^{IA}}{\partial \mathbf{v}_C} = \mathbf{0}_{1 \times 3} \quad (\text{C.40})$$

$$\frac{\partial \mathbf{z}_{range}^{IA}}{\partial \mathbf{r}_T} = \mathbf{u}_{rel}^T \quad (\text{C.41})$$

$$\frac{\partial \mathbf{z}_{range}^{IA}}{\partial \mathbf{v}_T} = \mathbf{0}_{1 \times 3} \quad (\text{C.42})$$

where

$$\mathbf{u}_{rel} = \frac{(\mathbf{r}_T - \mathbf{r}_C)}{|\mathbf{r}_T - \mathbf{r}_C|} \quad (\text{C.43})$$

C.3.1.3. *Bearing Measurements* Bearing measurements, which can be obtained from either the star-tracker or the lidar, are

$$\mathbf{z}_{bearing}^{IA} = \begin{bmatrix} \alpha \\ \delta \end{bmatrix} + \begin{bmatrix} b_\alpha \\ b_\delta \end{bmatrix} + \begin{bmatrix} \nu_\alpha \\ \nu_\delta \end{bmatrix} = \begin{bmatrix} \alpha \\ \delta \end{bmatrix} + \mathbf{b}_{bearing} + \boldsymbol{\nu}_{bearing} \quad (\text{C.44})$$

where the angles  $\alpha$  and  $\delta$ , the azimuth and elevation angles in the sensor frame with biases  $b_\alpha$  and  $b_\delta$ , and with noise characteristics

$$E(\boldsymbol{\nu}_{bearing}) = \mathbf{0}_{2 \times 1} \quad \text{and} \quad E(\boldsymbol{\nu}_{bearing} \boldsymbol{\nu}_{bearing}^T) = \begin{bmatrix} R_\alpha & 0 \\ 0 & R_\delta \end{bmatrix} \quad (\text{C.45})$$

We can express  $\alpha$  and  $\delta$  in terms of the line-of-sight vector which is defined as

$$\begin{bmatrix} \cos \alpha \cos \delta \\ \sin \alpha \cos \delta \\ \sin \delta \end{bmatrix} \triangleq \frac{1}{r} T_{SB}(q_{SB}) T_{BI}(q_{BI}) (\mathbf{r}_T - \mathbf{r}_C) \quad (\text{C.46})$$

where  $T_{SB}$  is the transformation matrix from the body(IMU) frame to the sensor frame (with  $q_{SB}$  being the quaternion associated with the transformation from body frame to sensor frame) and  $T_{BI}$  is the transformation matrix from the inertial frame to the body(IMU) frame (with  $q_{BI}$  being the quaternion associated with the transformation from the inertial frame to the body frame).

With this in hand, the measurement partials can be obtained (after a bit of manipulation [1]) to be

$$\frac{\partial \alpha}{\partial \mathbf{r}_C} = -\frac{\mathbf{u}_\alpha^T}{r} T_{SB}(q_{SB}) T_{BI}(q_{BI}) \quad (\text{C.47})$$

$$\frac{\partial \alpha}{\partial \mathbf{v}_C} = \mathbf{0}_{1 \times 3} \quad (\text{C.48})$$

$$\frac{\partial \alpha}{\partial \mathbf{r}_T} = \frac{\mathbf{u}_\alpha^T}{r} T_{SB}(q_{SB}) T_{BI}(q_{BI}) \quad (\text{C.49})$$

$$\frac{\partial \alpha}{\partial \mathbf{v}_T} = \mathbf{0}_{1 \times 3} \quad (\text{C.50})$$

and

$$\frac{\partial \delta}{\partial \mathbf{r}_C} = -\frac{\mathbf{u}_\delta^T}{r} T_{SB}(q_{SB}) T_{BI}(q_{BI}) \quad (\text{C.51})$$

$$\frac{\partial \delta}{\partial \mathbf{v}_C} = \mathbf{0}_{1 \times 3} \quad (\text{C.52})$$

$$\frac{\partial \delta}{\partial \mathbf{r}_T} = \frac{\mathbf{u}_\delta^T}{r} T_{SB}(q_{SB}) T_{BI}(q_{BI}) \quad (\text{C.53})$$

$$\frac{\partial \delta}{\partial \mathbf{v}_T} = \mathbf{0}_{1 \times 3} \quad (\text{C.54})$$

where

$$\mathbf{u}_\alpha = \frac{1}{\cos \delta} \begin{bmatrix} -\sin \alpha \\ \cos \alpha \\ 0 \end{bmatrix} \quad (\text{C.55})$$

$$\mathbf{u}_\delta = \begin{bmatrix} -\cos \alpha \sin \delta \\ -\cos \alpha \cos \delta \\ \cos \delta \end{bmatrix} \quad (\text{C.56})$$

### C.3.2. The Dual Inertial-Relative Measurement Formulations

C.3.2.1. *The Target Inertial State Ground Update* For the inertial-relative formulation, the measurement takes the following expression

$$\mathbf{z}_{TGU}^{IR} = \mathbf{x}_{rel} + \mathbf{x}_C + \boldsymbol{\nu}_{TGU} = \text{and } \boldsymbol{\nu}_{TGU} \sim N(\mathbf{0}, R_{TGU}) \quad (\text{C.57})$$

Since the target state is not a member of this filter's state-space, the measurement partials associated with this measurement for the inertial-relative formulation is

$$H_{TGU}^{IR} = [\mathbf{I} \quad \mathbf{I}] \quad (\text{C.58})$$

C.3.2.2. *Range Measurements* For the case of range measurements with the inertial-relative filter, the measurement equation can be written simply as

$$\mathbf{z}_{range}^{IA} = \sqrt{\mathbf{r}_{rel}^T \mathbf{r}_{rel}} + b_{range} + \nu_{range} \quad (\text{C.59})$$

with, as before, the range (measurement) noise statistics  $\nu_{range} \sim N(0, R_{range})$ . Since the relative state is a member of this filter's state-space, the measurement partials associated with this measurement for the inertial-relative measurement are

$$\frac{\partial \mathbf{z}_{range}^{IA}}{\partial \mathbf{r}_C} = \mathbf{0}_{1 \times 3} \quad (\text{C.60})$$

$$\frac{\partial \mathbf{z}_{range}^{IA}}{\partial \mathbf{v}_C} = \mathbf{0}_{1 \times 3} \quad (\text{C.61})$$

$$\frac{\partial \mathbf{z}_{range}^{IA}}{\partial \mathbf{r}_{rel}} = \frac{\mathbf{r}_{rel}^T}{|\mathbf{r}_{rel}|} \quad (\text{C.62})$$

$$\frac{\partial \mathbf{z}_{range}^{IA}}{\partial \mathbf{v}_{rel}} = \mathbf{0}_{1 \times 3} \quad (\text{C.63})$$

**C.3.2.3. Bearing Measurements** Bearing measurements, as in the previous formulation (in Eqs.(44) and (45)), can be expressed in terms of the line-of-sight vector and the relative position vector which can be expressed as follows

$$\begin{bmatrix} \cos \alpha \cos \delta \\ \sin \alpha \cos \delta \\ \sin \delta \end{bmatrix} = \frac{1}{r} T_{SB}(q_{SB})T_{BI}(q_{BI})\mathbf{r}_{rel} \quad (\text{C.64})$$

the quantities in Eq.(64) are defined in Section 3.1.3. The measurement partials are expressed as

$$\frac{\partial \alpha}{\partial \mathbf{r}_C} = \mathbf{0}_{1 \times 3} \quad (\text{C.65})$$

$$\frac{\partial \alpha}{\partial \mathbf{v}_C} = \mathbf{0}_{1 \times 3} \quad (\text{C.66})$$

$$\frac{\partial \alpha}{\partial \mathbf{r}_{rel}} = \frac{\mathbf{u}_\alpha^T}{r} T_{SB}(q_{SB})T_{BI}(q_{BI}) \quad (\text{C.67})$$

$$\frac{\partial \alpha}{\partial \mathbf{v}_{rel}} = \mathbf{0}_{1 \times 3} \quad (\text{C.68})$$

and

$$\frac{\partial \delta}{\partial \mathbf{r}_C} = \mathbf{0}_{1 \times 3} \quad (\text{C.69})$$

$$\frac{\partial \delta}{\partial \mathbf{v}_C} = \mathbf{0}_{1 \times 3} \quad (\text{C.70})$$

$$\frac{\partial \delta}{\partial \mathbf{r}_{rel}} = \frac{\mathbf{u}_\delta^T}{r} T_{SB}(q_{SB})T_{BI}(q_{BI}) \quad (\text{C.71})$$

$$\frac{\partial \delta}{\partial \mathbf{v}_{rel}} = \mathbf{0}_{1 \times 3} \quad (\text{C.72})$$

where

$$\mathbf{u}_\alpha = \frac{1}{\cos \delta} \begin{bmatrix} -\sin \alpha \\ \cos \alpha \\ 0 \end{bmatrix} \quad (\text{C.73})$$

$$\mathbf{u}_\delta = \begin{bmatrix} -\cos \alpha \sin \delta \\ -\sin \alpha \sin \delta \\ \cos \delta \end{bmatrix} \quad (\text{C.74})$$

#### C.4. Analysis of the Merits of the Inertial-Absolute and Inertial-Relative Filters

The Flight Day 1 rendezvous trajectory and models as described in [2] were used to analyze the two filter formulations. Both formulations had the same driving dynamics and measurement models – only the covariance propagation and covariance updates were different. With this in mind, the comparison of the two filter formulation as it related to covariance operations were analyzed.

**C.4.1. Covariance Propagation** First it must be pointed out that the propagation of the filter dynamics are identical between both filter parameterizations. That is to say, in each filter, the inertial absolute states of both the chaser vehicle and the target vehicle will be propagated. The difference arises in the propagation of the covariance matrices

associated with each of the filter parameterizations. In the *IA* filter, the covariances (and cross-covariances) of the chaser inertial states and the target inertial states are computed. It should be noted that the *dynamics* of the two vehicles' states are inherently *un*-correlated (see  $A_{IA}$  in Eq.(18)). In contrast, for the *IR* filter the covariances (and cross-covariances) the chaser inertial states and the inertial relative states (of the target with respect to the chaser) are computed. It should be noted that the *dynamics* in this filter's states are inherently *correlated* (see  $A_{IR}$  in Eq.(31)). Hence, there are inherently more non-zero computations (both additions and multiplications) involved<sup>6</sup>. Hence, there is more room for round-off errors in the covariance propagation in the *IR* filter. In order to see this, Table 1 contains an analysis of the propagation error as a function of the propagation interval. This propagation is carried out without process noise on either the chaser or target states. In addition, the chaser and the target states are uncorrelated at the initial time. Hence, since there are no measurements which could correlate the two vehicles' states, the correlation coefficients should remain zero (i.e.  $P_{C,T} = \mathbf{0}_{6 \times 6}$ ) throughout the interval. This was verified to be the case. Because of the reduced number of computations inherent in the *IA* filter formulation, it was assumed that the *IA* filter propagation was the 'truth' and the *IR* filter was compared to it.

$\Delta t(sec)$	$ P_{rel,rel} _2$	$ P_{T,T} _2$	$ \delta P_{rel,rel} _2/ P_{rel,rel} _2$	$ \delta P_{T,T} _2/ P_{T,T} _2$
100	1.010E7	4.150E3	1.529E-16	4.867E-13
1000	6.997E7	4.519E3	5.729E-9	8.869E-6
10000	1.239E9	7.943E4	1.926E-5	0.349

TABLE 1. Numerical Precision Comparison of the *IA* and *IR* filter formulations for propagation *without* process noise

$\Delta t(sec)$	$ P_{rel,rel} _2$	$ P_{T,T} _2$	$ \delta P_{rel,rel} _2/ P_{rel,rel} _2$	$ \delta P_{T,T} _2/ P_{T,T} _2$
100	1.010E7	4.150E3	1.236E-13	1.551E-7
1000	6.997E7	4.519E3	2.172E-8	2.320E-4
10000	1.239E9	7.943E4	1.905E-5	0.392

TABLE 2. Numerical Precision Comparison of the *IA* and *IR* filter formulations for propagation *with* process noise

It is clear that the additional non-zero multiplications and additions for the *IR* filter formulation compared to the *IA* formulation result in a build-up of round-off error. This has

<sup>6</sup>First, notice that in the *IR* filter, the term  $A_T - A_C$  will inherently cause a loss of precision. Second, assuming only the gravity gradient term in  $A$ , which is symmetric,  $A_T - A_C$  involves 6 additions/subtractions. The term  $(A_T - A_C)P_{CC}$  in Eqs.(30) and (31) involve 18 multiplications and 12 additions/subtractions. The term  $(A_T - A_C)P_{C,rel}$  (or  $P_{rel,C}(A_T - A_C)^T$ ) involve 27 multiplications and 18 additions/subtractions. The terms  $(A_T - A_C)P_{CC} + A_T P_{rel,C}$  and  $(A_T - A_C)P_{C,rel} + A_T P_{rel,rel}$  each involve 9 additions/subtractions. These total 45 multiplications and 54 additions/subtractions. This is doubled because of the symmetric nature of the covariance matrix. Therefore, there are 90 *additional* multiplications and 108 *additional* additions/subtractions *per function evaluation* in the *IR* filter formulation over the *IA* filter formulation. For a fourth-order Runge-Kutta integration method, the *IR* filter formulation results in an *additional* 360 multiplications and 432 additions/subtractions *per integration step*. Each of these operations results in a numerical loss of precision in finite-state machines.

the effect of reducing the propagation accuracy of the *IR* filter vis-à-vis the *IA* filter. The additional operations, in concert with the accompanying loss of precision, make a strong case for the use of the *IA* filter formulation.

**C.4.2. Measurement Update** It should be apparent from comparing Eqs. (35) and (56) that for the case of the target ground-update, there are 18 more multiplications (because of the identity matrix) for the *IR* formulation than the *IA* formulation for each target ground-update.

For all other relative measurements, there are more operations for the *IA* filter formulation than for the *IR* formulation. In fact, for the range measurements, there are more 54 multiplications and 54 more additions for the *IA* formulation than the *IR* formulation for *each* range measurement update.

For bearing measurements, there are more 108 multiplications and 108 more additions for the *IA* formulation than the *IR* formulation for *each* bearing measurement update.

So, for relative sensor measurements, clearly there are more multiplications and more additions for the *IA* filter formulation than for the *IR* formulation.

### C.5. Conclusions

While the inertial-absolute and inertial-relative filter formulations are mathematically equivalent, the implementation on finite-state machines influences the choice.

With regard to propagation of the covariance matrices, there are more computations for the *IR* filter than the (mathematically equivalent) *IA* filter. These additional computations, in concert with the types of operations, result in a (significant) loss of precision with regard to the propagation of the covariance matrices.

With regard to the measurement updates to the covariance matrices, for the case of relative navigation measurements, there are fewer non-zero operations for the *IR* filter than for the *IA* filter. This is one of the strengths of this filter (*IR*)formulation, and if the computation and precision with regard to the propagation of the covariance matrices were the same, the *IR* filter would be advantageous in terms of computations and precision.

After considering all these factors, the Inertial-Absolute filter is recommended for use on Orion.



## Bibliography

- [1] W.S. Agee and R.H. Turner. Triangular decomposition of a positive definite matrix plus a symmetric dyad with application to kalman filtering. *White Sands Missile Range Tech. Rep. No. 38*, 1972.
- [2] G. J. Bierman. *Factorization Methods for Discrete Sequential Estimation*. Academic Press, Dover Publications, New York, 1977, 2006.
- [3] Robert G. Brown and Patrick Y.C. Hwang. *Introduction to Random Signals and Applied Kalman Filtering*. John Wiley and Sons, Inc., New York, NY, 3rd edition, 1997.
- [4] A. E. Bryson and M. Frazier. Smoothing for linear and nonlinear dynamic systems. Technical Report TDR-63-119, Aeronautical Systems Division, Wright-Patterson Air Force Base, Ohio, Sept. 1962.
- [5] N.A. Carlson. Fast triangular factorization of the square root filter. *AIAA Journal*, 11(9):1259–1265, September 1973.
- [6] J. R. Carpenter and K. T. Alfriend. Navigation accuracy guidelines for orbital formation flying. *Journal of the Astronautical Sciences*, 53(2):207–219, 2006. Also appears as AIAA Paper 2003–5443.
- [7] J. Russell Carpenter. Navigation filter best practices. <https://mediaex-server.larc.nasa.gov/Academy/Catalog/Full/f1d0abb028d3491f8701da3fc64bcb2021>, January 2015. Accessed: January 4, 2018.
- [8] J. Russell Carpenter and Kevin Berry. Artificial damping for stable long-term orbital covariance propagation. In *Astrodynamics 2007*, volume 129 of *Advances in the Astronautical Sciences*, pages 1697–1707. Univelt, 2008.
- [9] J. Russell Carpenter and Emil R. Schiesser. Semimajor axis knowledge and gps orbit determination. *NAVIGATION: Journal of The Institute of Navigation*, 48(1):57–68, Spring 2001. Also appears as AAS Paper 99–190.
- [10] Russell Carpenter and Taesul Lee. A stable clock error model using coupled first- and second-order gauss-markov processes. In *AAS/AIAA 18th Spaceflight Mechanics Meeting*, volume 130, pages 151–162. Univelt, 2008.
- [11] W. H. Clohessy and R. S. Wiltshire. Terminal guidance for satellite rendezvous. *Journal of the Aerospace Sciences*, 27(5):653–658, 674, 1960.
- [12] W. F. Denham and S. Pines. Sequential estimation when measurement function nonlinearity is comparable to measurement error. *AIAA Journal*, 4(6):1071–1076, June 1966.
- [13] Cornelius J. Dennehy and J. Russell Carpenter. A summary of the rendezvous, proximity operations, docking, and undocking (rpodu) lessons learned from the defense advanced research project agency (darpa) orbital express (oe) demonstration system mission. NASA Technical Memorandum 2011-217088, NASA Engineering and Safety Center, 2011.
- [14] J. L. Farrell. Attitude determination by Kalman filtering. *Automatica*, 6:419–430, 1970.
- [15] James L. Farrell. Attitude determination by Kalman filtering. Contractor Report NASA-CR-598, NASA Goddard Space Flight Center, Washington, DC, Sept. 1964.
- [16] P. Ferguson and J. How. Decentralized estimation algorithms for formation flying spacecraft. In *AIAA Guidance, Navigation and Control Conference*, 2003.
- [17] Donald C. Fraser. *A New Technique for the Optimal Smoothing of Data*. Sc.D. thesis, Massachusetts Institute of Technology, Cambridge, Massachusetts, 1967.
- [18] Donald C. Fraser and James E. Potter. The optimum smoother as a combination of two optimum linear filters. *IEEE Transactions on Automatic Control*, AC-14(4):387–390, Aug. 1969.
- [19] Eliezer Gai, Kevin Daly, James Harrison, and Linda Lemos. Star-sensor-based attitude/attitude rate estimator. *Journal of Guidance, Control, and Dynamics*, 8(5):560–565, Sept.-Oct. 1985.
- [20] Arthur Gelb, editor. *Applied Optimal Estimation*. The MIT Press, Cambridge, MA, 1974.

- [21] David K. Geller. Orbital rendezvous: When is autonomy required? *Journal of Guidance, Control and Dynamics*, 30(4):974–981, July–August 2007.
- [22] Herbert Goldstein. *Classical Mechanics*. Addison-Wesley Publishing Company, Reading, MA, 2nd edition, 1980.
- [23] Gene Howard Golub and Charles F. Van Loan. *Matrix Computations*. JHU Press, 3rd edition, 1996.
- [24] M. Grigoriu. Response of dynamic systems to poisson white noise. *Journal of Sound and Vibration*, 195(3):375–389, 1996.
- [25] C. F. Hanak. Reducing the effects of measurement ordering on the gkf algorithm via a hybrid linear/extended kalman filter. Technical Report FltDyn-CEV-06-0107, NASA Johnson Space Center, August 2006.
- [26] G. W. Hill. Researches in the lunar theory. *American Journal of Mathematics*, 1(1):5–26, 129–147, 245–260, 1878.
- [27] Jonathan P. How, Louis S. Breger, Megan Mitchell, Kyle T. Alfriend, and Russell Carpenter. Differential semimajor axis estimation performance using carrier-phase differential global positioning system measurements. *Journal of Guidance, Control, and Dynamics*, 30(2):301–313, Mar-Apr 2007.
- [28] Peter C. Hughes. *Spacecraft Attitude Dynamics*. Wiley, New York, NY, 1986.
- [29] M. Idan. Estimation of Rodrigues parameters from vector observations. *IEEE Transactions on Aerospace and Electronic Systems*, 32(2):578–586, April 1996.
- [30] Andrew H. Jazwinski. *Stochastic Processes and Filtering Theory*. Academic Press, Inc., and Dover Publications, Inc., New York, NY, and Mineola, NY, 1970 and 2007.
- [31] S. C. Jenkins and D. K. Geller. State estimation and targeting for autonomous rendezvous and proximity operations. In *Proceedings of the AAS/AIAA Astrodynamics Specialists Conference*. Mackinac Island, MI, August 19–23 2007.
- [32] S. J. Julier and J. K. Uhlmann. A new extension of the kalman filter to nonlinear systems. In *Int. Symp. Aerospace/Defense Sensing, Simul. and Controls*, Orlando, FL, 1997.
- [33] Simon Julier and Jeffrey Uhlmann. Authors’ reply. *IEEE Transactions on Automatic Control*, 47(8):1408–1409, August 2002.
- [34] Simon Julier, Jeffrey Uhlmann, and Hugh F. Durrant-Whyte. A new method for the nonlinear transformation of means and covariances in filters and estimators. *IEEE Transactions on Automatic Control*, 45(3):477–482, March 2000.
- [35] John L. Junkins and J. D. Turner. *Optimal Spacecraft Rotational Maneuvers*. Elsevier, New York, NY, 1986.
- [36] R. E. Kalman. A new approach to linear filtering and prediction problems. *Transactions of the ASME – Journal of Basic Engineering*, 82:35–45, 1960.
- [37] Christopher D. Karlgaard and Hanspeter Schaub. Nonsingular attitude filtering using modified Rodrigues parameters. *The Journal of the Astronautical Sciences*, 57(4):777–791, Oct.–Dec. 2010.
- [38] B. A. Kriegsman and Y. C. Tau. Shuttle navigation system for entry and landing mission phases. *Journal of Spacecraft*, 12(4):213–219, April 1975.
- [39] W. M. Lear. Multi-phase navigation program for the space shuttle orbiter. Internal Note No. 73-FM-132, NASA Johnson Space Center, Houston, TX, 1973.
- [40] William M. Lear. Kalman Filtering Techniques. Technical Report JSC-20688, NASA Johnson Space Center, Houston, TX, 1985.
- [41] Deok-Jin Lee and Kyle T. Alfriend. Additive divided difference filtering for attitude estimation using modified rodrigues parameters. In John L. Crassidis et al., editors, *Proceedings of the F. Landis Markley Astronautics Symposium*, volume 132 (CD-ROM Supplement) of *Advances in the Astronautical Sciences*. American Astronautical Society, Univelt, 2008.
- [42] Tine Lefebvre, Herman Bruyninckx, and Joris De Schutter. Comment on ‘a new method for the nonlinear transformation of means and covariances in filters and estimators’. *IEEE Transactions on Automatic Control*, 47(8):1406–1408, August 2002.
- [43] Eugene J. Lefferts, F. Landis Markley, and Malcolm D. Shuster. Kalman filtering for spacecraft attitude estimation. *Journal of Guidance, Control, and Dynamics*, 5(5):417–429, Sept.–Oct. 1982. doi:10.2514/3.56190.
- [44] M. Mandic. Distributed estimation architectures and algorithms for formation flying spacecraft. Master’s thesis, Massachusetts Institute of Technology, 2006.

- [45] S. R. Marandi and V. J. Modi. A preferred coordinate system and the associated orientation representation in attitude dynamics. *Acta Astronautica*, 15(11):833–843, Nov. 1987.
- [46] F. Landis Markley. Unit quaternion from rotation matrix. *Journal of Guidance, Control, and Dynamics*, 31(2):440–442, March-April 2008.
- [47] F. Landis Markley. Lessons learned. *The Journal of the Astronautical Sciences*, 57(1 & 2):3–29, Jan.-June 2009.
- [48] F. Landis Markley and J. Russell Carpenter. Generalized linear covariance analysis. *The Journal of the Astronautical Sciences*, 57(1 & 2):233–260, Jan.-June 2009.
- [49] F. Landis Markley and John L. Crassidis. *Fundamentals of Spacecraft Attitude Determination and Control*, pages 46, 257–260, 263–269. Springer, New York, NY, 2014. doi:10.1007/978-1-4939-0802-8.
- [50] Peter S. Maybeck. *Stochastic Models, Estimation and Control, Vol. 1*. Academic Press, New York, NY, 1979.
- [51] Peter S. Maybeck. *Stochastic Models, Estimation and Control, Vol. 2*. Academic Press, New York, NY, 1979.
- [52] Eugene S. Muller and Peter M. Kachmar. A new approach to on-board orbit navigation. *Navigation: Journal of the Institute of Navigation*, 18(4):369–385, Winter 1971–72.
- [53] Magnus Nørgaard, Niels K. Poulsen, and Ole Ravn. Advances in derivative-free state estimation for nonlinear systems. Technical Report IMM-REP-1998-15, Technical University of Denmark, 2800 Lyngby, Denmark, April 7, 2000.
- [54] Magnus Nørgaard, Niels K. Poulsen, and Ole Ravn. New developments in state estimation for nonlinear estimation problems. *Automatica*, 36(11):1627–1638, November 2000.
- [55] Young W. Park, Jack P. Brazzel, J. Russell Carpenter, Heather D. Hinkel, and James H. Newman. Flight test results from real-time relative gps experiment on sts-69. In *SPACEFLIGHT MECHANICS 1996*, volume 93 of *Advances in the Astronautical Sciences*, pages 1277–1296, San Diego, CA, 1996. Univelt.
- [56] L. Perea, J. How, L. Breger, and P. Elosegui. Nonlinearity in sensor fusion: Divergence issues in ekf, modified truncated sof, and ukf. In *AIAA Guidance, Navigation and Control Conference and Exhibit*, Hilton Head, SC, August 20–23, 2007.
- [57] Mark E. Pittelkau. An analysis of the quaternion attitude determination filter. *The Journal of the Astronautical Sciences*, 51(1), Jan.-March 2003.
- [58] H. Plinval. Analysis of relative navigation architectures for formation flying spacecraft. Master’s thesis, Massachusetts Institute of Technology, 2006.
- [59] H. E. Rauch, F. Tung, and C. T. Striebel. Maximum likelihood estimates of linear dynamic systems. *AIAA Journal*, 3(8):1445–1450, Aug. 1965.
- [60] Reid G. Reynolds. Asymptotically optimal attitude filtering with guaranteed coverage. *Journal of Guidance, Control, and Dynamics*, 31(1):114–122, Jan.-Feb. 2008.
- [61] Olinde Rodrigues. Des lois géométriques qui régissent les déplacements d’un système solide dans l’espace, et de la variation des coordonnées provenant de ces déplacements considérés indépendamment des causes qui peuvent les produire. *Journal de Mathématiques Pures et Appliquées*, 5:380–440, 1840.
- [62] Hanspeter Schaub and John L. Junkins. *Analytical Mechanics of Aerospace Systems*. American Institute of Aeronautics and Astronautics, Inc., New York, NY, 2nd edition, 2009.
- [63] Emil Schiesser, Jack P. Brazzel, J. Russell Carpenter, and Heather D. Hinkel. Results of sts-80 relative gps navigation flight experiment. In *SPACEFLIGHT MECHANICS 1998*, volume 99 of *Advances in the Astronautical Sciences*, pages 1317–1334, San Diego, CA, 1998. Univelt.
- [64] S. F. Schmidt. Application of state-space methods to navigation problems. In C. T. Leondes, editor, *Advances in Control Systems: Theory and Applications*, volume 3, pages 293–340. Academic Press, New York, 1966.
- [65] Stanley F. Schmidt. The kalman filter - its recognition and development for aerospace applications. *Journal of Guidance, Control, and Dynamics*, 4(1):4–7, 2016/01/09 1981.
- [66] John H. Seago, Jacob Griesback, James W. Woodburn, and David A. Vallado. Sequential orbit-estimation with sparse tracking. In *Space Flight Mechanics 2011*, volume 140 of *Advances in the Astronautical Sciences*, pages 281–299. Univelt, 2011.
- [67] Malcolm D. Shuster. A survey of attitude representations. *The Journal of the Astronautical Sciences*, 41(4):439–517, Oct.-Dec. 1993.

- [68] John Stuelpnagel. On the parametrization of the three-dimensional rotation group. *SIAM Review*, 6(4):422–430, Oct. 1964.
- [69] Byron D. Tapley, Bob E. Schutz, and George R. Born. *Statistical Orbit Determination*. Academic Press, 2004.
- [70] C.L Thornton and G.J. Bierman. Gram-schmidt algorithms for covariance propagation. *IEEE Conference on Decision and Control*, pages 489–498, 1975.
- [71] C.L Thornton and G.J. Bierman. Numerical comparison of discrete kalman filtering algorithms: An orbit determination case study. *JPL Technical Memorandum*, 33-771, June 1976.
- [72] Oldrich Vasicek. An equilibrium characterization of the term structure. *J. Financial Economics*, 5(2):177–188, November 1977.
- [73] M. C. Wang and G. E. Uhlenbeck. On the theory of brownian motion ii. In N. Wax, editor, *Selected Papers on Noise and Stochastic Processes*, pages 113–132. Dover, 1954.
- [74] James R. Wertz, editor. *Spacecraft Attitude Determination and Control*. Kluwer Academic Publishers, The Netherlands, 1978.
- [75] Thomas F. Wiener. *Theoretical Analysis of Gimballess Inertial Reference Equipment Using Delta-Modulated Instruments*. Ph.D. dissertation, Massachusetts Institute of Technology, Cambridge, MA, 1962.
- [76] J. R. Wright. Sequential orbit determination with auto-correlated gravity modeling errors. *Journal of Guidance and Control*, 4(3):304–309, May–June 1980.
- [77] James R. Wright. Optimal orbit determination. In Kyle T. Alfriend et al., editor, *Space Flight Mechanics 2002*, volume 112 of *Advances in the Astronautical Sciences*, pages 1123–134. Univelt, 2002.
- [78] Renato Zanetti, Kyle J. DeMars, and Robert H. Bishop. Underweighting nonlinear measurements. *Journal of Guidance, Control and Dynamics*, 33(5):1670–1675, September–October 2010.
- [79] Renato Zanetti and Christopher D’Souza. Recursive implementations of the consider filter. *Proceedings of the 2012 AAS Jer-Nan Juang Symposium*, 2012.

**REPORT DOCUMENTATION PAGE**

*Form Approved  
OMB No. 0704-0188*

The public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Department of Defense, Washington Headquarters Services, Directorate for Information Operations and Reports (0704-0188), 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

**PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.**

<b>1. REPORT DATE (DD-MM-YYYY)</b> 18-04-2018		<b>2. REPORT TYPE</b> Technical Publication		<b>3. DATES COVERED (From - To)</b>	
<b>4. TITLE AND SUBTITLE</b> Navigation Filter Best Practices				<b>5a. CONTRACT NUMBER</b>	
				<b>5b. GRANT NUMBER</b>	
				<b>5c. PROGRAM ELEMENT NUMBER</b>	
<b>6. AUTHOR(S)</b> Carpenter, J. R.; D'Souza, C. N.				<b>5d. PROJECT NUMBER</b>	
				<b>5e. TASK NUMBER</b>	
				<b>5f. WORK UNIT NUMBER</b> 869021.03.04.01.03	
<b>7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)</b> NASA NASA Engineering and Safety Center Hampton, Virginia 23681				<b>8. PERFORMING ORGANIZATION REPORT NUMBER</b> L-20926	
<b>9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)</b> National Aeronautics and Space Administration Washington, DC 20546-0001				<b>10. SPONSOR/MONITOR'S ACRONYM(S)</b> NASA	
				<b>11. SPONSOR/MONITOR'S REPORT NUMBER(S)</b> NASA/TP-2018-219822	
<b>12. DISTRIBUTION/AVAILABILITY STATEMENT</b> Unclassified-Unlimited Subject Category 17 (Space Communications, Spacecraft Communications, Command and Tracking: Navigation Systems) Availability: NASA STI Program (757) 864-9658					
<b>13. SUPPLEMENTARY NOTES</b> An electronic version can be found at <a href="http://ntrs.nasa.gov">http://ntrs.nasa.gov</a> .					
<b>14. ABSTRACT</b> This work identifies best practices for onboard navigation filtering algorithms. These best practices have been collected by NASA practitioners of the art and science of navigation over the first 50 years of the Space Age. While this is a NASA document concerned with space navigation, it is likely that many of the principles would apply equally to the wider navigation community.					
<b>15. SUBJECT TERMS</b> Navigation, Kalman Filter, Estimation					
<b>16. SECURITY CLASSIFICATION OF:</b>			<b>17. LIMITATION OF ABSTRACT</b>	<b>18. NUMBER OF PAGES</b>	<b>19a. NAME OF RESPONSIBLE PERSON</b>
<b>a. REPORT</b>	<b>b. ABSTRACT</b>	<b>c. THIS PAGE</b>			STI Information Desk ( <a href="mailto:help@sti.nasa.gov">help@sti.nasa.gov</a> )
U	U	U	UU	149	<b>19b. TELEPHONE NUMBER (Include area code)</b> (757) 864-9658