

# Kalman and Bayesian Filters in Python

Roger R Labbe Jr

January 3, 2016



# Contents

<b>Preface</b>	<b>11</b>
0.1 Kalman and Bayesian Filters . . . . .	11
0.2 Motivation for this Book . . . . .	12
0.3 Reading Online . . . . .	13
0.4 PDF Version . . . . .	13
0.5 Downloading the book . . . . .	14
0.6 Installation and Software Requirements . . . . .	14
0.7 My Libraries and Modules . . . . .	15
0.8 Thoughts on Python and Coding Math . . . . .	15
0.9 License . . . . .	16
0.10 Contact . . . . .	16
0.11 Resources . . . . .	16
<b>1 The g-h Filter</b>	<b>17</b>
1.1 Building Intuition via Thought Experiments . . . . .	17
1.2 The g-h Filter . . . . .	27
1.3 Notation . . . . .	30
1.4 NumPy arrays . . . . .	30
1.4.1 Exercise - Create arrays . . . . .	33
1.4.2 Solution . . . . .	33
1.5 Exercise: Write Generic Algorithm . . . . .	34
1.5.1 Solution and Discussion . . . . .	34
1.6 Choice of g and h . . . . .	35
1.7 Exercise: create measurement function . . . . .	35
1.7.1 Solution . . . . .	36
1.7.2 Discussion . . . . .	36
1.8 Exercise: Bad Initial Conditions . . . . .	36
1.8.1 Solution and Discussion . . . . .	37
1.9 Exercise: Extreme Noise . . . . .	37
1.9.1 Solution and Discussion . . . . .	37
1.10 Exercise: The Effect of Acceleration . . . . .	38
1.10.1 Solution and Discussion . . . . .	38
1.11 Exercise: Varying g . . . . .	39
1.11.1 Solution and Discussion . . . . .	39
1.12 Varying h . . . . .	41
1.13 Interactive Example . . . . .	42
1.14 Don't Lie to the Filter . . . . .	43
1.15 Tracking a Train . . . . .	44
1.16 g-h Filters with FilterPy . . . . .	48
1.17 Summary . . . . .	49
1.18 References . . . . .	50

<b>2 Discrete Bayes Filter</b>	<b>51</b>
2.1 Tracking a Dog . . . . .	51
2.2 Extracting Information from Multiple Sensor Readings . . . . .	53
2.3 Noisy Sensors . . . . .	54
2.4 Incorporating Movement Data . . . . .	56
2.5 Terminology . . . . .	57
2.6 Adding Noise to the Prediction . . . . .	57
2.7 Generalizing with Convolution . . . . .	60
2.8 Integrating Measurements and Movement Updates . . . . .	62
2.9 The Effect of Bad Sensor Data . . . . .	65
2.10 Drawbacks and Limitations . . . . .	67
2.11 Tracking and Control . . . . .	68
2.11.1 Simulating the Train Behavior . . . . .	69
2.12 Bayes Theorem . . . . .	72
2.13 Total Probability Theorem . . . . .	73
2.14 Summary . . . . .	74
2.15 References . . . . .	75
<b>3 Gaussian Probabilities</b>	<b>77</b>
3.1 Introduction . . . . .	77
3.2 Mean, Variance, and Standard Deviations . . . . .	77
3.2.1 Random Variables . . . . .	77
3.3 Probability Distribution . . . . .	78
3.3.1 The Mean, Median, and Mode of a Random Variable . . . . .	79
3.4 Expected Value of a Random Variable . . . . .	80
3.4.1 Variance of a Random Variable . . . . .	80
3.4.2 Why the Square of the Differences . . . . .	84
3.5 Gaussians . . . . .	85
3.6 Nomenclature . . . . .	86
3.7 Gaussian Distributions . . . . .	87
3.8 The Variance and Belief . . . . .	89
3.9 The 68-95-99.7 Rule . . . . .	90
3.10 Interactive Gaussians . . . . .	91
3.11 Computational Properties of the Gaussian . . . . .	91
3.12 Computing Probabilities with <code>scipy.stats</code> . . . . .	92
3.13 Fat Tails . . . . .	93
3.14 Summary and Key Points . . . . .	96
3.15 References . . . . .	96
<b>4 One Dimensional Kalman Filters</b>	<b>97</b>
4.1 One Dimensional Kalman Filters . . . . .	97
4.2 Tracking A Dog . . . . .	97
4.2.1 Simulating the Sensor Measurement . . . . .	97
4.2.2 Simulating the Dog's movement . . . . .	98
4.2.3 Implementation in Python . . . . .	98
4.3 Math with Gaussians . . . . .	103
4.3.1 Interactive Example . . . . .	108
4.4 Implementing the Update Step . . . . .	109
4.5 Implementing Predictions . . . . .	111
4.5.1 Exercise: Modify Variance Values . . . . .	117
4.5.2 Tracking Animation . . . . .	117
4.6 The Kalman Filter is a g-h Filter . . . . .	117
4.7 Full Description of the Algorithm . . . . .	119
4.8 Comparison with g-h Filter and discrete Bayes Filter . . . . .	121

4.9	Implementation in a Class . . . . .	123
4.10	Introduction to Designing a Filter . . . . .	124
4.10.1	Animation . . . . .	126
4.11	Exercise: Plot Histogram of Measurements . . . . .	127
4.11.1	Solution . . . . .	127
4.11.2	Discussion . . . . .	128
4.12	Example: Extreme Amounts of Noise . . . . .	128
4.13	Example: Incorrect Process Variance . . . . .	129
4.14	Example: Bad Initial Estimate . . . . .	130
4.15	Example: Large Noise and Bad Initial Estimate . . . . .	131
4.16	Exercise: Interactive Plots . . . . .	133
4.16.1	Solution . . . . .	133
4.17	Exercise - Nonlinear Systems . . . . .	134
4.17.1	Solution . . . . .	135
4.17.2	Discussion . . . . .	135
4.18	Exercise - Noisy Nonlinear Systems . . . . .	136
4.18.1	Solution . . . . .	136
4.18.2	Discussion . . . . .	136
4.19	Fixed Gain Filters . . . . .	137
4.20	Derivation From Bayes Theorem (Optional) . . . . .	137
4.21	Summary . . . . .	138
<b>5</b>	<b>Multivariate Gaussians - Modeling Uncertainty in Multiple Dimensions</b>	<b>141</b>
5.1	Introduction . . . . .	141
5.2	Multivariate Normal Distributions . . . . .	141
5.3	Correlation and Covariance . . . . .	142
5.4	Multivariate Normal Distribution Equation . . . . .	146
5.4.1	Pearson's Correlation Coefficient . . . . .	154
5.5	Using Correlations to Improve Estimates . . . . .	155
5.6	Multiplying Multidimensional Gaussians . . . . .	156
5.7	Hidden Variables . . . . .	161
5.8	Summary . . . . .	165
5.9	References . . . . .	165
<b>6</b>	<b>Multivariate Kalman Filters - Working with Multiple State Variables</b>	<b>167</b>
6.1	Introduction . . . . .	167
6.2	Newton's Equations of Motion . . . . .	167
6.3	Kalman Filter Algorithm . . . . .	168
6.4	Tracking a Dog . . . . .	169
6.4.1	Initialization - Choose State Variables and Set Initial Conditions . . . . .	170
6.4.2	Predict Step . . . . .	172
6.4.3	Update Step . . . . .	176
6.4.4	Implementing the Kalman Filter . . . . .	179
6.5	The Kalman Filter Equations . . . . .	183
6.5.1	Prediction Equations . . . . .	183
6.5.2	Update Equations . . . . .	186
6.5.3	An Example not using FilterPy . . . . .	188
6.5.4	Summary . . . . .	190
6.6	Exercise: Show Effect of Hidden Variables . . . . .	192
6.6.1	Solution . . . . .	192
6.6.2	Discussion . . . . .	193
6.7	How Velocity is Calculated . . . . .	194
6.8	Adjusting the Filter . . . . .	195
6.9	A Detailed Examination of the Covariance Matrix . . . . .	199

6.10 Question: Explain Ellipse Differences . . . . .	203
6.10.1 Solution . . . . .	203
6.11 Filter Initialization . . . . .	205
6.12 Batch Processing . . . . .	206
6.13 Smoothing the Results . . . . .	207
6.14 Exercise: Compare Velocities . . . . .	208
6.14.1 Solution . . . . .	208
6.15 Discussion and Summary . . . . .	209
6.16 References . . . . .	210
<b>7 Kalman Filter Math</b> . . . . .	<b>211</b>
7.1 Modeling a Dynamic System . . . . .	211
7.2 State-Space Representation of Dynamic Systems . . . . .	213
7.2.1 Forming First Order Equations from Higher Order Equations . . . . .	214
7.2.2 First Order Differential Equations In State-Space Form . . . . .	214
7.2.3 Finding the Fundamental Matrix for Time Invariant Systems . . . . .	215
7.2.4 The Matrix Exponential . . . . .	215
7.2.5 Time Invariance . . . . .	216
7.2.6 Linear Time Invariant Theory . . . . .	218
7.2.7 Numerical Solutions . . . . .	218
7.3 Design of the Process Noise Matrix . . . . .	219
7.3.1 Continuous White Noise Model . . . . .	219
7.3.2 Piecewise White Noise Model . . . . .	221
7.3.3 Using FilterPy to Compute Q . . . . .	223
7.3.4 Simplification of Q . . . . .	224
7.4 Numeric Integration of Differential Equations . . . . .	224
7.4.1 Euler's Method . . . . .	225
7.4.2 Runge Kutta Methods . . . . .	227
7.5 Bayesian Filtering . . . . .	229
7.6 Converting the Multivariate Equations to the Univariate Case . . . . .	231
7.7 Converting Kalman Filter to a g-h Filter . . . . .	233
7.8 References . . . . .	235
<b>8 Designing Kalman Filters</b> . . . . .	<b>237</b>
8.1 Introduction . . . . .	237
8.2 Tracking a Robot . . . . .	237
8.2.1 Step 1: Choose the State Variables . . . . .	238
8.2.2 <b>Step 2:</b> Design State Transition Function . . . . .	239
8.2.3 Step 3: Design the Process Noise Matrix . . . . .	240
8.2.4 <b>Step 4:</b> Design the Control Function . . . . .	240
8.2.5 <b>Step 5:</b> Design the Measurement Function . . . . .	240
8.2.6 <b>Step 6:</b> Design the Measurement Noise Matrix . . . . .	241
8.2.7 Initial Conditions . . . . .	241
8.2.8 Implement the Filter . . . . .	242
8.3 The Effect of Filter Order . . . . .	245
8.3.1 Zero Order Kalman Filter . . . . .	247
8.3.2 First Order Kalman Filter . . . . .	247
8.3.3 Second Order Kalman Filter . . . . .	248
8.3.4 Evaluating the Performance . . . . .	249
8.4 Evaluating Filter Performance . . . . .	262
8.4.1 Normalized Estimated Error Squared (NEES) . . . . .	262
8.4.2 Likelihood Function . . . . .	264
8.4.3 NIS . . . . .	265
8.5 Control Inputs . . . . .	266

8.6	Sensor Fusion . . . . .	267
8.6.1	Exercise: Can you Kalman Filter GPS outputs? . . . . .	272
8.6.2	Exercise: Prove that the position sensor improves the filter . . . . .	275
8.7	Nonstationary Processes . . . . .	277
8.7.1	Sensor fusion: Different Data Rates . . . . .	278
8.8	Tracking a Ball . . . . .	280
8.8.1	Step 1: Choose the State Variables . . . . .	282
8.8.2	<b>Step 2:</b> Design State Transition Function . . . . .	283
8.8.3	<b>Step 3:</b> Design the Control Input Function . . . . .	283
8.8.4	<b>Step 4:</b> Design the Measurement Function . . . . .	284
8.8.5	<b>Step 5:</b> Design the Measurement Noise Matrix . . . . .	285
8.8.6	Step 6: Design the Process Noise Matrix . . . . .	285
8.8.7	Step 7: Design the Initial Conditions . . . . .	285
8.9	Tracking a Ball in Air . . . . .	287
8.9.1	Implementing Air Drag . . . . .	287
8.10	References . . . . .	294
<b>9</b>	<b>Nonlinear Filtering</b>	<b>295</b>
9.1	Introduction . . . . .	295
9.2	The Problem with Nonlinearity . . . . .	296
9.3	An Intuitive Look at the Problem . . . . .	297
9.4	The Effect of Nonlinear Functions on Gaussians . . . . .	298
9.5	A 2D Example . . . . .	303
9.6	The Algorithms . . . . .	305
9.7	Summary . . . . .	306
9.8	References . . . . .	307
<b>10</b>	<b>The Unscented Kalman Filter</b>	<b>309</b>
10.1	Sigma Points - Sampling from a Distribution . . . . .	310
10.2	Choosing Sigma Points . . . . .	312
10.3	The Unscented Transform . . . . .	314
10.3.1	Unscented Transform Example . . . . .	315
10.4	The Unscented Kalman Filter . . . . .	316
10.4.1	Predict Step . . . . .	317
10.4.2	Update Step . . . . .	317
10.5	Van der Merwe's Scaled Sigma Point Algorithm . . . . .	318
10.5.1	Sigma Point Computation . . . . .	319
10.5.2	Weight Computation . . . . .	319
10.5.3	Reasonable Choices for the Parameters . . . . .	320
10.6	Using the UKF . . . . .	320
10.7	Tracking a Flying Airplane . . . . .	323
10.7.1	Tracking Manuevering Aircraft . . . . .	328
10.7.2	Sensor Fusion . . . . .	330
10.7.3	Multiple Position Sensors . . . . .	333
10.8	Effects of Sensor Error and Geometry . . . . .	336
10.9	Exercise: Explain Filter Performance . . . . .	336
10.9.1	Solution . . . . .	336
10.10	Implementation of the UKF . . . . .	339
10.10.1	Weights . . . . .	340
10.10.2	Sigma Points . . . . .	340
10.10.3	Predict Step . . . . .	342
10.10.4	Update Step . . . . .	342
10.10.5	FilterPy's Implementation . . . . .	343
10.11	Batch Processing . . . . .	343

10.12 Smoothing the Results . . . . .	345
10.13 Choosing the Sigma Parameters . . . . .	347
10.14 Robot Localization - A Fully Worked Example . . . . .	349
10.14.1 Robot Motion Model . . . . .	349
10.14.2 Design the State Variables . . . . .	350
10.14.3 Design the System Model . . . . .	350
10.14.4 Design the Measurement Model . . . . .	351
10.14.5 Design Measurement Noise . . . . .	352
10.14.6 Implementation . . . . .	352
10.14.7 Steering the Robot . . . . .	356
10.15 Discussion . . . . .	357
10.16 References . . . . .	358
<b>11 The Extended Kalman Filter</b>	<b>359</b>
11.1 Linearizing the Kalman Filter . . . . .	359
11.2 Example: Tracking a Flying Airplane . . . . .	361
11.2.1 Design the State Variables . . . . .	362
11.2.2 Design the Process Model . . . . .	362
11.2.3 Design the Measurement Model . . . . .	363
11.2.4 Design Process and Measurement Noise . . . . .	365
11.2.5 Implementation . . . . .	366
11.3 Using SymPy to compute Jacobians . . . . .	368
11.4 Robot Localization . . . . .	368
11.4.1 Robot Motion Model . . . . .	369
11.4.2 Design the State Variables . . . . .	370
11.4.3 Design the System Model . . . . .	370
11.4.4 Design the Measurement Model . . . . .	372
11.4.5 Design Measurement Noise . . . . .	373
11.4.6 Implementation . . . . .	373
11.4.7 Discussion . . . . .	379
11.5 UKF vs EKF . . . . .	380
<b>12 Particle Filters</b>	<b>383</b>
12.1 Motivation . . . . .	383
12.2 Monte Carlo Sampling . . . . .	384
12.3 Generic Particle Filter Algorithm . . . . .	384
12.4 Probability distributions via Monte Carlo . . . . .	386
12.5 The Particle Filter . . . . .	388
12.5.1 Predict Step . . . . .	389
12.5.2 Update Step . . . . .	390
12.5.3 Computing the State Estimate . . . . .	391
12.5.4 Particle Resampling . . . . .	391
12.6 SIR Filter - A Complete Example . . . . .	392
12.6.1 Effect of Sensor Errors on the Filter . . . . .	395
12.6.2 Filter Degeneracy From Inadequate Samples . . . . .	396
12.7 Importance Sampling . . . . .	398
12.8 Resampling Methods . . . . .	399
12.8.1 Multinomial Resampling . . . . .	400
12.8.2 Residual Resampling . . . . .	401
12.8.3 Stratified Resampling . . . . .	402
12.8.4 Systematic Resampling . . . . .	403
12.8.5 Choosing a Resampling Algorithm . . . . .	403
12.9 Summary . . . . .	405
12.10 References . . . . .	405

<b>13 Smoothing</b>	<b>407</b>
13.1 Introduction . . . . .	407
13.2 An Overview of How Smoothers Work . . . . .	408
13.3 Types of Smoothers . . . . .	409
13.4 Fixed-Point Smoothing . . . . .	410
13.5 Fixed-Interval Smoothing . . . . .	410
13.6 Fixed-Lag Smoothing . . . . .	413
13.7 References . . . . .	416
<b>14 Adaptive Filtering</b>	<b>417</b>
14.1 Introduction . . . . .	417
14.2 Maneuvering Targets . . . . .	417
14.3 Detecting a Maneuver . . . . .	424
14.4 Adjustable Process Noise . . . . .	426
14.4.1 Continuous Adjustment . . . . .	427
14.4.2 Continuous Adjustment - Standard Deviation Version . . . . .	429
14.5 Fading Memory Filter . . . . .	434
14.6 Noise Level Switching . . . . .	437
14.7 Variable State Dimension . . . . .	437
14.8 Multiple Model Estimation . . . . .	437
14.8.1 A Two Filter Adaptive Filter . . . . .	438
14.9 MMAE . . . . .	440
14.9.1 Limitations of the MMAE Filter . . . . .	442
14.10 Interacting Multiple Models (IMM) . . . . .	443
14.11 References . . . . .	444
<b>A Installation, Python, NumPy, and FilterPy</b>	<b>445</b>
A.1 Installing the SciPy Stack . . . . .	445
A.2 Installing FilterPy . . . . .	446
A.2.1 Manual Install of the SciPy stack . . . . .	446
A.3 Installing/downloading and running the book . . . . .	447
A.4 Using IPython Notebook . . . . .	448
A.5 SymPy . . . . .	448
<b>B Symbols and Notations</b>	<b>451</b>
B.1 Labbe . . . . .	451
B.2 Wikipedia . . . . .	451
B.3 Brookner . . . . .	452
B.4 Gelb . . . . .	452
B.5 Brown . . . . .	452
B.6 Zarchan . . . . .	452
B.7 Bar-Shalom . . . . .	452
B.8 Thrun . . . . .	452
<b>C H Infinity filter</b>	<b>453</b>
<b>D Ensemble Kalman Filters</b>	<b>455</b>
D.1 The Algorithm . . . . .	456
D.1.1 Initialize Step . . . . .	457
D.1.2 Predict Step . . . . .	457
D.1.3 Update Step . . . . .	457
D.2 Implementation and Example . . . . .	459
D.3 Outstanding Questions . . . . .	462
D.4 References . . . . .	463



# Preface

Introductory textbook for Kalman filters and Bayesian filters. The book is written using Jupyter Notebook so that read the book in your browser and also run and modify the code, seeing the results inside the book. What better way to learn?

## 0.1 Kalman and Bayesian Filters

Sensors are noisy. The world is full of data and events that we want to measure and track, but we cannot rely on sensors to give us perfect information. The GPS in my car reports altitude. Each time I pass the same point in the road it reports a slightly different altitude. My kitchen scale gives me different readings if I weigh the same object twice.

In simple cases the solution is obvious. If my scale gives slightly different readings I can just take a few readings and average them. Or I can replace it with a more accurate scale. But what do we do when the sensor is very noisy, or the environment makes data collection difficult? We may be trying to track the movement of a low flying aircraft. We may want to create an autopilot for a drone, or ensure that our farm tractor seeded the entire field. I do computer vision, and I need to track moving objects in images, and the computer vision algorithms create very noisy and unreliable results.

This book teaches you how to solve these sorts of filtering problems. I use many different algorithms, but they are all based on Bayesian probability. In simple terms Bayesian probability determines what is likely to be true based on past information.

If I asked you the heading of my car at this moment you would have no idea. You'd proffer a number between  $1^\circ$  and  $360^\circ$  degrees, and have a 1 in 360 chance of being right. Now suppose I told you that 2 seconds ago its heading was  $243^\circ$ . In 2 seconds my car could not turn very far so you could make a far more accurate prediction. You are using past information to more accurately infer information about the present or future.

The world is also noisy. That prediction helps you make a better estimate, but it also subject to noise. I may have just braked for a dog or swerved around a pothole. Strong winds and ice on the road are external influences on the path of my car. In control literature we call this noise though you may not think of it that way.

There is more to Bayesian probability, but you have the main idea. Knowledge is uncertain, and we alter our beliefs based on the strength of the evidence. Kalman and Bayesian filters blend our noisy and limited knowledge of how a system behaves with the noisy and limited sensor readings to produce the best possible estimate of the state of the system. Our principle is to never discard information.

Say we are tracking an object and a sensor reports that it suddenly changed direction. Did it really turn, or is the data noisy? It depends. If this is a jet fighter we'd very inclined to believe the report of a sudden maneuver. If it is a freight train on a straight track we would discount it. We'd further modify our belief depending on how accurate the sensor is. Our beliefs depend on the past and on our knowledge of the system we are tracking and on the characteristics of the sensors.

The Kalman filter was invented by Rudolf Emil Kálmán to solve this sort of problem in a mathematically optimal way. Its first use was on the Apollo missions to the moon, and since then it has been used in an

enormous variety of domains. There are Kalman filters in aircraft, on submarines, on cruise missiles. Wall street uses them to track the market. They are used in robots, in IoT (Internet of Things) sensors, and in laboratory instruments. Chemical plants use them to control and monitor reactions. They are used in medical equipment to do things like medical imaging or to remove noise from cardiac signals. If it involves a sensor and/or time-series data, a Kalman filter or a close relative to the Kalman filter is usually involved.

## 0.2 Motivation for this Book

I'm a software engineer that spent almost two decades in aerospace, and so I have always been 'bumping elbows' with the Kalman filter, but never implemented one. They've always had a fearsome reputation for difficulty. The theory is beautiful, but quite difficult to learn if you are not already well trained in topics such as signal processing, control theory, probability and statistics, and guidance and control theory. As I moved into solving tracking problems with computer vision the need to implement them myself became urgent.

There are excellent textbooks in the field, such as Grewal and Andrew's Kalman Filtering. But sitting down and trying to read many of these books is a dismal and trying experience if you do not have the necessary background. Typically the first few chapters fly through several years of undergraduate math, blithely referring you to textbooks on Itō calculus, and presenting an entire semester's worth of statistics in a few brief paragraphs. They are textbooks for an upper undergraduate or graduate level course, and an invaluable reference to researchers and professionals, but the going is truly difficult for the more casual reader. Notation is introduced without explanation, different texts use different words and variables names for the same concept, and the books are almost devoid of examples or worked problems. I often found myself able to parse the words and comprehend the mathematics of a definition, but had no idea as to what real world phenomena these words and math were attempting to describe. "But what does that mean?" was my repeated thought. Here are typical examples which once puzzled me:

$$\begin{aligned}\hat{x}_k &= \Phi_k \hat{x}_{k-1} + G_k u_{k-1} + K_k [z_k - H\Phi_k \hat{x}_{k-1} - HG_k u_{k-1}] \\ \mathbf{P}_{k|k} &= (I - \mathbf{K}_k \mathbf{H}_k) \text{cov}(\mathbf{x}_k - \hat{\mathbf{x}}_{k|k-1}) (I - \mathbf{K}_k \mathbf{H}_k)^T + \mathbf{K}_k \text{cov}(\mathbf{v}_k) \mathbf{K}_k^T\end{aligned}$$

However, as I began to finally understand the Kalman filter I realized the underlying concepts are quite straightforward. If you know a few simple probability rules, and have some intuition about how we fuse uncertain knowledge the concepts of the Kalman filter are accessible. Kalman filters have a reputation for difficulty, but shorn of much of the formal terminology the beauty of the subject and of their math became clear to me, and I fell in love with the topic.

As I began to understand the math and theory more difficulties appeared. A book or paper will make some statement of fact and presents a graph as proof. Unfortunately, why the statement is true is not clear to me, or I cannot reproduce the plot. Or maybe I wonder "is this true if R=0?" Or the author provides pseudocode at such a high level that the implementation is not obvious. Some books offer Matlab code, but I do not have a license to that expensive package. Finally, many books end each chapter with many useful exercises. Exercises which you need to understand if you want to implement Kalman filters for yourself, but exercises with no answers. If you are using the book in a classroom, perhaps this is okay, but it is terrible for the independent reader. I loathe that an author withholds information from me, presumably to avoid 'cheating' by the student in the classroom.

All of this impedes learning. If you are designing a Kalman filter for a aircraft or missile you must thoroughly master of all of the mathematics and topics in a typical Kalman filter textbook. I just want to track an image on a screen, or write some code for my Arduino project. I want to know how the plots in the book are made, and chose different parameters than the author chose. I want to run simulations. I want to inject more noise in the signal and see how a filter performs. There are thousands of opportunities for using Kalman filters in everyday code, and yet this fairly straightforward topic is the provenance of rocket scientists and academics.

I wrote this book to address all of those needs. This is not the book for you if you design military radars. Go get a Masters or PhD at a great STEM school, because you'll need it. This book is for the hobbyist, the

curious, and the working engineer that needs to filter or smooth data.

This book is interactive. While you can read it online as static content, I urge you to use it as intended. It is written using Jupyter Notebook (formally known as IPython Notebook). This allows me to combine text, math, Python, and Python output in one place. Every plot, every piece of data in this book is generated from Python inside the notebook. Want to double the value of a parameter? Just change the parameter's value, and press CTRL-ENTER. A new plot or printed output will appear.

This book has exercises, but it also has the answers. I trust you. If you just need an answer, go ahead and read the answer. If you want to internalize this knowledge, try to implement the exercise before you read the answer. Since the book is interactive, you enter and run your solution inside the book - you don't have to move to a different environment, or deal with importing a bunch of stuff before starting.

If you are a hobbyist this book should provide everything you need. If you are serious about Kalman filters you'll need more. My intention is to introduce enough of the concepts and mathematics to make the textbooks and papers approachable.

This book is free. I've spent several thousand dollars on Kalman filtering books. I cannot believe they are within the reach of someone in a depressed economy or a financially struggling student. I have gained so much from free software like Python, and free books like those from Allen B. Downey [1]. It's time to repay that. So, the book is free, it is hosted on free servers at GitHub, and it uses only free and open software such as IPython and MathJax.

## 0.3 Reading Online

GitHub

The book is hosted on GitHub, and you can read any chapter by clicking on its name. GitHub statically renders Jupyter Notebooks. You will not be able to run or alter the code, but you can read all of the content.

The GitHub pages for this project are at

<https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python>

binder

I am experimentally trying a new service, binder. binder serves interactive notebooks online, so you can run the code and change the code within your browser without downloading the book or installing Jupyter. I have not experimented with it much. Please raise an issue on my GitHub page if anything fails. I'm not officially supporting this as binder is in beta, but I'm quite excited about the possibilities.

Use this link to access the book via binder:

<http://mybinder.org/repo/rabbe/Kalman-and-Bayesian-Filters-in-Python>

nbviewer

The nbviewer website will render any Notebook in a static format. I find it does a slightly better job than the GitHub renderer, but it is slightly harder to use. It accesses GitHub directly; whatever I have checked into GitHub will be rendered by nbviewer.

You may access this book via nbviewer here:

[http://nbviewer.ipython.org/github/rabbe/Kalman-and-Bayesian-Filters-in-Python/blob/master/table\\_of\\_contents.ipynb](http://nbviewer.ipython.org/github/rabbe/Kalman-and-Bayesian-Filters-in-Python/blob/master/table_of_contents.ipynb)

## 0.4 PDF Version

I periodically generate a PDF of the book from the notebooks. You can access it here:

[https://drive.google.com/file/d/0By\\_SW19c1BfhSVFzNHc0SjduNzg/view?usp=sharing](https://drive.google.com/file/d/0By_SW19c1BfhSVFzNHc0SjduNzg/view?usp=sharing)

## 0.5 Downloading the book

This book is interactive and I recommend using it in that form. If you install IPython on your computer and then clone this book you will be able to run all of the code in the book inside your browser. You can perform experiments, see how filters react to different data, see how different filters react to the same data, and so on. I find this sort of immediate feedback vital. You do not have to wonder “what happens if”. Try it and see!

The GitHub pages for this project are at

```
https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python
```

You can clone it to your hard drive with the command

```
git clone --depth=1 https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python.git
```

This will create a directory named `Kalman-and-Bayesian-Filters-in-Python`. The `depth` parameter just gets you the latest version. Unless you need to see my entire commit history this is a lot faster and saves space. Navigate to the directory, and run Jupyter Notebook with the command

```
ipython notebook
```

This will open a browser window showing the contents of the base directory. The book is organized into chapters. Each chapter is named `xx-name.ipynb`, where `xx` is the chapter number. `.ipynb` is the Notebook file extension.

Admittedly this is a cumbersome interface to a book. I am following in the footsteps of several other projects that are re-purposing Jupyter Notebook to generate entire books. I feel the slight annoyances have a huge payoff - instead of having to download a separate code base and run it in an IDE while you try to read a book, all of the code and text is in one place. If you want to alter the code, you may do so and immediately see the effects of your change. If you find a bug, you can make a fix, and push it back to my repository so that everyone in the world benefits. And, of course, you will never encounter a problem I face all the time with traditional books - the book and the code are out of sync with each other, and you are left scratching your head as to which source to trust.

## 0.6 Installation and Software Requirements

Installation of the required software is described in detail in the Installation appendix, which you can read online here:

```
http://nbviewer.ipython.org/github/rlabbe/Kalman-and-Bayesian-Filters-in-Python/blob/master/Appendix_A_Installation.ipynb
```

You will need IPython 3.0 or greater, Python 2.7+ or Python 3.4+, NumPy, SciPy, SymPy, and Matplotlib. You will also need my open source library FilterPy.

The easiest way to get all of those except for FilterPy is by installing a free Scientific Python distribution, such as the Anaconda distribution <https://store.continuum.io/cshop/anaconda/> (the link says “shop”, but it is free). I am not naming them out of favoritism, I am merely documenting my environment. Should you have trouble running any of the code, perhaps knowing this will help you.

Once you have Python installed install FilterPy by typing

```
pip install filterpy
```

at the command line.

## 0.7 My Libraries and Modules

I am writing an open source Bayesian filtering Python library called FilterPy. Installation instructions are given above.

FilterPy is hosted GitHub at (<https://github.com/rlabbe/filterpy>) but the `pip` installed version should serve your needs.

Code that is specific to the book is stored with the book in the subdirectory `/code`. It contains code for formatting the book. It also contains python files with names like `xxx_internal.py`. I use these to store functions that are useful for a specific chapter. This allows me to hide Python code that is not particularly interesting to read - I may be generating a plot or chart, and I want you to focus on the contents of the chart, not the mechanics of how I generate that chart with Python. If you are curious as to the mechanics of that, just go and browse the source.

Some chapters introduce functions that are useful for the rest of the book. Those functions are initially defined within the Notebook itself, but the code is also stored in a Python file in `/code` that is imported if needed in later chapters. I do document when I do this where the function is first defined, but this is still a work in progress. I try to avoid this because then I always face the issue of code in the directory becoming out of sync with the code in the book. However, Jupyter Notebook does not give us a way to refer to code cells in other notebooks, so this is the only mechanism I know of to share functionality across notebooks.

There is an undocumented directory called `/experiments`. This is where I write and test code prior to putting it in the book. There is some interesting stuff in there, and feel free to look at it. As the book evolves I plan to create examples and projects, and a lot of this material will end up there. Small experiments will eventually just be deleted. If you are just interested in reading the book you can safely ignore this directory.

The directory `/code` contains a css file containing the style guide for the book. The default look and feel of Jupyter Notebook is rather plain. I have followed the examples set by books such as [Probabilistic Programming and Bayesian Methods for Hackers](#) [4]. I have also been very influenced by Professor Lorena Barba's fantastic work, [available here](#) [5]. I owe all of my look and feel to the work of these projects.

## 0.8 Thoughts on Python and Coding Math

Most Kalman filtering and other engineering texts are written by mathematicians or academics. When there is software(rarely), it is not production quality. Take Paul Zarchan's book [Fundamentals of Kalman Filtering](#) as an example. This is a fantastic book and it belongs in your library, and is one of the few books that provides full source for every example and chart. But the code is Fortran without any subroutines beyond calls to functions like `MATMUL`. Kalman filters are re-implemented throughout the book. The same listing mixes simulation with filtering code, making it hard to distinguish them. Some chapters implement the same filter in subtly different ways, and uses bold text to highlight the few lines that changed. If Runge Kutta is needed it is embedded in the code, without comments.

There's a better way. If I want to perform Runge Kutta I call `ode45`, I do not embed an Runge Kutta implementation in my code. I don't want to implement Runge Kutta multiple times and debug it several times. If I do find a bug, I can fix it once and be assured that it now works across all my different projects. And, it is readable. It is rare that I care about the implementation of Runge Kutta.

This is a textbook on Kalman filtering, and you can argue that we do care about the implementation of Kalman filters. That is true, but you will find out the code that performs the filtering amounts to 10 or so lines of code. The code to implement the math is fairly trivial. Most of the work that Kalman filter requires is the design of the matrices that get fed into the math engine.

A possible downside is that the equations that perform the filtering are hidden behind functions, which we could argue is a loss in a pedagogical text. I argue the converse. I want you to learn how to use Kalman filters in the real world, for real projects, and you shouldn't be cutting and pasting established algorithms all over the place.

I use Python classes. I mostly use classes as a way to organize the data that the filters require, not to implement OO features such as inheritance. For example, the `KalmanFilter` class stores matrices called `x`, `P`, `R`, `Q`. I've seen procedural libraries for Kalman filters, and they require the programmer to maintain all of those matrices. This perhaps isn't so bad for a toy program, but program a bank of Kalman filters and you will not enjoy having to manage all of those matrices and other associated data. I have derived from these classes occasionally in my own work, and find it handy, but I don't want to force OO on people as I know many do not like it.

## 0.9 License

Kalman Filters and Random Signals in Python by Roger Labbe is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

<http://creativecommons.org/licenses/by-nc-sa/4.0/>

Based on the work at <https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python>.

## 0.10 Contact

[rlabbejr@gmail.com](mailto:rlabbejr@gmail.com)

## 0.11 Resources

- [1] <http://www.greenteapress.com/>
- [2] <http://ipython.org/ipython-doc/rel-1.0.0/interactive/nbconvert.html>
- [3] <https://store.continuum.io/cshop/anaconda/>
- [4] [http://nbviewer.ipython.org/github/CamDavidsonPilon/Probabilistic-Programming-and-Bayesian-Methods-for-Hackers/blob/master/Chapter1\\_Introduction/Chapter1.ipynb](http://nbviewer.ipython.org/github/CamDavidsonPilon/Probabilistic-Programming-and-Bayesian-Methods-for-Hackers/blob/master/Chapter1_Introduction/Chapter1.ipynb)
- [5] <https://github.com/barbagroup/CFDPython>

# Chapter 1

## The g-h Filter

### 1.1 Building Intuition via Thought Experiments

Imagine that we live in a world without scales - the devices you stand on to weigh yourself. One day at work a coworker comes running up to you and announces her invention of a ‘scale’ to you. After she explains, you eagerly stand on it and announce the results: “172 lbs”. You are ecstatic - for the first time in your life you know what you weigh. More importantly, dollar signs dance in your eyes as you imagine selling this device to weight loss clinics across the world! This is fantastic!

Another coworker hears the commotion and comes over to find out what has you so excited. You explain the invention and once again step onto the scale, and proudly proclaim the result: “161 lbs.” And then you hesitate, confused.

“It read 172 lbs a few seconds ago” you complain to your coworker.

“I never said it was accurate,” she replies.

Sensors are inaccurate. This is the motivation behind a huge body of work in filtering, and solving this problem is the topic of this book. I could just provide the solutions that have been developed over the last half century, but these solutions developed by asking very basic, fundamental questions into the nature of what we know and how we know it. Before we attempt the math, let’s follow that journey of discovery, and see if it does not inform our intuition about filtering.

#### Try Another Scale

Is there any way we can improve upon this result? The obvious, first thing to try is get a better sensor. Unfortunately, your co-worker informs you that she has built 10 scales, and they all operate with about the same accuracy. You have her bring out another scale, and you weigh yourself on one, and then on the other. The first scale (A) reads “160 lbs”, and the second (B) reads “170 lbs”. What can we conclude about your weight?

Well, what are our choices?

- We could choose to only believe A, and assign 160lbs to our weight estimate.
- We could choose to only believe B, and assign 170lbs to our weight.
- We could choose a number less than either A or B.
- We could choose a number greater than either A or B.
- We could choose a number between A and B.

The first two choices are plausible, but we have no reason to favor one scale over the other. Why would we choose to believe A instead of B? We have no reason for such a belief. The third and fourth choices are irrational. The scales are admittedly not very accurate, but there is no reason at all to choose a number

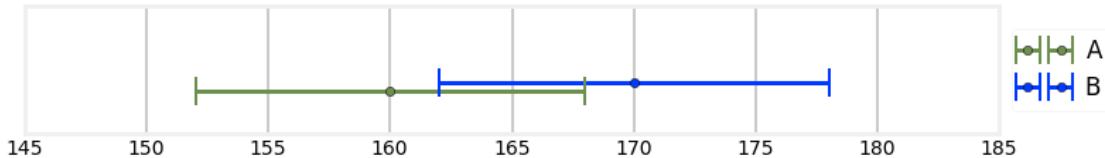
outside of the range of what they both measured. The final choice is the only reasonable one. If both scales are inaccurate, and as likely to give a result above my actual weight as below it, more often than not probably the answer is somewhere between A and B.

In mathematics this concept is formalized as expected value, and we will cover it in depth later. For now ask yourself what would be the ‘usual’ thing to happen if we made one million separate readings. Some of the times both scales will read too low, sometimes both will read too high, and the rest of the time they will straddle the actual weight. If they straddle the actual weight then certainly we should choose a number between A and B. If they don’t straddle then we don’t know if they are both too high or low, but by choosing a number between A and B we at least mitigate the effect of the worst measurement. For example, suppose our actual weight is 180 lbs. 160 lbs is a big error. But if we choose a weight between 160 lbs and 170 lbs our estimate will be better than 160 lbs. The same argument holds if both scales returned a value greater than the actual weight.

We will deal with this more formally later, but for now I hope it is clear that our best estimate is the average of A and B.  $\frac{160+170}{2} = 165$ .

We can look at this graphically. I have plotted the measurements of A and B with an assumed error of  $\pm 8$  lbs. The overlap falls between 160 and 170 so the only weight that makes sense must lie within 160 and 170 pounds. a lot of the plotting code is not particularly useful to read, so for each chapter I have placed the uninteresting code in a file named xxx\_internal. I import this file and call whatever function I need. If you want to read the code, the file exists in the code subdirectory.

```
In [2]: import gh_internal as gh
gh.plot_errorbar1()
```

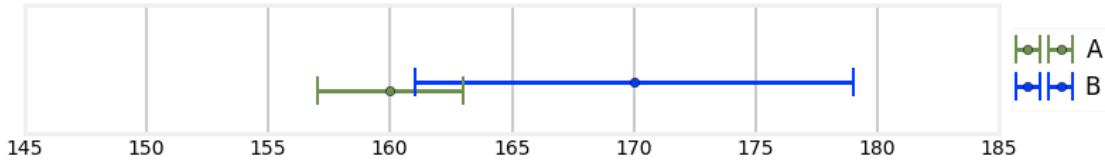


So 165 lbs looks like a reasonable estimate, but there is more information here that we might be able to take advantage of. The only weights that are possible lie in the intersection between the error bars of A and B. For example, a weight of 161 lbs is impossible because scale B could not give a reading of 170 lbs with a maximum error of 8 pounds. Likewise a weight of 171 lbs is impossible because scale A could not give a reading of 160 lbs with a maximum error of 8 lbs. In this example the only possible weights lie in the range of 162 to 168 lbs.

That doesn’t yet allow us to find a better weight estimate, but let’s play ‘what if’ some more. What if we are now told that A is three times more accurate than B? Consider the 5 options we listed above. It still makes no sense to choose a number outside the range of A and B, so we will not consider those. It perhaps seems more compelling to choose A as our estimate - after all, we know it is more accurate, why not use it instead of B? Can B possibly improve our knowledge over A alone?

The answer, perhaps counter intuitively, is yes, it can. First, let’s look at the same measurements of A=160 and B=170, but with the error of A  $\pm 3$  lbs and the error of B is 3 times as much,  $\pm 9$  lbs.

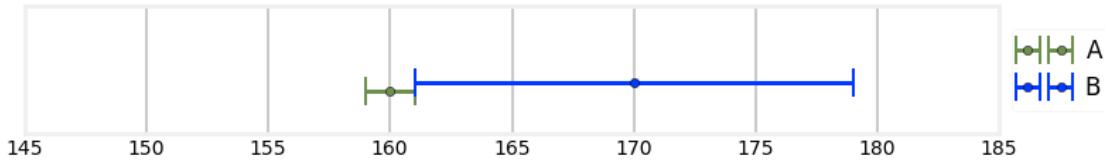
```
In [3]: gh.plot_errorbar2()
```



The overlap of the error bars of A and B are the only possible true weight. This overlap is smaller than the error in A alone. More importantly, in this case we can see that the overlap doesn't include 160 lbs or 165 lbs. If we only used the measurement from A because it is more accurate than B we would give an estimate of 160 lbs. If we average A and B we would get 165 lbs. Neither of those weights are possible given our knowledge of the accuracy of the scales. By including the measurement of B we would give an estimate somewhere between 161 lbs and 163 lbs, the limits of the intersections of the two error bars.

Let's take this to the extreme limits. Assume we know scale A is accurate to 1 lb. In other words, if we truly weigh 170 lbs, it could report 169, 170, or 171 lbs. We also know that scale B is accurate to 9 lbs. We do a weighing on each scale, and get A=160, and B=170. What should we estimate our weight to be? Let's look at that graphically.

In [4]: `gh.plot_errorbar3()`



Here we can see that the only possible weight is 161 lbs. This is an important result. With two relatively inaccurate sensors we are able to deduce an extremely accurate result.

So two sensors, even if one is less accurate than the other, is better than one.

However, we have strayed from our problem. No customer is going to want to buy multiple scales, and besides, we initially started with an assumption that all scales were equally (in)accurate. This insight of using all measurements regardless of accuracy will play a large role later, so don't forget it.

What if I have one scale, but I weigh myself many times? We concluded that if we had two scales of equal accuracy we should average the results of their measurements. What if I weigh myself 10,000 times with one scale? We have already stated that the scale is equally likely to return a number too large as it is to return one that is too small. It is not that hard to prove that the average of a large number of weights will be very close to the actual weight, but let's write a simulation for now.

```
In [5]: import numpy as np
measurements = np.random.uniform(160, 170, size=10000)
print('Average of measurements is {:.4f}'.format(
    measurements.mean()))
```

Average of measurements is 165.0549

The exact number printed depends on your random number generator, but it should be very close to 165.

This code makes one assumption that probably isn't true - that the scale is as likely to read 160 as 165 for a true weight of 165 lbs. This is almost never true. Real sensors are more likely to get readings nearer the true value, and are less and less likely to get readings the further away from the true value it gets. We will cover this in detail in the Gaussian chapter. For now, I will use without further explanation the `numpy.random.normal()` function, which will produce more values nearer 165 lbs, and fewer further away. Take it on faith for now that this will produce noisy measurements very similar to how a real scale would.

```
In [6]: measurements = np.random.normal(165, 5, size=10000)
        print('Average of measurements is {:.4f}'.format(
            measurements.mean()))
```

```
Average of measurements is 165.0200
```

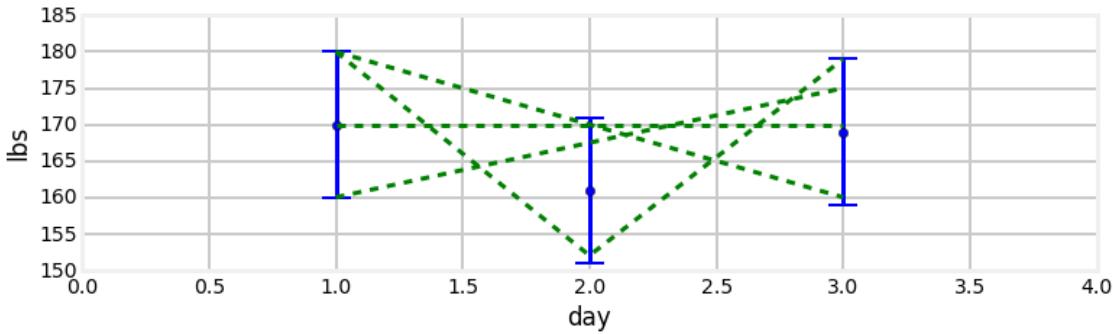
The answer again is very close to 165.

Okay, great, we have an answer to our sensor problem! But it is not a very practical answer. No one has the patience to weigh themselves ten thousand, or even a dozen times.

So, let's play 'what if' again. What if you measured your weight once a day, and got the readings 170, 161, and then 169. Did you gain weight, lose weight, or is this all just noisy measurements?

We really can't say. The first measurement was 170, and the last was 169, implying a 1 lb loss. But if the scale is only accurate to 10 lbs, that is explainable by noise. I could have actually gained weight; maybe my weight on day one was 165 lbs, and on day three it was 172. It is possible to get those weight readings with that weight gain. My scale tells me I am losing weight, and I am actually gaining weight! Let's look at that in a chart. I've plotted the weighings along with the error bars, and then some possible weight gain/losses that could be explained by those measurements in dotted green lines.

```
In [7]: gh.plot_hypothesis()
```

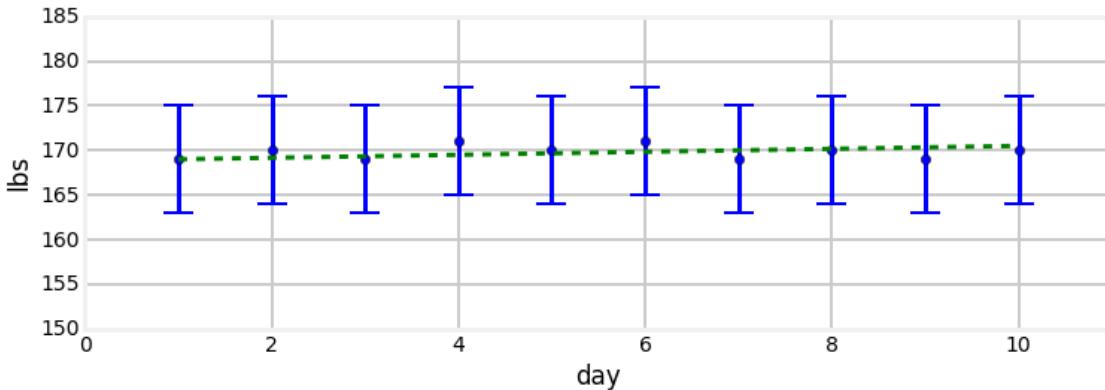


As we can see there is an extreme range of weight changes that could be explained by these three measurements. Shall we give up? No. Recall that we are talking about measuring a humans' weight. There is no way for a human to weigh 180 lbs on day 1, and 160 lbs on day 3, or to lose 30 lbs in one day only to gain it back the next (we will assume no amputations or other trauma has happened to the person). The behavior of the physical system we are measuring should influence how we interpret the measurements.

Suppose I take a different scale, and I get the following measurements: 169, 170, 169, 171, 170, 171, 169, 170, 169, 170. What does your intuition tell you? It is possible, for example, that you gained 1 lb each day, and the noisy measurements just happens to look like you stayed the same weight. Equally, you could have

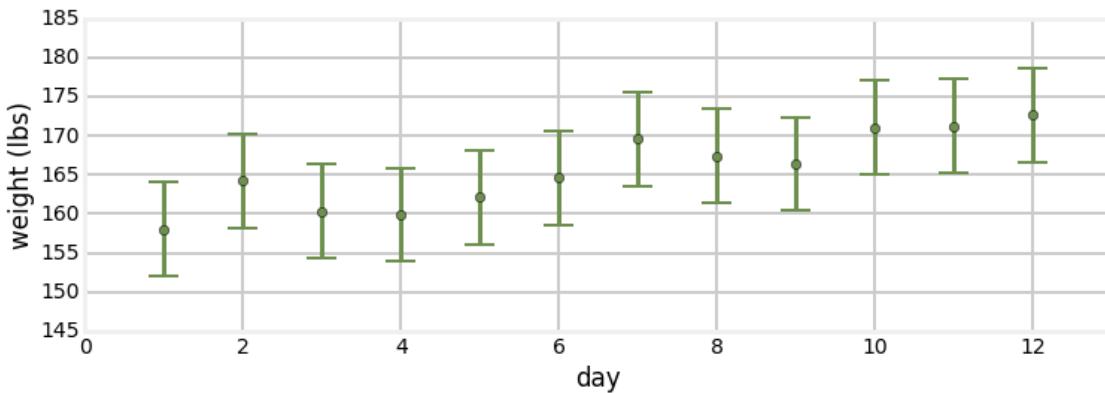
lost 1 lb a day and gotten the same readings. But is that likely? How likely is it to flip a coin and get 10 heads in a row? Not very likely. We can't prove it based solely on these readings, but it seems pretty likely that my weight held steady. In the chart below I've plotted the measurements with error bars, and a likely true weight in dashed green. This dashed line is not meant to be the 'correct' answer to this problem, merely one that is reasonable and could be explained by the measurement.

In [8]: `gh.plot_hypothesis2()`



Another what if: what if the readings were 158.0, 164.2, 160.3, 159.9, 162.1, 164.6, 169.6, 167.4, 166.4, 171.0? Let's look at a chart of that and then answer some questions.

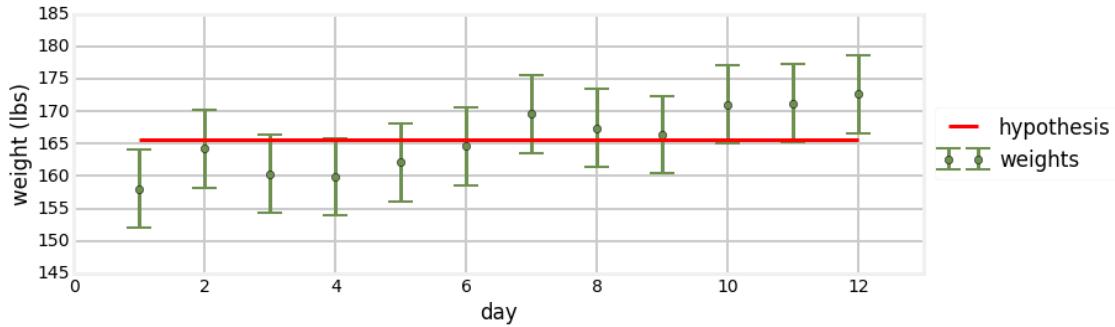
In [9]: `gh.plot_hypothesis3()`



Does it 'seem' likely that I lost weight and this is just really noisy data? Not really. Does it seem likely that I held the same weight? Again, no. This data trends upwards over time; not evenly, but definitely upwards. We can't be sure, but that surely looks like a weight gain, and a significant weight gain at that. Let's test this assumption with some more plots. It is often easier to 'eyeball' data in a chart versus a table.

So let's look at two hypotheses. First, let's assume our weight did not change. To get that number we agreed that we should average the measurements. Let's look at that.

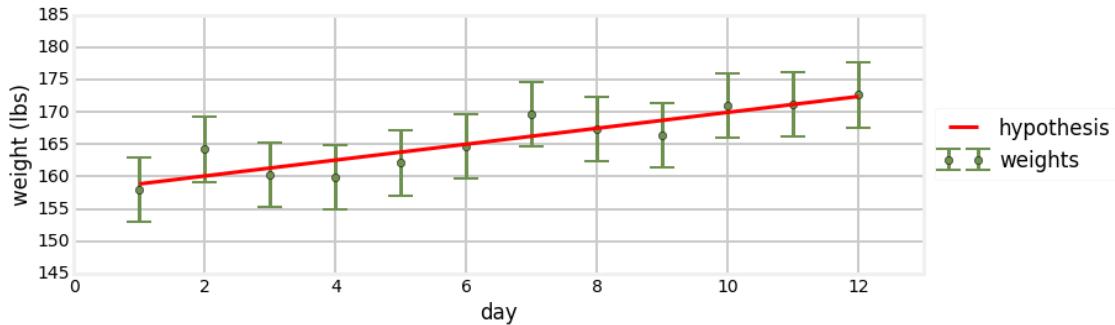
In [10]: `gh.plot_hypothesis4()`



That doesn't look very convincing. In fact, we can see that there is no horizontal line that we could draw that is inside all of the error bars.

Now, let's assume we gained weight. How much? I don't know, but numpy does! We want to draw a line through the measurements that looks 'about' right. numpy has functions that will do this according to a rule called "least squares fit". Let's not worry about the details of that computation, or why we are writing our own filter if numpy provides one, and plot the results.

In [11]: `gh.plot_hypothesis5()`



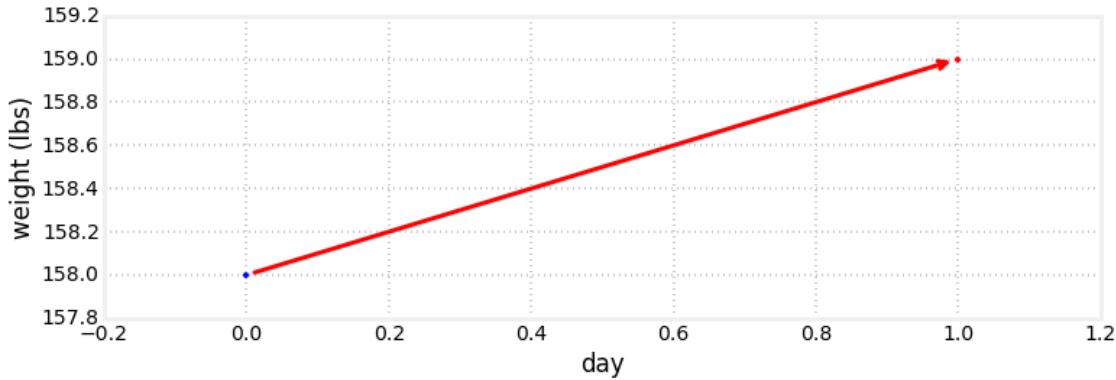
This looks much better, at least to my eyes. Notice now the hypothesis lies very close to each measurement, whereas in the previous plot the hypothesis was often quite far from the measurement. It seems far more likely to be true that I gained weight than I didn't gain any weight. Did I actually gain 13 lbs? Who can say? That seems impossible to answer.

"But is it impossible?" pipes up a coworker.

Let's try something crazy. Let's assume that I know I am gaining about one lb a day. It doesn't matter how I know that right now, assume I know it is approximately correct. Maybe I am eating a 6000 calorie a day diet, which would result in such a weight gain. Or maybe there is another way to estimate the weight gain. Let's see if we can make use of such information if it was available without worrying about the source of that information yet.

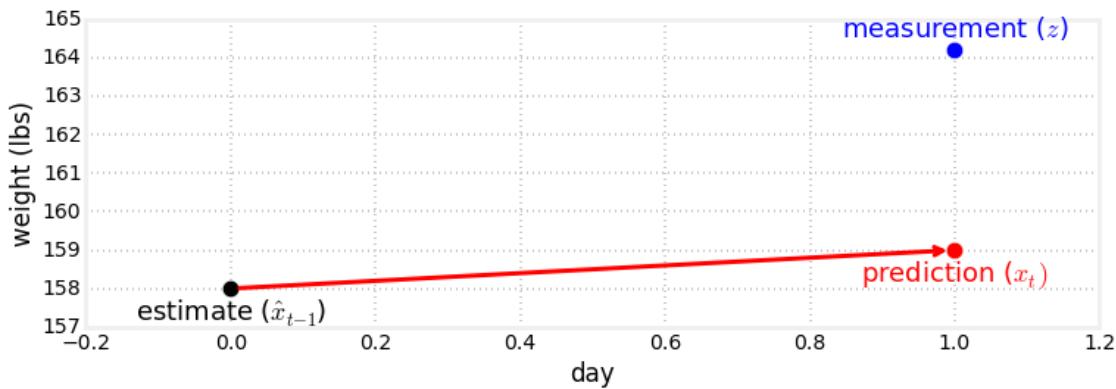
The first measurement was 158. We have no way of knowing any different, so let's accept that as our estimate. If our weight today is 158, what will it be tomorrow? Well, we think we are gaining weight at 1 lb/day, so our prediction is 159, like so:

In [12]: `gh.plot_estimate_chart_1()`



Okay, but what good is this? Sure, we could assume the 1 lb/day is accurate, and predict our weight for the next 10 days, but then why use a scale at all if we don't incorporate its readings? So let's look at the next measurement. We step on the scale again and it displays 164.2 lbs.

In [13]: `gh.plot_estimate_chart_2()`



We have a problem. Our prediction doesn't match our measurement. But, that is what we expected, right? If the prediction was always exactly the same as the measurement, it would not be capable of adding any information to the filter.

The key insight to this entire book is in the next paragraph. Read it carefully!

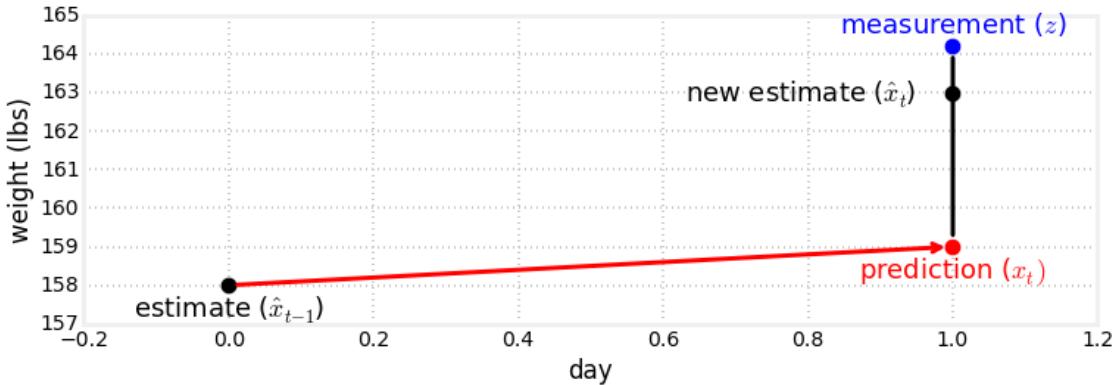
So what do we do? If we only take data from the measurement then the prediction will not affect the result. If we only take data from the prediction then the measurement will be ignored. If this is to work we need to take some kind of blend of the prediction and measurement (I've italicized the key point).

Blending two values - this sounds a lot like the two scale problem earlier. Using the same reasoning as before we can see that the only thing that makes sense is to choose a number between the prediction and the measurement. For example, an estimate of 165 makes no sense, nor does 157. Our estimates should lie between 159 (the prediction) and 164.2 (the measurement).

Should it be half way? Maybe, but in general it seems like we might know that our prediction is more or less accurate compared to the measurements. Probably the accuracy of our prediction differs from the accuracy

of the scale. Recall what we did when scale A was much more accurate than scale B - we scaled the answer to be closer to A than B. Let's look at that in a chart.

In [14]: `gh.plot_estimate_chart_3()`



Now let's try a randomly chosen number to scale our estimate:  $\frac{4}{10}$ . Our estimate will be four tenths the measurement and the rest will be from the prediction. In other words, we are expressing a belief here, a belief that the prediction is somewhat more likely to be correct than the measurement. We compute that as

$$\text{new\_estimate} = \text{prediction} + \frac{4}{10}(\text{measurement} - \text{prediction})$$

The difference between the measurement and prediction is called the residual, which is depicted by the black vertical line in the plot above. This will become an important value to use later on, as it is an exact computation of the difference between measurements and the filter's output. Smaller residuals imply better performance.

Let's code that and see the results when we test it against the series of weights from above. We have to take into account one other factor. Weight gain has units of lbs/time, so to be general we will need to add a time step  $t$ , which we will set to 1 (day).

I hand generated the weight data to correspond to a true starting weight of 160 lbs, and a weight gain of 1 lb per day. In other words on the first day (day zero) the true weight is 160lbs, on the second day (day one, the first day of weighing) the true weight is 161 lbs, and so on.

We need to make a guess for the initial weight. It is too early to talk about initialization strategies, so for now I will assume 160 lbs.

```
In [15]: import book_plots
import matplotlib.pyplot as plt

weights = [158.0, 164.2, 160.3, 159.9, 162.1, 164.6,
           169.6, 167.4, 166.4, 171.0, 171.2, 172.6]

time_step = 1 # day
scale_factor = 4/10

def predict_using_gain_guess(weight, gain_rate, do_print=True):
    # store the filtered results
```

```

estimates, predictions = [weight], []

# most filter literature uses 'z' for measurements
for z in weights:
    # predict new position
    prediction = weight + gain_rate * time_step

    # update filter
    weight = prediction + scale_factor * (z - prediction)

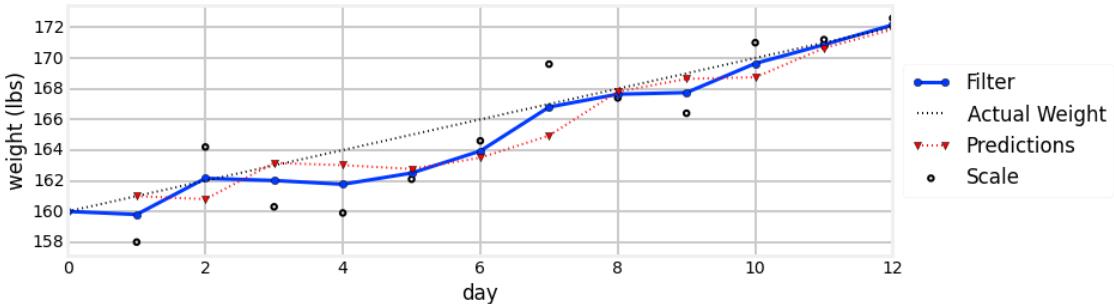
    # save
    estimates.append(weight)
    predictions.append(prediction)
    if do_print:
        gh.print_results(estimates, prediction, weight)

# plot results
gh.plot_gh_results(weights, estimates, predictions)

initial_guess = 160.
predict_using_gain_guess(weight=initial_guess, gain_rate=1)

previous: 160.00, prediction: 161.00 estimate 159.80
previous: 159.80, prediction: 160.80 estimate 162.16
previous: 162.16, prediction: 163.16 estimate 162.02
previous: 162.02, prediction: 163.02 estimate 161.77
previous: 161.77, prediction: 162.77 estimate 162.50
previous: 162.50, prediction: 163.50 estimate 163.94
previous: 163.94, prediction: 164.94 estimate 166.80
previous: 166.80, prediction: 167.80 estimate 167.64
previous: 167.64, prediction: 168.64 estimate 167.75
previous: 167.75, prediction: 168.75 estimate 169.65
previous: 169.65, prediction: 170.65 estimate 170.87
previous: 170.87, prediction: 171.87 estimate 172.16

```



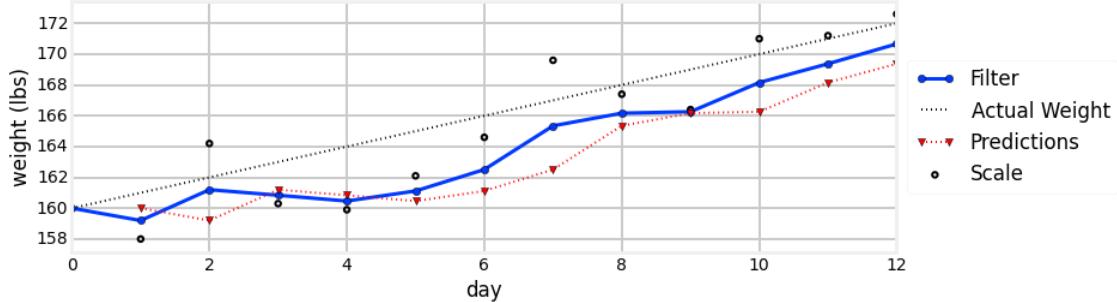
That is pretty good! There is a lot of data here, so let's talk about how to interpret it. The thick blue line shows the estimate from the filter. It starts at day 0 with the initial guess of 160 lbs. The red line shows the prediction that is made from the previous day's weight. So, on day one the previous weight was 160 lbs, the weight gain is 1 lb, and so the first prediction is 161 lbs. The estimate on day one is then part way between the prediction and measurement at 159.8 lbs. Above the chart is a print out of the previous weight,

predicted weight, and new estimate for each day. Finally, the thin black line shows the actual weight gain of the person being weighed.

The estimates are not a straight line, but they are straighter than the measurements and somewhat close to the trend line we created. Also, it seems to get better over time.

This may strike you as quite silly; of course the data will look good if we assume the conclusion, that our weight gain is around 1 lb/day! Let's see what the filter does if our initial guess is bad. Let's see what happens if I predict that there is no weight gain.

```
In [16]: predict_using_gain_guess(initial_guess, 0, do_print=False)
```



That is not so impressive. Clearly a filter that requires us to correctly guess a rate of change is not very useful. Even if our initial guess was correct, the filter will fail as soon as that rate of change changes. If I stop overeating the filter will have extremely difficulty in adjusting to that change.

But, ‘what if’? What if instead of leaving the weight gain at the initial guess of 1 lb (or whatever), we compute it from the existing measurements and estimates. On day one our estimate for the weight is:

$$(160 + 1) + \frac{4}{10}(158 - 161) = 159.8$$

On the next day we measure 164.2, which implies a weight gain of 4.4 lbs (since  $164.2 - 159.8 = 4.4$ ), not 1. Can we use this information somehow? It seems plausible. After all, the weight measurement itself is based on a real world measurement of our weight, so there is useful information. Our estimate of our weight gain may not be perfect, but it is surely better than just guessing our gain is 1 lb. Data is better than a guess, even if it is noisy.

So, should we set the new gain/day to 4.4 lbs? Yesterday we thought the weight gain was 1 lb, today we think it is 4.4 lbs. We have two numbers, and want to combine them somehow. Hmm, sounds like our same problem again. Let's use our same tool, and the only tool we have so far - pick a value part way between the two. This time I will use another arbitrarily chosen number,  $\frac{1}{3}$ . The equation is identical as for the weight estimate except we have to incorporate time because this is a rate (gain/day):

$$\text{new gain} = \text{old gain} + \frac{1}{3} \frac{\text{measurement} - \text{predicted weight}}{1 \text{ day}}$$

```
In [17]: weight = 160 # initial guess
           gain_rate = 1.0 # initial guess

           time_step = 1
           weight_scale = 4/10
```

```

gain_scale = 1/3
estimates = [weight]
predictions = []

for z in weights:
    # prediction step
    weight = weight + gain_rate*time_step
    gain_rate = gain_rate
    predictions.append(weight)

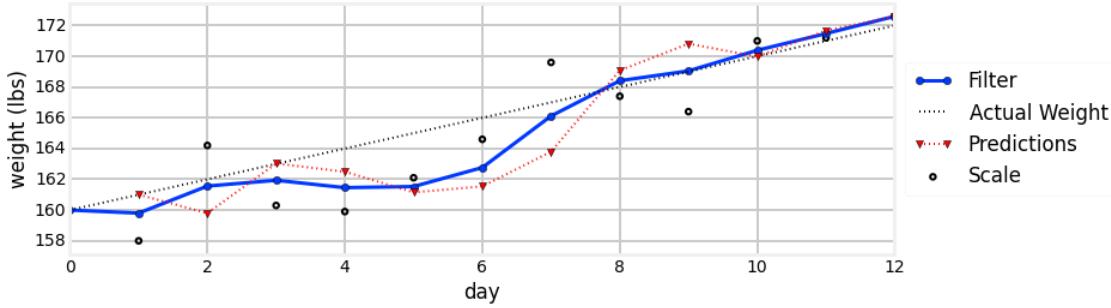
    # update step
    residual = z - weight

    gain_rate = gain_rate + gain_scale * (residual/time_step)
    weight = weight + weight_scale * residual

    estimates.append(weight)

gh.plot_gh_results(weights, estimates, predictions);

```



**note:** the semi-colon on the last line suppresses a string similar to <matplotlib.text.Text at 0x80f5160> being printed in the output. The notebook prints the value of the last line in the input cell unless you use a semi-colon as a terminator.

I think this is starting to look really good. We used no methodology for choosing our scaling factors of  $\frac{4}{10}$  and  $\frac{1}{3}$  (actually, they are poor choices for this problem), and we ‘luckily’ choose 1 lb/day as our initial guess for the weight gain, but otherwise all of the reasoning followed from very reasonable assumptions.

One final point before we go on. In the prediction step I wrote the line

```
gain_rate = gain_rate
```

This obviously has no effect, and can be removed. I wrote this to emphasize that in the prediction step you need to predict next value for all variables, both `weight` and `gain_rate`. In this case we are assuming that the gain does not vary, but when we generalize this algorithm we will remove that assumption.

## 1.2 The g-h Filter

This algorithm is known as the g-h filter.  $g$  and  $h$  refer to the two scaling factors that we used in our example.  $g$  is the scaling we used for the measurement (weight in our example), and  $h$  is the scaling for the change in

measurement over time (lbs/day in our example).

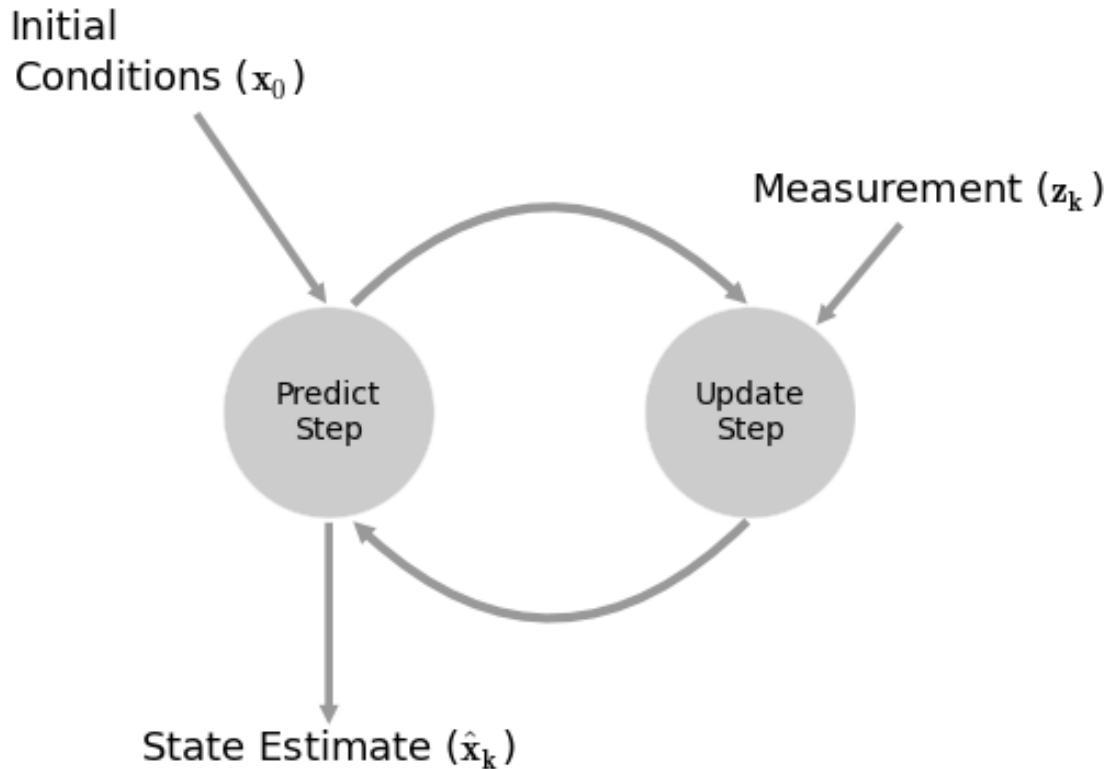
This filter is the basis for a huge number of filters, including the Kalman filter. In other words, the Kalman filter is a form of the g-h filter, which I will prove later in the book. So is the Least Squares filter, which you may have heard of, and so is the Benedict-Bordner filter, which you probably have not. Each filter has a different way of assigning values to  $g$  and  $h$ , but otherwise the algorithms are identical. For example, the  $\alpha\beta$  filter assigns a constant to  $g$  and  $h$ , constrained to a certain range of values. Other filters such as the Kalman will vary  $g$  and  $h$  dynamically at each time step.

**Let me repeat the key points as they are so important.** If you do not understand these you will not understand the rest of the book. If you do understand them, then the rest of the book will unfold naturally for you as mathematical elaborations to various ‘what if’ questions we will ask about  $g$  and  $h$ . The math may look profoundly different, but the algorithm will be exactly the same.

- Multiple data points are more accurate than one data point, so throw nothing away no matter how inaccurate it is.
- Always choose a number part way between two data points to create a more accurate estimate.
- Predict the next measurement and rate of change based on the current estimate and how much we think it will change.
- The new estimate is then chosen as part way between the prediction and next measurement.

Let's look at a visual depiction of the algorithm.

In [18]: `gh.create_predict_update_chart()`



Let me introduce some more formal terminology. The system is the object that we want to estimate. In this chapter the system is whatever we are trying to weigh. Some texts call this the plant. That terminology comes from control system theory. [https://en.wikipedia.org/wiki/Plant\\_\(control\\_theory\)](https://en.wikipedia.org/wiki/Plant_(control_theory))

The state of the system is the current configuration or values of that system that is of interest to us. We are interested only in the weight reading. If I put a 100 kg weight on the scale, the state is 100kg. We define the state based on what is relevant to us. The color of the scale is irrelevant to us so we do not include those values in the state. A QA engineer for the manufacturer might include color in the state so that she can track and control the manufacturing process. The state estimate is our filter's estimate of the state. For example, for the 100 kg weight our estimate might be 99.327 kg due to sensor errors.

We use a process model to mathematically model the system. In this chapter our process model is the assumption that my weight today is yesterday's weight plus my weight gain for the last day. The process model does not model or otherwise account for the sensors. Another example would be a process model for an automobile. The process model might be "distance equals velocity times time. This model is not perfect as the velocity of a car can vary over a non-zero amount of time, the tires can slip on the road, and so on. The system error or process error is the error in this model. We never know this value exactly; if we did we could refine our model to have zero error. Some texts use plant model and plant error. You may also see system model. They all mean the same thing.

The predict step is known as system propagation. It uses the process model to form a new state estimate. Because of the process error this prediction is normally imperfect. Assuming we are tracking data over time, we say we propagate the state into the future. Some texts call this the evolution.

The update step is known as the measurement update. One iteration of the system propagation and measurement update is known as an epoch.

Now let's explore a few different problem domains to better understand this algorithm. Consider the problem of trying to track a train on a track. The track constrains the position of the train to a very specific region. Furthermore, trains are large and slow. It takes them many minutes to slow down or speed up significantly. So, if I know that the train is at kilometer marker 23 km at time t and moving at 18 kph, I can be extremely confident in predicting its position at time t + 1 second. And why is that important? Suppose we can only measure its position with an accuracy of  $\pm 250$  meters. The train is moving at 18 kph, which is 5 meters per second. So at t+1 second the train will be at 23.005 km yet the measurement could be anywhere from 22.755 km to 23.255 km. So if the next measurement says the position is at 23.4 we know that must be wrong. Even if at time t the engineer slammed on the brakes the train will still be very close to 23.005 km because a train cannot slow down very much in 1 second. If we were to design a filter for this problem (and we will a bit further in the chapter!) we would want to design a filter that gave a very high weighting to the prediction vs the measurement.

Now consider the problem of tracking a thrown ball. We know that a ballistic object moves in a parabola in a vacuum when in a gravitational field. But a ball thrown on the surface of the Earth is influenced by air drag, so it does not travel in a perfect parabola. Baseball pitchers take advantage of this fact when they throw curve balls. Let's say that we are tracking the ball inside a stadium using computer vision. The accuracy of the computer vision tracking might be modest, but predicting the ball's future positions by assuming that it is moving on a parabola is not extremely accurate either. In this case we'd probably design a filter that gave roughly equal weight to the measurement and the prediction.

Now consider trying to track a child's balloon in a hurricane. We have no legitimate model that would allow us to predict the balloon's behavior except over very brief time scales (we know the balloon cannot go 10 miles in 1 second, for example). In this case we would design a filter that emphasized the measurements over the predictions.

Most of this book is devoted to expressing the concerns in the last three paragraphs mathematically, which then allows us to find an optimal solution. In this chapter we will merely be assigning different values to  $g$  and  $h$  in a more intuitive, and thus less optimal way. But the fundamental idea is to blend somewhat inaccurate measurements with somewhat inaccurate models of how the systems behaves to get a filtered estimate that is better than either information source by itself.

We can express this as an algorithm:

### Initialization

1. Initialize the state of the filter
2. Initialize our belief in the state

### Predict

1. Use system behavior to predict state at the next time step
2. Adjust belief to account for the uncertainty in prediction

### Update

1. Get a measurement and associated belief about its accuracy
2. Compute residual between estimated state and measurement
3. New estimate is somewhere on the residual line

We will use this same algorithm throughout the book, albeit with some modifications.

## 1.3 Notation

I'll begin to introduce the notations and variable names used in the literature. Some of this was already used in the above charts. Measurement is typically denoted  $z$  and that is what we will use in this book (some literature uses  $y$ ). Subscript  $k$  indicates the time step, so  $z_k$  is the data for this time step. A bold font denotes a vector or matrix. So far we have only considered having one sensor, and hence one sensor measurement, but in general we may have  $n_s$  sensors and  $n_m$  measurements.  $\mathbf{x}$  denotes our state, and is bold to denote that it is a vector. For example, for our scale example, it represents both the initial weight and initial weight gain rate, like so:

$$\mathbf{x} = \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

For a weight of 62 kg with a gain of 0.3 kg/day we have

$$\mathbf{x} = \begin{bmatrix} 62 \\ 0.3 \end{bmatrix}$$

So, the algorithm is simple. The state is initialized with  $\mathbf{x}_0$ , the initial estimate. We then enter a loop, predicting the state for time or step  $k$  from the values from time/step  $k - 1$ . We then get the measurement  $z_k$  and choose some intermediate point between the measurements and prediction, creating the estimate  $\mathbf{x}_k$ .

## 1.4 NumPy arrays

We will find NumPy's array data structure to be indispensable through the rest of the book, so let's take the time to learn about them now. I will teach you enough to get started; refer to NumPy's documentation if you want to become an expert with them.

`numpy.array` implements a one or more dimensional array. Its type is `numpy.ndarray`, and we will refer to this as an ndarray for short. You can construct it with any list like object. The following constructs a 1-D array from a list:

```
In [5]: import numpy as np
x = np.array([1, 2, 3])
print(type(x))
x

<class 'numpy.ndarray'>
```

```
Out[5]: array([1, 2, 3])
```

You can also use tuples:

```
In [4]: x = np.array((4,5,6))
x
```

```
Out[4]: array([4, 5, 6])
```

```
In [20]: x = np.array([[1, 2, 3],
                     [4, 5, 6]])
print(x)
```

```
[[1 2 3]
 [4 5 6]]
```

You can create arrays of 3 or more dimensions, but we have no need for that here, and so I will not elaborate.

By default the arrays use the data type of the values in the list; if there are multiple types then it will choose the type that most accurately represents all the values. So, for example, if your list contains a mix of `int` and `float` the data type of the array would be of type `float`. You can override this with the `dtype` parameter.

```
In [21]: x = np.array([1, 2, 3], dtype=float)
print(x)
```

```
[ 1.  2.  3.]
```

You can access the array elements using subscript location:

```
In [7]: x = np.array([[1, 2, 3],
                     [4, 5, 6]])

print(x[1,2])
```

```
6
```

You can access a column or row by using slices. A ‘`:`’ used as a subscript is shorthand for all data in that row or column. So `x[:,0]` returns an array of all data in the first column (the 0 specifies the first column):

```
In [8]: x[:, 0]
```

```
Out[8]: array([1, 4])
```

We can get the second row with:

In [9]: `x[1, :]`

Out[9]: `array([4, 5, 6])`

You can perform matrix addition with the `+` operator, but matrix multiplication requires the `dot` method or function. The `*` operator performs element-wise multiplication, which is **not** what you want for linear algebra.

```
In [22]: x = np.array([[1., 2.,
                     [3., 4.]])
    print('addition:\n', x+x)
    print('\nelement-wise multiplication\n', x*x)
    print('\nmultiplication\n', np.dot(x,x))
    print('\ndot is also a member of np.array\n', x.dot(x))

addition:
[[ 2.  4.]
 [ 6.  8.]]

element-wise multiplication
[[ 1.  4.]
 [ 9. 16.]]

multiplication
[[ 7. 10.]
 [15. 22.]]

dot is also a member of np.array
[[ 7. 10.]
 [15. 22.]]
```

You can get the transpose with `.T`, and the inverse with `numpy.linalg.inv`.

```
In [23]: print('transpose\n', x.T)
        print('\ninverse\n', np.linalg.inv(x))

transpose
[[ 1.  3.]
 [ 2.  4.]]

inverse
[[-2.   1. ]
 [ 1.5 -0.5]]
```

Finally, there are helper functions like `zeros` to create a matrix of all zeros, `ones` to get all ones, and `eye` to get the identity matrix. If you want a multidimensional array, use a tuple to specify the shape.

```
In [24]: print('zeros\n', np.zeros(7))
        print('\nzeros(3x2)\n', np.zeros((3, 2)))
        print('\neye\n', np.eye(3))

zeros
[ 0.  0.  0.  0.  0.  0.  0.]
```

```

zeros(3x2)
[[ 0.  0.]
 [ 0.  0.]
 [ 0.  0.]]

eye
[[ 1.  0.  0.]
 [ 0.  1.  0.]
 [ 0.  0.  1.]]

```

NumPy is part of the SciPy package. Both are well documented. You can find a tutorial and links to the full documentation at the following location:

<https://docs.scipy.org/doc/numpy-dev/user/quickstart.html>

### 1.4.1 Exercise - Create arrays

I want you to create a NumPy array of 10 elements with each element containing 1/10. There are several ways to do this; try to implement as many as you can think of.

In [25]: *# your solution*

### 1.4.2 Solution

Here are three ways to do this. The first one is the one I want you to know. I used the ‘/’ operator to divide all of the elements of the array with 10. We will shortly use this to convert the units of an array from meters to km.

```

In [26]: print(np.ones(10) / 10.)
          print(np.array([.1, .1, .1, .1, .1, .1, .1, .1, .1]))
          print(np.array([.1]*10))

[ 0.1  0.1  0.1  0.1  0.1  0.1  0.1  0.1  0.1]
[ 0.1  0.1  0.1  0.1  0.1  0.1  0.1  0.1  0.1]
[ 0.1  0.1  0.1  0.1  0.1  0.1  0.1  0.1  0.1]

```

Here is one I haven’t covered yet. The function `numpy.asarray()` will convert its argument to an ndarray if it isn’t already one. If it is, the data is unchanged. This is a handy way to write a function that can accept either Python lists or ndarrays, and it is very efficient if the type is already ndarray as nothing new is created.

```

In [27]: def one_tenth(x):
          x = np.asarray(x)
          return x / 10

          print(one_tenth([1, 2, 3]))           # I work!
          print(one_tenth(np.array([4, 5, 6]))) # so do I!

[ 0.1  0.2  0.3]
[ 0.4  0.5  0.6]

```

## 1.5 Exercise: Write Generic Algorithm

In the example above, I explicitly coded this to solve the weighing problem that we've been discussing throughout the chapter. For example, the variables are named "weight\_scale", "gain", and so on. I did this to make the algorithm easy to follow - you can easily see that we correctly implemented each step. But, that is code written for exactly one problem, and the algorithm is the same for any problem. So let's rewrite the code to be generic - to work with any problem. Use this function signature:

```
def g_h_filter(data, x0, dx, g, h, dt):
    """
    Performs g-h filter on 1 state variable with a fixed g and h.

    'data' contains the data to be filtered.
    'x0' is the initial value for our state variable
    'dx' is the initial change rate for our state variable
    'g' is the g-h's g scale factor
    'h' is the g-h's h scale factor
    'dt' is the length of the time step
    """

```

Return the data as a NumPy array, not a list. Test it by passing in the same weight data as before, plot the results, and visually determine that it works.

### 1.5.1 Solution and Discussion

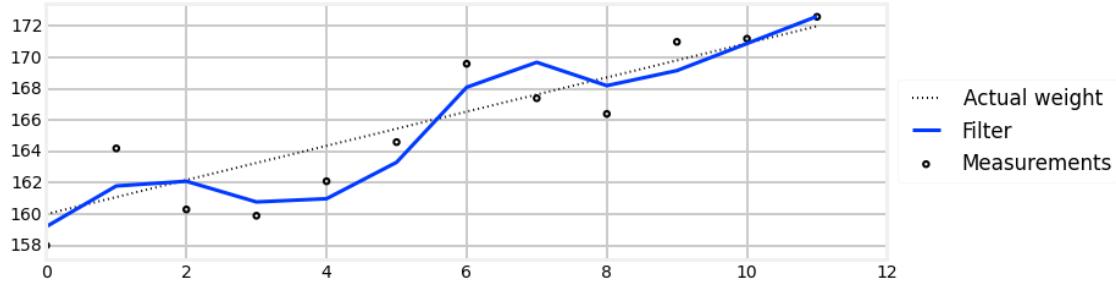
```
In [28]: from gh_internal import plot_g_h_results

def g_h_filter(data, x0, dx, g, h, dt=1.):
    x = x0
    results = []
    for z in data:
        #prediction step
        x_est = x + (dx*dt)
        dx = dx

        # update step
        residual = z - x_est
        dx = dx + h * (residual) / dt
        x = x_est + g * residual
        results.append(x)

    return np.array(results)

book_plots.plot_track([0, 11], [160, 172], label='Actual weight')
data = g_h_filter(data=weights, x0=160, dx=1, g=6./10, h=2./3, dt=1.)
plot_g_h_results(weights, data);
```



## 1.6 Choice of g and h

The g-h filter is not one filter - it is a classification for a family of filters. Eli Brookner in [Tracking and Kalman Filtering Made Easy](#) lists 11, and I am sure there are more. Not only that, but each type of filter has numerous subtypes. Each filter is differentiated by how  $g$  and  $h$  are chosen. So there is no ‘one size fits all’ advice that I can give here. Some filters set  $g$  and  $h$  as constants, others vary them dynamically. The Kalman filter varies them dynamically at each step. Some filters allow  $g$  and  $h$  to take any value within a range, others constrain one to be dependent on the other by some function  $f()$ , where  $g = f(h)$ .

The topic of this book is not the entire family of g-h filters; more importantly, we are interested in the [Bayesian](#) aspect of these filters, which I have not addressed yet. Therefore I will not cover selection of  $g$  and  $h$  in depth. [Tracking and Kalman Filtering Made Easy](#) is an excellent resource for that topic. If this strikes you as an odd position for me to take, recognize that the typical formulation of the Kalman filter does not use  $g$  and  $h$  at all. The Kalman filter is a g-h filter because it mathematically reduces to this algorithm. When we design the Kalman filter we use design criteria that can be mathematically reduced to  $g$  and  $h$ , but the Kalman filter form is usually a much more powerful way to think about the problem. Don’t worry if this is not too clear right now, it will be much clearer later after we develop the Kalman filter theory.

It is worth seeing how varying  $g$  and  $h$  affects the results, so we will work through some examples. This will give us strong insight into the fundamental strengths and limitations of this type of filter, and help us understand the behavior of the rather more sophisticated Kalman filter.

## 1.7 Exercise: create measurement function

Now let’s write a function that generates noisy data for us. In this book I model a noisy signal as the signal plus [white noise](#). We’ve not yet covered the statistics to fully understand the definition of white noise. In essence, think of it as data that randomly varies higher and lower than the signal with no pattern. We say that it is a serially uncorrelated random variable with zero mean and finite variance. If you don’t follow that, you will by the end of the [Gaussians](#) chapter. You may not be successful at this exercise if you have no knowledge of statistics. If so, just read the solution and discussion.

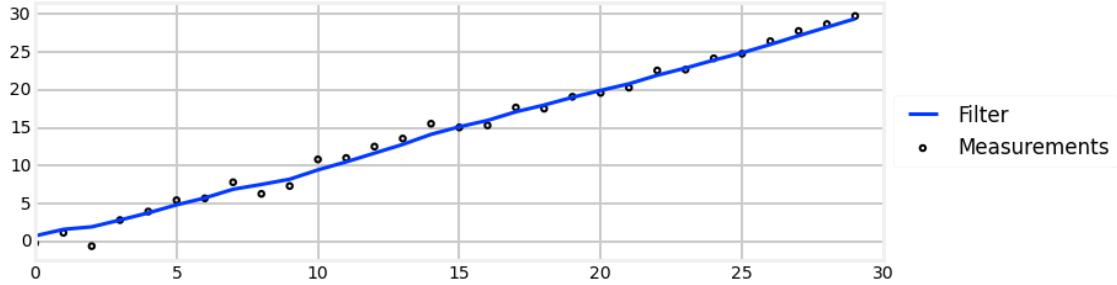
White noise can be generated by `numpy.random.randn()`. We want a function that we call with the starting value, the amount of change per step, the number of steps, and the amount of noise we want to add. It should return a list of the data. Test it by creating 30 points, filtering it with `g_h_filter()`, and plot the results with `plot_g_h_results()`.

In [29]: # your code here

### 1.7.1 Solution

```
In [30]: from numpy.random import randn
def gen_data(x0, dx, count, noise_factor):
    return [x0 + dx*i + randn()*noise_factor for i in range(count)]

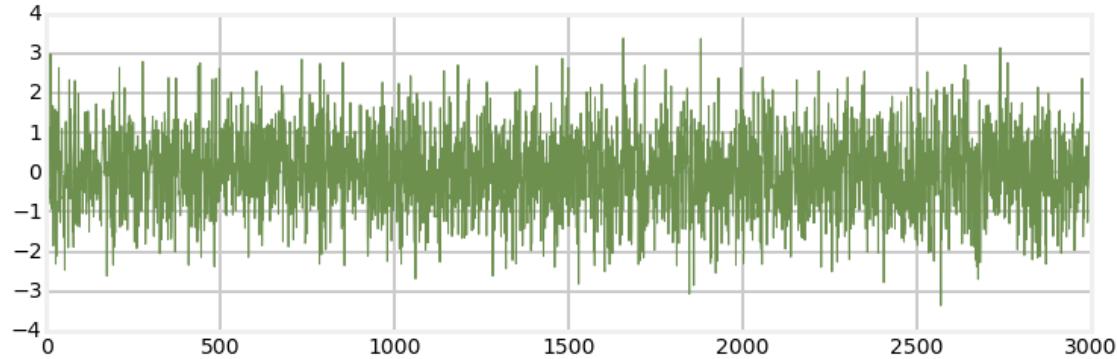
measurements = gen_data(0, 1, 30, 1)
data = g_h_filter(data=measurements, x0=0, dx=1, dt=1, g=.2, h=0.02)
plot_g_h_results(measurements, data)
```



### 1.7.2 Discussion

`randn()` returns random numbers centered around 0 - it is just as likely to be greater than zero as under zero. It varies by one standard deviation - don't worry if you don't know what that means. I've plotted 3000 calls to `randn()` - you can see that the values are centered around zero and mostly range from a bit under -1 to a bit more than +1, though occasionally they get much larger.

```
In [31]: plt.plot([randn() for _ in range(3000)], lw=1);
```



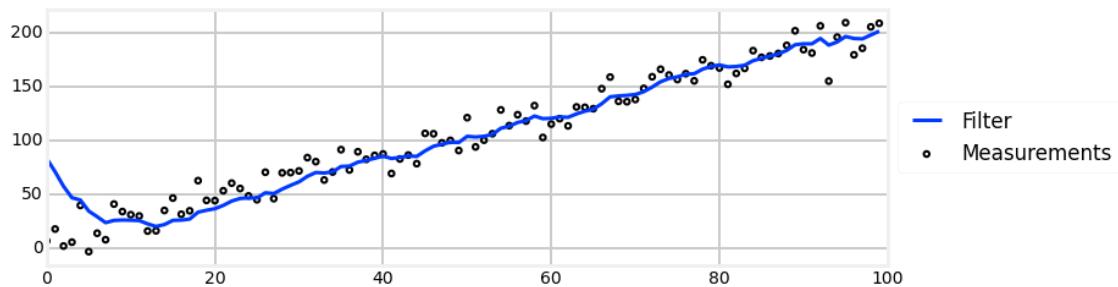
## 1.8 Exercise: Bad Initial Conditions

Now write code that uses `gen_data` and `g_h_filter` to filter 100 data points that starts at 5, has a derivative of 2, a noise scaling factor of 10, and uses  $g=0.2$  and  $h=0.02$ . Set your initial guess for  $x$  to be 100.

In [32]: # your code here

### 1.8.1 Solution and Discussion

```
In [33]: zs = gen_data(x0=5, dx=2, count=100, noise_factor=10)
         data = g_h_filter(data=zs, x0=100., dx=2., dt=1., g=0.2, h=0.01)
         plot_g_h_results(measurements=zs, filtered_data=data)
```



The filter starts out with estimates that are far from the measured data due to the bad initial guess of 100. You can see that it ‘rings’ before settling in on the measured data. ‘Ringing’ means that the signal overshoots and undershoots the data in a sinusoidal type pattern. This is a very common phenomena in filters, and a lot of work in filter design is devoted to minimizing ringing. That is a topic that we are not yet prepared to address, but I wanted to show you the phenomenon.

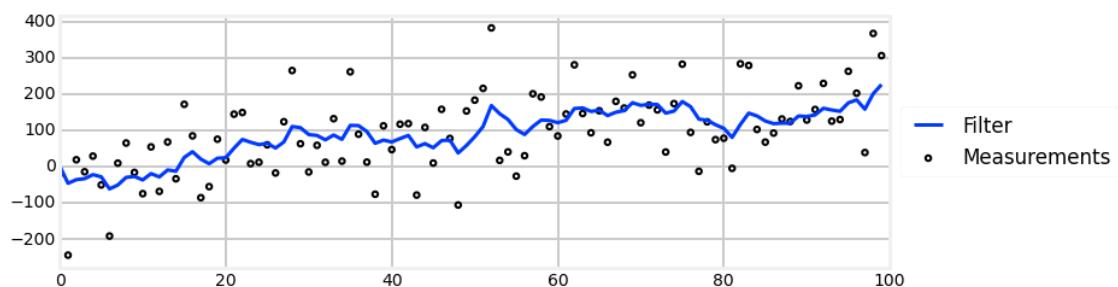
## 1.9 Exercise: Extreme Noise

Rerun the same test, but this time use a noise factor of 100. Remove the initial condition ringing by changing the initial condition from 100 down to 5.

In [34]: # your code here

### 1.9.1 Solution and Discussion

```
In [35]: zs = gen_data(x0=5, dx=2, count=100, noise_factor=100)
         data = g_h_filter(data=zs, x0=5., dx=2., g=0.2, h=0.02)
         plot_g_h_results(measurements=zs, filtered_data=data)
```



This doesn't look so wonderful to me. We can see that perhaps the filtered signal varies less than the noisy signal, but it is far from the straight line. If we were to plot just the filtered result no one would guess that the signal with no noise starts at 5 and increments by 2 at each time step. And while in locations the filter does seem to reduce the noise, in other places it seems to overshoot and undershoot.

At this point we don't know enough to really judge this. We added **a lot** of noise; maybe this is as good as filtering can get. However, the existence of the multitude of chapters beyond this one should suggest that we can do much better than this suggests.

## 1.10 Exercise: The Effect of Acceleration

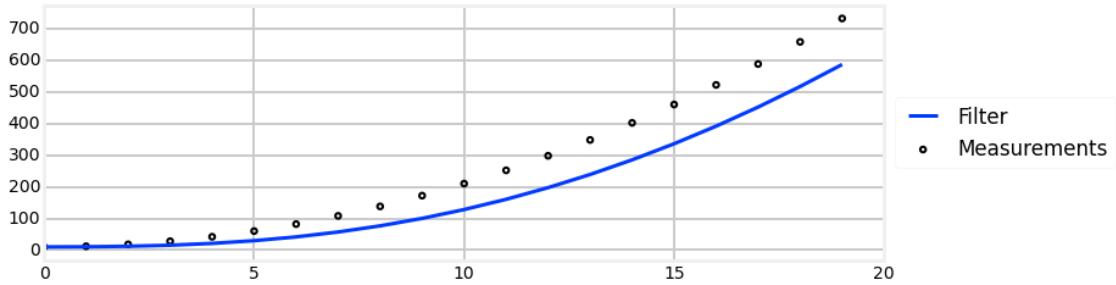
Write a new data generation function that adds in a constant acceleration factor to each data point. In other words, increment  $dx$  as you compute each data point so that the velocity ( $dx$ ) is ever increasing. Set the noise to 0,  $g = 0.2$  and  $h = 0.02$  and plot the results. Explain what you see.

In [36]: # your code here

### 1.10.1 Solution and Discussion

```
In [37]: def gen_data(x0, dx, count, noise_factor, accel=0):
    zs = []
    for i in range(count):
        zs.append(x0 + dx*i + randn()*noise_factor)
        dx += accel
    return zs

predictions = []
zs = gen_data(x0=10, dx=0, count=20, noise_factor=0, accel=2)
data = g_h_filter(data=zs, x0=10, dx=0, g=0.2, h=0.02)
plt.xlim([0, 20])
plot_g_h_results(measurements=zs, filtered_data=data)
```



Each prediction lags behind the signal. If you think about what is happening this makes sense. Our model assumes that velocity is constant. The g-h filter computes the first derivative of  $x$  (we use  $\dot{x}$  to denote the derivative) but not the second derivative  $\ddot{x}$ . So we are assuming that  $\ddot{x} = 0$ . At each prediction step we predict the new value of  $x$  as  $x + \dot{x} * t$ . But because of the acceleration the prediction must necessarily fall

behind the actual value. We then try to compute a new value for  $\dot{x}$ , but because of the  $h$  factor we only partially adjust  $\dot{x}$  to the new velocity. On the next iteration we will again fall short.

Note that there is no adjustment to  $g$  or  $h$  that we can make to correct this problem. This is called the lag error or systemic error of the system. It is a fundamental property of g-h filters. Perhaps your mind is already suggesting solutions or workarounds to this problem. As you might expect, a lot of research has been devoted to this problem, and we will be presenting various solutions to this problem in this book. > The ‘take home’ point is that the filter is only as good as the mathematical model used to express the system.

## 1.11 Exercise: Varying g

Now let’s look at the effect of varying  $g$ . Before you perform this exercise, recall that  $g$  is the scale factor for choosing between the measurement and prediction. What do you think of a large value of  $g$  will be? A small value?

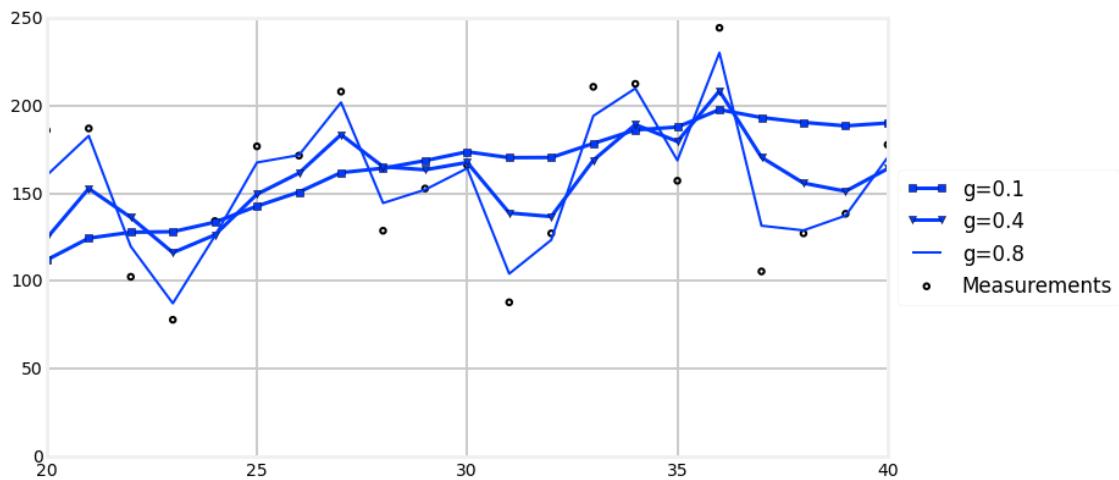
Now, let the `noise_factor=50` and `dx=5`. Plot the results of  $g = 0.1$ ,  $0.4$ , and  $0.8$ .

In [38]: # your code here

### 1.11.1 Solution and Discussion

```
In [39]: import book_format
np.random.seed(100)
zs = gen_data(x0=5, dx=5, count=50, noise_factor=50)
data1 = g_h_filter(data=zs, x0=0., dx=5., dt=1., g=0.1, h=0.01)
data2 = g_h_filter(data=zs, x0=0., dx=5., dt=1., g=0.4, h=0.01)
data3 = g_h_filter(data=zs, x0=0., dx=5., dt=1., g=0.8, h=0.01)

with book_format.figsize(y=6):
    book_plots.plot_measurements(zs, lw=1, color='k')
    book_plots.plot_filter(data1, label='g=0.1', marker='s')
    book_plots.plot_filter(data2, label='g=0.4', marker='v')
    book_plots.plot_filter(data3, label='g=0.8', lw=2)
    book_plots.show_legend()
    book_plots.set_limits([20,40], [0, 250])
```

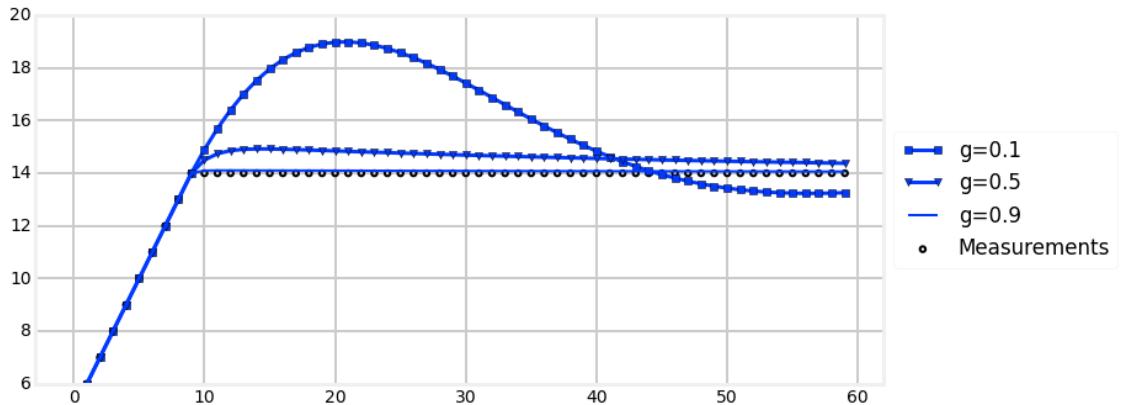


It is clear that as  $g$  is larger we more closely follow the measurement instead of the prediction. When  $g = 0.8$  we follow the signal almost exactly, and reject almost none of the noise. One might naively conclude that  $g$  should always be very small to maximize noise rejection. However, that means that we are mostly ignoring the measurements in favor of our prediction. What happens when the signal changes not due to noise, but an actual state change? Let's have a look. I will create data that has  $\dot{x} = 1$  for 9 steps before changing to  $\dot{x} = 0$ .

```
In [40]: zs = [5, 6, 7, 8, 9, 10, 11, 12, 13, 14]
for i in range(50):
    zs.append(14)

data1 = g_h_filter(data=zs, x0=4., dx=1., dt=1., g=0.1, h=0.01)
data2 = g_h_filter(data=zs, x0=4., dx=1., dt=1., g=0.5, h=0.01)
data3 = g_h_filter(data=zs, x0=4., dx=1., dt=1., g=0.9, h=0.01)

with book_format.figsize(y=5):
    book_plots.plot_measurements(zs)
    book_plots.plot_filter(data1, label='g=0.1', marker='s')
    book_plots.plot_filter(data2, label='g=0.5', marker='v')
    book_plots.plot_filter(data3, label='g=0.9', lw=2)
    book_plots.show_legend()
    plt.ylim([6, 20])
```



Here we can see the effects of ignoring the signal. We not only filter out noise, but legitimate changes in the signal as well.

Maybe we need a ‘Goldilocks’ filter, where  $g$  is not too large, not too small, but just right? Well, not exactly. As alluded to earlier, different filters choose  $g$  and  $h$  in different ways depending on the mathematical properties of the problem. For example, the Benedict-Bordner filter was invented to minimize the transient error in this example, where  $\dot{x}$  makes a step jump. We will not discuss this filter in this book, but here are two plots chosen with different allowable pairs of  $g$  and  $h$ . This filter design minimizes transient errors for step jumps in  $\dot{x}$  at the cost of not being optimal for other types of changes in  $\dot{x}$ .

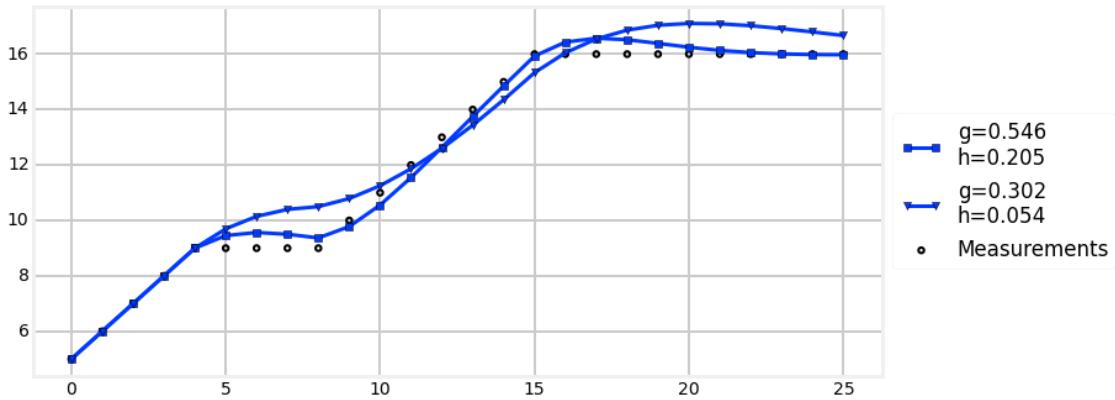
```
In [41]: zs = [5,6,7,8,9,9,9,9,10,11,12,13,14,
           15,16,16,16,16,16,16,16,16,16,16]
```

```

data1 = g_h_filter(data=zs, x0=4., dx=1., dt=1., g=.302, h=.054)
data2 = g_h_filter(data=zs, x0=4., dx=1., dt=1., g=.546, h=.205)

with book_format.figsize(y=5):
    book_plots.plot_measurements(zs)
    book_plots.plot_filter(data2, label='g=0.546\nh=0.205', marker='s')
    book_plots.plot_filter(data1, label='g=0.302\nh=0.054', marker='v')
    book_plots.show_legend()

```



## 1.12 Varying $h$

Now let's leave  $g$  unchanged and investigate the effect of modifying  $h$ . We know that  $h$  affects how much we favor the measurement of  $\dot{x}$  vs our prediction. But what does this mean? If our signal is changing a lot (quickly relative to the time step of our filter), then a large  $h$  will cause us to react to those transient changes rapidly. A smaller  $h$  will cause us to react more slowly.

We will look at three examples. We have a noiseless measurement that slowly goes from 0 to 1 in 50 steps. Our first filter uses a nearly correct initial value for  $\dot{x}$  and a small  $h$ . You can see from the output that the filter output is very close to the signal. The second filter uses the very incorrect guess of  $\dot{x} = 2$ . Here we see the filter 'ringing' until it settles down and finds the signal. The third filter uses the same conditions but it now sets  $h = 0.5$ . If you look at the amplitude of the ringing you can see that it is much smaller than in the second chart, but the frequency is greater. It also settles down a bit quicker than the second filter, though not by much.

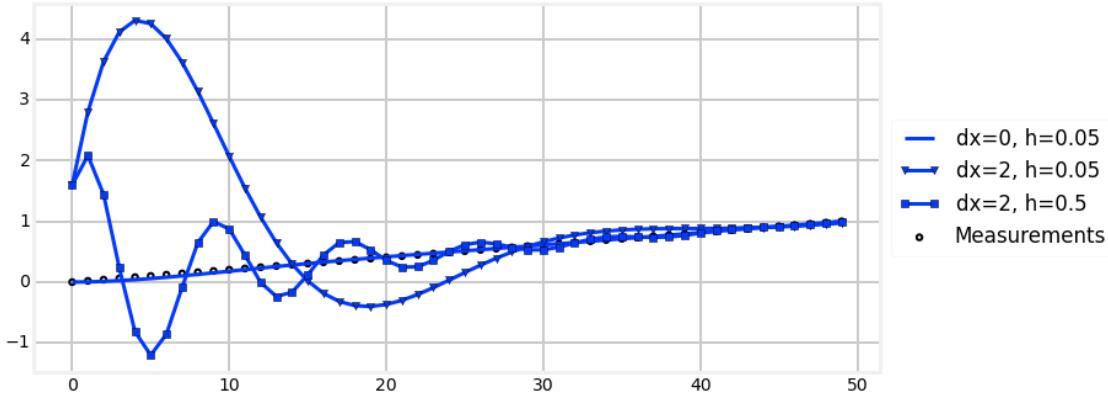
In [42]: `zs = np.linspace(0, 1, 50)`

```

data1 = g_h_filter(data=zs, x0=0, dx=0., dt=1., g=.2, h=0.05)
data2 = g_h_filter(data=zs, x0=0, dx=2., dt=1., g=.2, h=0.05)
data3 = g_h_filter(data=zs, x0=0, dx=2., dt=1., g=.2, h=0.5)

with book_format.figsize(y=5):
    book_plots.plot_measurements(zs)
    book_plots.plot_filter(data1, label='dx=0, h=0.05')
    book_plots.plot_filter(data2, label='dx=2, h=0.05', marker='v')
    book_plots.plot_filter(data3, label='dx=2, h=0.5', marker='s')
    book_plots.show_legend()

```



## 1.13 Interactive Example

For those of you running this in Jupyter Notebook I've written an interactive version of the filter so you can see the effect of changing  $\dot{x}$ ,  $g$  and  $h$  in real time. As you adjust the sliders for  $\dot{x}$ ,  $g$  and  $h$  the date will be refiltered and the results plotted for you.

If you really want to test yourself, read the next paragraph and try to predict the results before you move the sliders.

Some things to try include setting  $g$  and  $h$  to their minimum values. See how perfectly the filter tracks the data! This is only because we are perfectly predicting the weight gain. Adjust  $\dot{x}$  to larger or smaller than 5. The filter should diverge from the data and never reacquire it. Start adding back either  $g$  or  $h$  and see how the filter snaps back to the data. See what the difference in the line is when you add only  $g$  vs only  $h$ . Can you explain the reason for the difference? Then try setting  $g$  greater than 1. Can you explain the results? Put  $g$  back to a reasonable value (such as 0.1), and then make  $h$  very large. Can you explain these results? Finally, set both  $g$  and  $h$  to their largest values.

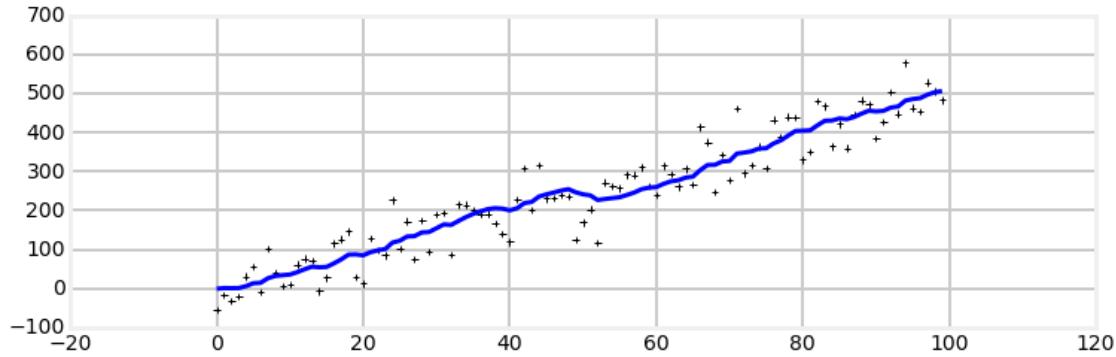
If you want to explore with this more, change the value of the array `zs` to the values used in any of the charts above and rerun the cell to see the result.

```
In [43]: from IPython.html.widgets import interact, FloatSlider

zs = gen_data(x0=5, dx=5, count=100, noise_factor=50)

def interactive_gh(x, dx, g, h):
    data = g_h_filter(data=zs, x0=x, dx=dx, g=g, h=h)
    plt.scatter(range(len(zs)), zs, edgecolor='k', facecolor='none',
                marker='+', lw=1)
    plt.plot(data, color='b')
    plt.show()

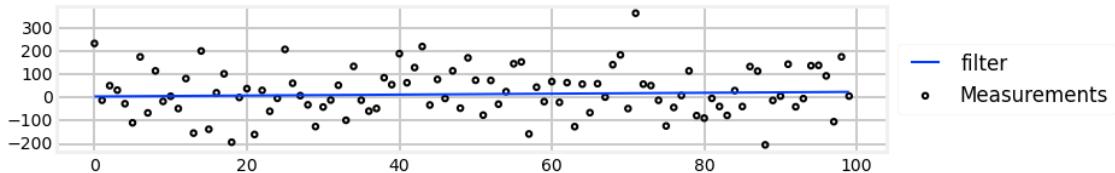
interact(interactive_gh,
         x=FloatSlider(value=0., min=-50, max=50.),
         dx=FloatSlider(value=5., min=-50., max=50.),
         g=FloatSlider(value=0.1, min=0.01, max=2, step=.02),
         h=FloatSlider(value=0.02, min=0.0, max=0.5, step=0.01));
```



## 1.14 Don't Lie to the Filter

You are free to set  $g$  and  $h$  to any value. Here is a filter that performs perfectly despite extreme noise.

```
In [44]: zs = gen_data(x0=5, dx=.2, count=100, noise_factor=100)
        data = g_h_filter(data=zs, x0=5, dx=.2, dt=1., g=0, h=0)
        with book_format.figsize(y=2.):
            book_plots.plot_measurements(zs)
            book_plots.plot_filter(data, label='filter', lw=2)
            book_plots.show_legend()
```

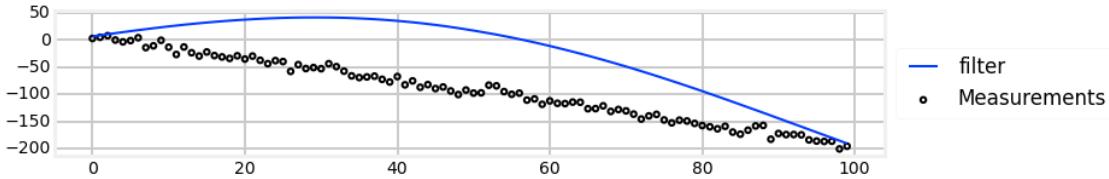


I brilliantly extracted a straight line out of very noisy data! Maybe I shouldn't try to collect my Fields Medal in mathematics just yet. I did this by setting both  $g$  and  $h$  to 0. What does this do? It makes the filter ignore the measurements, and so for each update it computes the new position as  $x + \Delta x \Delta t$ . Of course the result is a straight line if we ignore the measurements.

A filter that ignores measurements is useless. I know you would never set both  $g$  and  $h$  to zero as that takes a special kind of genius that only I possess, but I promise that if you are not careful you will set them lower than they should be. You can always make great looking results from test data. When you try your filter on different data you will be disappointed in the results because you finely tuned the constants for a specific data set.  $g$  and  $h$  must reflect the real world behavior of the system you are filtering, not the behavior of one specific data set. In later chapters we will learn a lot about how to do that. For now I can only say be careful, or you will be getting perfect results with your test data, but results like this once you switch to real data:

```
In [45]: zs = gen_data(x0=5, dx=-2, count=100, noise_factor=5)
        data = g_h_filter(data=zs, x0=5, dx=2., dt=1., g=.005, h=0.001)
```

```
with book_format.figsize(y=2.):
    book_plots.plot_measurements(zs)
    book_plots.plot_filter(data, label='filter', lw=2)
    book_plots.show_legend()
```



## 1.15 Tracking a Train

We are ready for a practical example. Earlier in the chapter we talked about tracking a train. Trains are heavy and slow, thus they cannot change speed quickly. They are on a track, so they cannot change direction except by slowing to a stop and then reversing course. Hence, we can conclude that if we already know the train's approximate position and velocity then we can predict its position in the near future with a great deal of accuracy. A train cannot change its velocity much in a second or two.

So let's write a filter for a train. Its position is expressed as its position on the track in relation to some fixed point which we say is 0 km. I.e., a position of 1 means that the train is 1 km away from the fixed point. Velocity is expressed as meters per second. We perform measurement of position once per second, and the error is  $\pm 500$  meters. How should we implement our filter?

First, let's simulate the situation without a filter. We will assume that the train is currently at kilometer 23, and moving at 15 m/s. We can code this as

```
pos = 23*1000
vel = 15
```

Now we can compute the position of the train at some future time, assuming no change in velocity, with

```
def compute_new_position(pos, vel, dt=1):
    return pos + (vel * dt)
```

We can simulate the measurement by adding in some random noise to the position. Here our error is 500m, so the code might look like:

```
def measure_position(pos):
    return pos + random.randn()*500
```

Let's put that in a cell and plot the results of 100 seconds of simulation. I will use NumPy's `asarray` function to convert the data into an NumPy array. This will allow me to divide all of the elements of the array at once by using the '/' operator.

In [46]: `from numpy.random import randn`

```
def compute_new_position(pos, vel, dt=1):
```

```

""" dt is the time delta in seconds."""
return pos + (vel * dt)

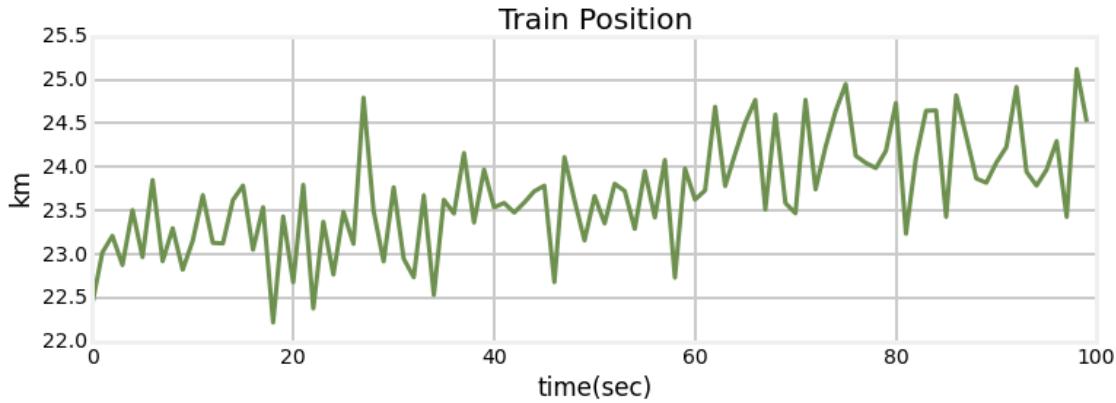
def measure_position(pos):
    return pos + randn()*500

def gen_train_data(pos, vel, count):
    zs = []
    for t in range(count):
        pos = compute_new_position(pos, vel)
        zs.append(measure_position(pos))
    return np.asarray(zs)

pos, vel = 23*1000, 15
zs = gen_train_data(pos, vel, 100)

plt.plot(zs / 1000.) # convert to km
book_plots.set_labels('Train Position', 'time(sec)', 'km')

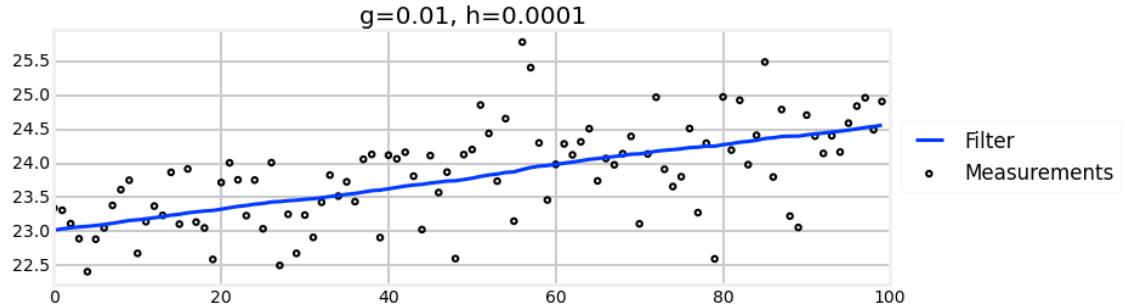
```



We can see from the chart how poor the measurements are. No real train could ever move like that.

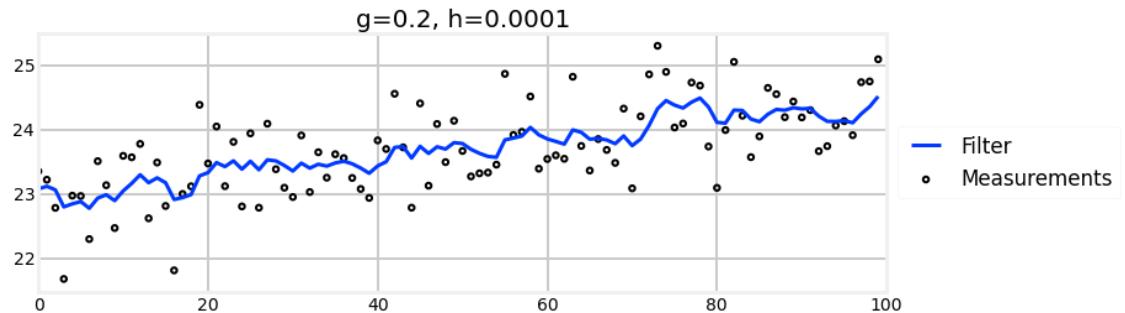
So what should we set  $g$  and  $h$  to if we want to filter this data? We have not developed the theory for this, but let's try to get a reasonable answer by the seat of our pants. We know that the measurements are very inaccurate, so we don't want to give them much weight at all. To do this we need to choose a very small  $g$ . We also know that trains can not accelerate or decelerate quickly, so we also want a very small  $h$ . For example:

```
In [47]: zs = gen_train_data(pos, 15, 100)
data = g_h_filter(data=zs, x0=pos, dx=15., dt=1., g=.01, h=0.0001)
plot_g_h_results(zs/1000, data/1000, 'g=0.01, h=0.0001')
```



That is pretty good for an initial guess. Lets make  $g$  larger to see the effect.

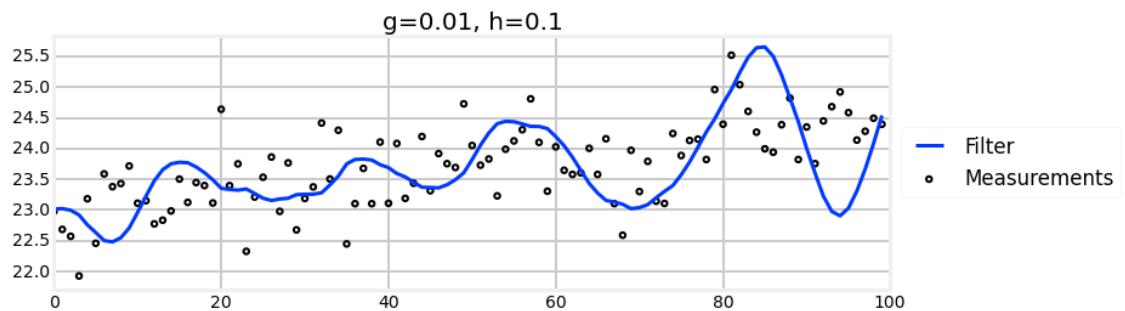
```
In [48]: zs = gen_train_data(pos, 15, 100)
        data = g_h_filter(data=zs, x0=pos, dx=15., dt=1., g=.2, h=0.0001)
        plot_g_h_results(zs/1000, data/1000, 'g=0.2, h=0.0001')
```



We made  $g=0.2$  and we can see that while the train's position is smoothed, the estimated position (and hence velocity) fluctuates a lot in a very tiny frame, far more than a real train can do. So empirically we know that we want  $g \ll 0.2$ .

Now let's see the effect of a poor choice for  $h$ .

```
In [49]: zs = gen_train_data(pos, 15, 100)
        data = g_h_filter(data=zs, x0=pos, dx=15., dt=1., g=0.01, h=0.1)
        plot_g_h_results(zs/1000, data/1000, 'g=0.01, h=0.1')
```

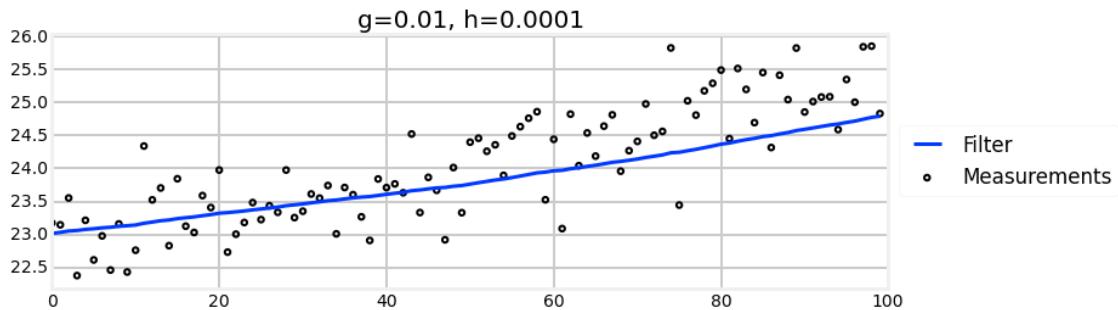


Here the position changes smoothly thanks to the small  $g$ , but the large  $h$  makes the filter very reactive to the measurements. This happens because in the course of a few seconds the rapidly changing measurement implies a very large velocity change, and a large  $h$  tells the filter to react to those changes quickly. Trains cannot change velocity quickly, so the filter is not doing a good job of filtering the data - the filter is changing velocity faster than a train can.

Finally, let's add some acceleration to the train. I don't know how fast a train can actually accelerate, but let's say it is accelerating at  $0.2 \text{ m/sec}^2$ .

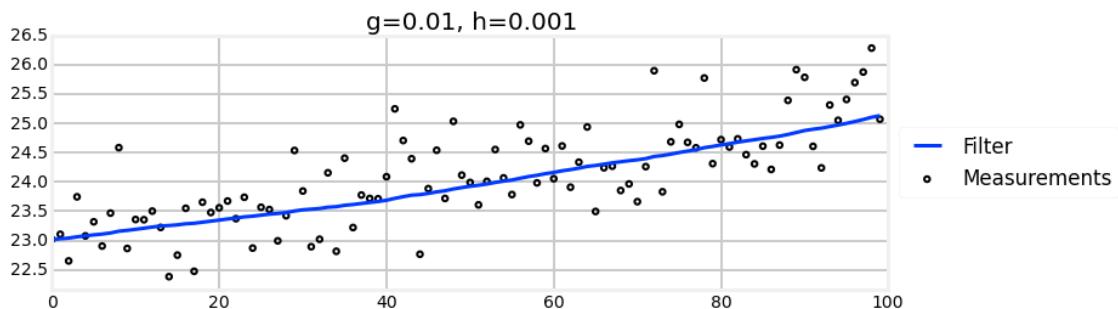
```
In [50]: def gen_train_data_with_acc(pos, vel, count):
    zs = []
    for t in range(count):
        pos = compute_new_position(pos, vel)
        vel += 0.2
        zs.append(measure_position(pos))
    return np.asarray(zs)
```

```
In [51]: zs = gen_train_data_with_acc(pos, 15, 100)
data = g_h_filter(data=zs, x0=pos, dx=15., dt=1., g=.01, h=0.0001)
plot_g_h_results(zs/1000, data/1000, 'g=0.01, h=0.0001')
```



Here we see that the filter is not quite tracking the train anymore due to the acceleration. We can fiddle with  $h$  to let it track better, at the expense of a less smooth filtered estimate.

```
In [52]: zs = gen_train_data_with_acc(pos, 15, 100)
data = g_h_filter(data=zs, x0=pos, dx=15., dt=1., g=.01, h=0.001)
plot_g_h_results(zs/1000, data/1000, 'g=0.01, h=0.001')
```



There are two lessons to be learned here. First, use the  $h$  term to respond to changes in velocity that you are not modeling. But, far more importantly, there is a trade off here between responding quickly and accurately to changes in behavior and producing ideal output for when the system is in a steady state that you have. If the train never changes velocity we would make  $h$  extremely small to avoid having the filtered estimate unduly affected by the noise in the measurement. But in an interesting problem there are almost always changes in state, and we want to react to them quickly. The more quickly we react to them, the more we are affected by the noise in the sensors.

I could go on, but my aim is not to develop g-h filter theory here so much as to build insight into how combining measurements and predictions leads to a filtered solution. There is extensive literature on choosing  $g$  and  $h$  for problems such as this, and there are optimal ways of choosing them to achieve various goals. As I explained earlier it is easy to ‘lie’ to the filter when experimenting with test data like this. In the subsequent chapters we will learn how the Kalman filter solves this problem in the same basic manner, but with far more sophisticated mathematics.

## 1.16 g-h Filters with FilterPy

FilterPy is an open source filtering library that I wrote. It has all of the filters in this book, along with others. It is rather easy to program your own g-h filter, but as we progress we will rely on FilterPy more. As a quick introduction, let’s look at the g-h filter in FilterPy.

If you do not have FilterPy installed just issue the following command from the command line.

```
pip install filterpy
```

Read Appendix A for more information on installing or downloading FilterPy from GitHub.

To use the g-h filter import it and create an object from the class `GHFilter`.

```
In [53]: from filterpy.gh import GHFilter
f = GHFilter(x=0., dx=0., dt=1., g=.8, h=.2)
```

When you construct the object you specify the initial value and rate of change for the signal (`x` and ‘`dx`’), the time step between updates(`dt`) and the two filter parameter (`g` and `h`). `dx` must have the same units of `x/dt` - if `x` is in meters and `dt` is in seconds then `dx` must be in meters per second.

To run the filter call `update`, passing the measurement in the parameter `z`.

```
In [54]: f.update(z=1.2)
```

```
Out[54]: (0.96, 0.24)
```

`update()` returns the new value of `x` and `dx` in a tuple, but you can also access them from the object.

```
In [55]: print(f.x, f.dx)
```

```
0.96 0.24
```

You can dynamically alter `g` and `h`.

```
In [56]: print(f.update(z=2.1, g=.85, h=.15))
```

```
(1.965, 0.375)
```

You can filter a sequence of measurements in a batch.

```
In [57]: print(f.batch_filter([3., 4., 5.]))
```

```
[[ 1.965  0.375]
 [ 2.868  0.507]
 [ 3.875  0.632]
 [ 4.901  0.731]]
```

You can filter multiple independent variables. If you are tracking an aircraft you'll need to track it in 3D space. Use NumPy arrays for `x`, `dx`, and the measurements.

```
In [58]: x_0 = np.array([1., 10., 100.])
dx_0 = np.array([10., 12., .2])

f_air = GHFilter(x=x_0, dx=dx_0, dt=1., g=.8, h=.2)
f_air.update(z=np.array((2., 11., 102.)))
print('x =', f_air.x)
print('dx =', f_air.dx)

x = [ 3.8    13.2   101.64]
dx = [ 8.2    9.8    0.56]
```

The class `GHFilterOrder` allows you to create a filter of order 0, 1, or 2. A g-h filter is order 1. The g-h-k filter, which we haven't talked about, also tracks accelerations. Both classes have functionality required by real applications, such as computing the Variance Reduction Factor (VRF), which we haven't discussed in this chapter. I could fill a book just on the theory and applications of g-h filters, but we have other goals in this book. If you are interested, explore the FilterPy code and do some further reading.

The documentation for FilterPy is at <https://filterpy.readthedocs.org/>.

## 1.17 Summary

I encourage you to experiment with this filter to develop your understanding of how it reacts. It shouldn't take too many attempts to come to the realization that ad-hoc choices for  $g$  and  $h$  do not perform very well. A particular choice might perform well in one situation, but very poorly in another. Even when you understand the effect of  $g$  and  $h$  it can be difficult to choose proper values. In fact, it is extremely unlikely that you will choose values for  $g$  and  $h$  that is optimal for any given problem. Filters are designed, not selected ad hoc.

In some ways I do not want to end the chapter here, as there is a significant amount that we can say about selecting  $g$  and  $h$ . But the g-h filter in this form is not the purpose of this book. Designing the Kalman filter requires you to specify a number of parameters - indirectly they do relate to choosing  $g$  and  $h$ , but you will never refer to them directly when designing Kalman filters. Furthermore,  $g$  and  $h$  will vary at every time step in a very non-obvious manner.

There is another feature of these filters we have barely touched upon - Bayesian statistics. You will note that the term 'Bayesian' is in the title of this book; this is not a coincidence! For the time being we will leave  $g$  and  $h$  behind, largely unexplored, and develop a very powerful form of probabilistic reasoning about filtering. Yet suddenly this same g-h filter algorithm will appear, this time with a formal mathematical edifice that allows us to create filters from multiple sensors, to accurately estimate the amount of error in our solution, and to control robots.

## 1.18 References

- [1] NASA Kalman Filtering Presentation

[http://nescacademy.nasa.gov/review/downloadfile.php?file=FundamentalsofKalmanFiltering\\_Presentation.pdf&id=199](http://nescacademy.nasa.gov/review/downloadfile.php?file=FundamentalsofKalmanFiltering_Presentation.pdf&id=199)

# Chapter 2

## Discrete Bayes Filter

The Kalman filter belongs to a family of filters called Bayesian filters. Most textbook treatments of the Kalman filter present the Bayesian formula, perhaps shows how it factors into the Kalman filter equations, but mostly keeps the discussion at a very abstract level.

That approach requires a fairly sophisticated understanding of several fields of mathematics, and it still leaves much of the work of understanding and forming an intuitive grasp of the situation in the hands of the reader.

I will use a different way to develop the topic, to which I owe the work of Dieter Fox and Sebastian Thrun a great debt. It depends on building an intuition on how Bayesian statistics work by tracking an object through a hallway - they use a robot, I use a dog. I like dogs, and they are less predictable than robots which imposes interesting difficulties for filtering. The first published example of this that I can find seems to be Fox 1999 [1], with a fuller example in Fox 2003 [2]. Sebastian Thrun also uses this formulation in his excellent Udacity course Artificial Intelligence for Robotics [3]. In fact, if you like watching videos, I highly recommend pausing reading this book in favor of first few lessons of that course, and then come back to this book for a deeper dive into the topic.

Let's now use a simple thought experiment, much like we did with the g-h filter, to see how we might reason about the use of probabilities for filtering and tracking.

### 2.1 Tracking a Dog

Let's begin with a simple problem. We have a dog friendly workspace, and so people bring their dogs to work. Occasionally the dogs wander out of offices and down the halls. We want to be able to track them. So during a hackathon somebody created a little sonar sensor to attach to the dog's collar. It emits a signal, listens for the echo, and based on how quickly an echo comes back we can tell whether the dog is in front of an open doorway or not. It also senses when the dog walks, and reports in which direction the dog has moved. It connects to our network via wifi and sends an update once a second.

I want to track my dog Simon, so I attach the device to his collar and then fire up Python, ready to write code to track him through the building. At first blush this may appear impossible. If I start listening to the sensor of Simon's collar I might read 'door', 'hall', 'hall', and so on. How can I use that information to determine where Simon is?

To keep the problem small enough to plot easily we will assume that there are only 10 positions in a single hallway to consider, which we will number 0 to 9, where 1 is to the right of 0, 2 is to the right of 1, and so on. For reasons that will be clear later, we will also assume that the hallway is circular or rectangular. If you move right from position 9, you will be at position 0.

When I begin listening to the sensor I have no reason to believe that Simon is at any particular position in

the hallway. He is equally likely to be in any position. There are 10 positions, so the probability that he is in any given position is 1/10.

Let's represent our belief of his position at any time in a NumPy array.

```
In [2]: import numpy as np
belief = np.array([1/10]*10)
np.set_printoptions(precision=3, linewidth=50)
print(belief)
```

```
[ 0.1  0.1  0.1  0.1  0.1  0.1  0.1  0.1  0.1  0.1]
```

In Bayesian statistics this is called our prior. It basically means the probability prior to incorporating measurements or other information. More completely, this is called the prior probability distribution, but that is a mouthful and so it is normally shorted to prior. A probability distribution is a collection of all possible probabilities for an event. Probability distributions always sum to 1 because something had to happen; the distribution lists all the different somethings and the probability of each.

I'm sure you've used probabilities before - as in "the probability of rain today is 30%". The last paragraph sounds like more of that. But Bayesian statistics was a revolution in probability because it treats the probability as a belief about a single event. Let's take an example. I know if I flip a fair coin a infinitely many times 50% of the flips will be heads and 50% tails. That is standard probability, not Bayesian, and is called frequentist statistics to distinguish it from Bayes. Now, let's say I flip the coin one more time. Which way do I believe it landed? Frequentist probability has nothing to say about that; it will merely state that 50% of coin flips land as heads. Bayes treats this as a belief about a single event - the strength of my belief that this specific coin flip is heads is 50%.

There are more differences, but for now recognize that Bayesian statistics reasons about our belief about a single event, whereas frequentists reason about collections of events. In this rest of this chapter, and most of the book, when I talk about the probability of something I am referring to the probability that some specific thing is true. When I do that I'm taking the Bayesian approach.

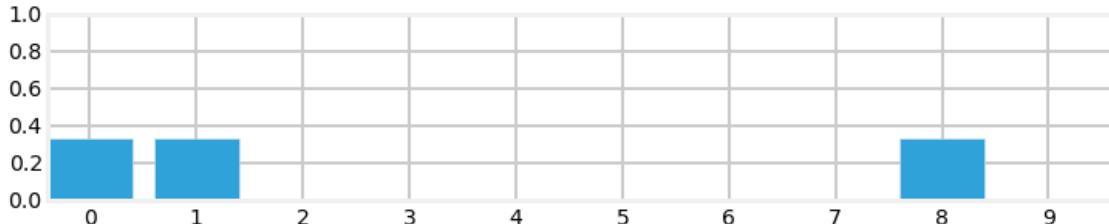
Now let's create a map of the hallway in another list. Suppose there are first two doors close together, and then another door quite a bit further down the hallway. We will use 1 to denote a door, and 0 to denote a wall:

```
In [3]: hallway = np.array([1, 1, 0, 0, 0, 0, 0, 0, 1, 0])
```

So I start listening to Simon's transmissions on the network, and the first data I get from the sensor is "door". For the moment assume the sensor always returns the correct answer. From this I conclude that he is in front of a door, but which one? I have no idea. I have no reason to believe he is in front of the first, second, or third door. But what I can do is assign a probability to each door. All doors are equally likely, and there are three of them, so I assign a probability of 1/3 to each door.

```
In [4]: from book_format import figsize, set_figsize
import book_plots as bp
import matplotlib.pyplot as plt

belief = np.array([1./3, 1./3, 0, 0, 0, 0, 0, 0, 1./3, 0])
set_figsize(y=2)
bp.bar_plot(belief)
```



This distribution is called a categorical distribution, which is the term used to describe a discrete distribution describing the probability of observing  $n$  outcomes. It is a multimodal distribution because we have multiple beliefs about the position of our dog. Of course we are not saying that we think he is simultaneously in three different locations, merely that so far we have narrowed down our knowledge in his position to be one of these three locations. My (Bayesian) belief is that there is a 33.3% chance of being at door 0, 33.3% at door 1, and a 33.3% chance of being at door 8.

A few words about the mode of a distribution. This terms from elementary statistics. Given a set of numbers, such as  $\{1, 2, 2, 2, 3, 3, 4\}$ , the mode is the number that occurs most often. For this set the mode is 2. A set can contain more than one mode. The set  $\{1, 2, 2, 2, 3, 3, 4, 4, 4\}$  contains the modes 2 and 4, because both occur three times. We say the former set is unimodal, and the latter is multimodal.

I hand coded the `belief` array in the code above. How would we implement this in code? Well, hallway represents each door as a 1, and wall as 0, so we will multiply the hallway variable by the percentage, like so;

```
In [5]: pbelief = hallway * (1/3)
        print(pbelief)

[ 0.333  0.333  0.      0.      0.      0.      0.
  0.333  0.     ]
```

## 2.2 Extracting Information from Multiple Sensor Readings

Let's put Python aside and think about the problem a bit. Suppose we were to read the following from Simon's sensor:

- door
- move right
- door

Can we deduce where Simon is at the end of that sequence? Of course! Given the hallway's layout there is only one place where you can be in front of a door, move once to the right, and be in front of another door, and that is at the left end. Therefore we can confidently state that Simon is in front of the second doorway. If this is not clear, suppose Simon had started at the second or third door. After moving to the right, his sensor would have returned 'wall'. That doesn't match the sensor readings, so we know he didn't start there. We can continue with that logic for all the remaining starting positions. Therefore the only possibility is that he is now in front of the second door. We implement this in Python with:

```
In [6]: belief = np.array([0., 1., 0., 0., 0., 0., 0., 0., 0., 0.])
```

Obviously I carefully constructed the hallway layout and sensor readings to give us an exact answer quickly. Real problems will not be so clear cut. But this should trigger your intuition - the first sensor reading only gave us very low probabilities (0.333) for Simon's location, but after a position update and another sensor reading we knew much more about where he is. You might suspect, correctly, that if you had a very long hallway with a large number of doors that after several sensor readings and positions updates we would either be able to know where Simon was, or have the possibilities narrowed down to a small number of possibilities. For example, suppose we had a long sequence of "door, right, door, right, wall, right, wall, right, door, right, door, right, wall, right, wall, right, wall, right, door". Simon could only have started in a location where his movements had a door sequence of [1,1,0,0,1,1,0,0,0,0,1] in the hallway. There might be only one match for that, or at most a few. Either way we will be far more certain about his position then when we started.

We could implement this solution now, but instead let us consider a real world complication to the problem.

## 2.3 Noisy Sensors

Perfect sensors are rare. Perhaps the sensor would not detect a door if Simon sat in front of it while scratching himself, or it might report there is a door if he is facing towards the wall instead of down the hallway. So in practice when I get a report 'door' I cannot assign 1/3 as the probability for each door. I have to assign something less than 1/3 to each door, and then assign a small probability to each blank wall position.

At this point it doesn't matter exactly what numbers we assign; let us say that the probability of the sensor being right is 3 times more likely to be right than wrong. How would we do this?

At first this may seem like an insurmountable problem. If the sensor is noisy it casts doubt on every piece of data. How can we conclude anything if we are always unsure?

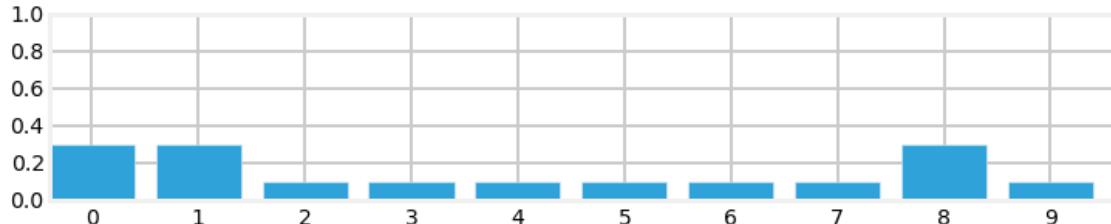
The answer, as with the problem above, is probabilities. We are already comfortable assigning a probabilistic belief about the location of the dog; now we have to incorporate the additional uncertainty caused by the sensor noise. Lets say we get a reading of 'door'. We already said that the sensor is three times as likely to be correct as incorrect, so we should scale the probability distribution by 3 where there is a door. If we do that the result will no longer be a probability distribution, but we will learn how to correct that in a moment.

Let's look at that in Python code. Here I use the variable `z` to denote the measurement as that is the customary choice in the literature (`y` is also commonly used).

```
In [7]: def update(map_, belief, z, correct_scale):
    for i, val in enumerate(map_):
        if val == z:
            belief[i] *= correct_scale

    belief = np.array([0.1] * 10)
    reading = 1 # 1 is 'door'
    update(hallway, belief, z=reading, correct_scale=3.)
    print('sum =', sum(belief))
    bp.bar_plot(belief)

sum = 1.6
```



We can see that this is not a probability distribution because it does not sum to 1.0. But we can see that the code is doing mostly the right thing - the doors are assigned a number (0.3) that is 3 times higher than the walls (0.1). All we need to do is normalize the result so that the probabilities correctly sum to 1.0. Normalization is done by dividing each element by the sum of all elements in the list.

Also, it is a bit odd to be talking about “3 times as likely to be right as wrong”. We are working in probabilities, so let’s specify the probability of the sensor being correct, and computing the scale factor from that.

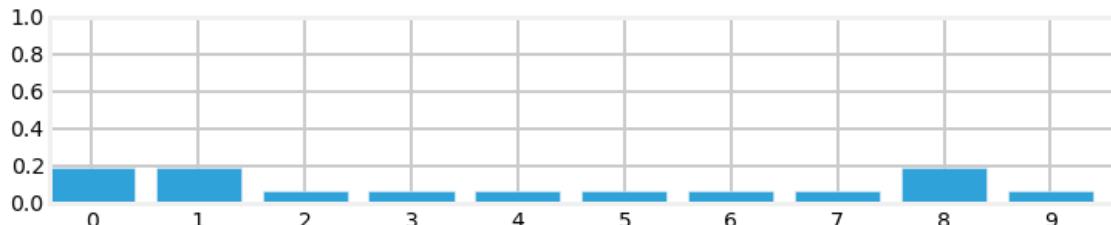
```
In [8]: def normalize(distribution):
    assert distribution.dtype.kind == 'f'
    """ Normalize distribution so it sums to 1.0 """
    distribution /= sum(distribution.astype(float))

    def update(map_, belief, z, prob_correct):
        scale = prob_correct / (1. - prob_correct)
        for i, val in enumerate(map_):
            if val == z:
                belief[i] *= scale
        normalize(belief)

    belief = np.array([0.1] * 10)
    update(hallway, belief, 1, prob_correct=.75)

    print('sum =', sum(belief))
    print('probability of door =', belief[0])
    print('probability of wall =', belief[2])
    bp.bar_plot(belief)

sum = 1.0
probability of door = 0.1875
probability of wall = 0.0625
```



We can see from the output that the sum is now 1.0, and that the probability of a door vs wall is still three times larger. The result also fits our intuition that the probability of a door must be less than 0.333, and that the probability of a wall must be greater than 0.0. Finally, it should fit our intuition that we have not yet been given any information that would allow us to distinguish between any given door or wall position, so all door positions should have the same value, and the same should be true for wall positions.

This result is called the posterior, which is short for posterior probability distribution. All this means is a probability distribution after incorporating the measurement information (posterior means ‘after’ in this context). To review, the prior is the probability distribution before including the measurement’s information. Another term is the likelihood - this is the probability for the evidence (the measurement in this case) being true. That gives us this equation:

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{normalization}}$$

It is very important to learn and internalize these terms as most of the literature uses them exclusively.

## 2.4 Incorporating Movement Data

Recall how quickly we were able to find an exact solution to our dog’s position when we incorporated a series of measurements and movement updates. However, that occurred in a fictional world of perfect sensors. Might we be able to find an exact solution even in the presence of noisy sensors?

Unfortunately, the answer is no. Even if the sensor readings perfectly match an extremely complicated hallway map we could not say that we are 100% sure that the dog is in a specific position - there is, after all, the possibility that every sensor reading was wrong! Naturally, in a more typical situation most sensor readings will be correct, and we might be close to 100% sure of our answer, but never 100% sure. This may seem complicated, but lets go ahead and program the math, which as we have seen is quite simple.

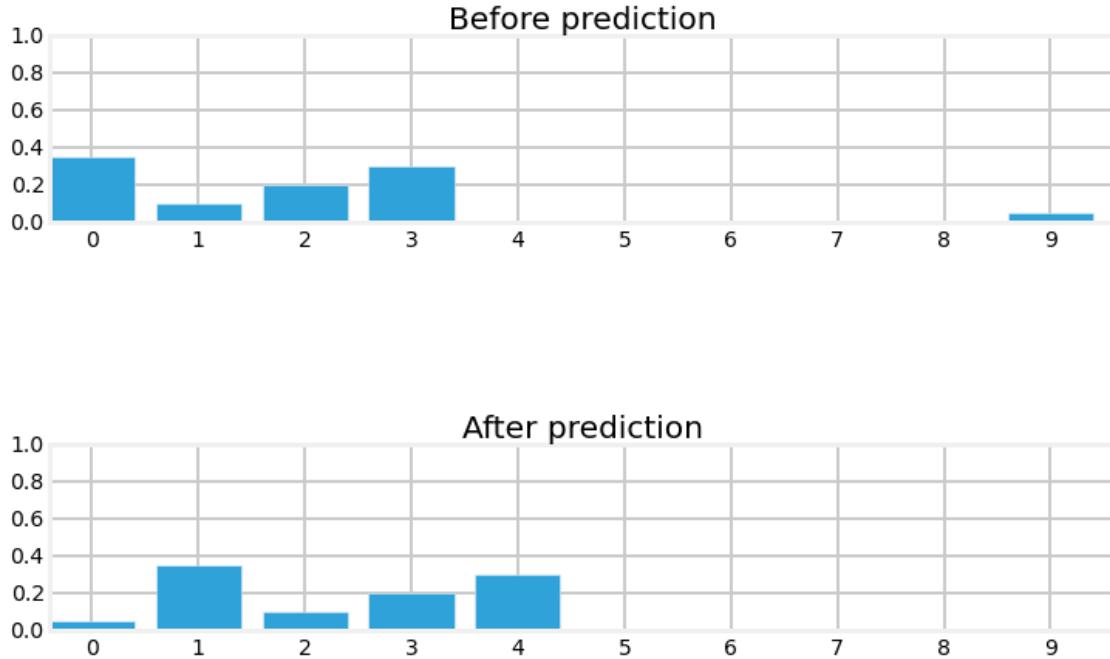
First let’s deal with the simple case - assume the movement sensor is perfect, and it reports that the dog has moved one space to the right. How would we alter our **belief** array?

I hope after a moment’s thought it is clear that we should shift all the values one space to the right. If we previously thought there was a 50% chance of Simon being at position 3, then after the move to the right we should believe that there is a 50% chance he is at position 4. So let’s implement that. Recall that the hallway is circular, so we will use modulo arithmetic to perform the shift correctly.

```
In [9]: def perfect_predict(belief, move):
    """ move the position by 'move' spaces, where positive is
    to the right, and negative is to the left
    """
    n = len(belief)
    result = np.zeros(n)
    for i in range(n):
        result[i] = belief[(i-move) % n]
    belief[:] = result # copy back to original array

    belief = np.array([.35, .1, .2, .3, 0, 0, 0, 0, 0, .05])
    bp.bar_plot(belief, title='Before prediction')
    plt.show()

    perfect_predict(belief, 1)
    bp.bar_plot(belief, title='After prediction')
```



We can see that we correctly shifted all values one position to the right, wrapping from the end of the array back to the beginning.

## 2.5 Terminology

I introduced this terminology in the last chapter, but let's review to help solidify your knowledge.

The system is what we are trying to model or filter. Here the system is our dog. The state is its current configuration or value. In this chapter the state is our dog's position. We rarely know the actual state, so we say our filters produce the estimated state of the system. In practice this often gets called the state, so be careful to understand if the reference is to the state of the filter (which is the estimated state), or the state of the system (which is the actual state).

One cycle of prediction and updating with a measurement is called the state or system evolution, which is short for time evolution [7]. Another term is system propagation. This term refers to how the state of the system changes over time. For filters, this time is usually discrete. For our dog tracking, the system state is the position of the dog, so the state evolution is the position after a discrete amount of time has passed.

We model the system behavior with the process model. Here, our process model is that the dog moves one position at each time step. This is not a particularly accurate model of how dog's behave. The error in the model is called the system error or process error.

More terminology - the prediction is our new prior. Time has moved forward and we made a prediction without benefit of knowing the measurements.

## 2.6 Adding Noise to the Prediction

We want to solve real world problems, and we have already stated that all sensors have noise. Therefore the code above must be wrong since it assumes we have perfect measurements. What if the sensor reported that

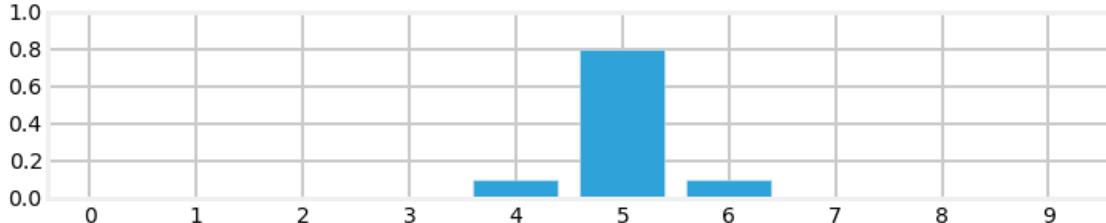
our dog moved one space, but he actually moved two spaces, or zero? Once again this may initially sound like an insurmountable problem, but let's model it and see what happens. Since this is an example we will create a simple noise model for the sensor - later in the book we will handle more difficult errors.

We will say that when the sensor sends a movement update, it is 80% likely to be right, and it is 10% likely to overshoot one position to the right, and 10% likely to undershoot to the left. That is, if we say the movement was 4 (meaning 4 spaces to the right), the dog is 80% likely to have moved 4 spaces to the right, 10% to have moved 3 spaces, and 10% to have moved 5 spaces.

This is slightly harder than the math we have done so far, but it is still tractable. Each result in the array now needs to incorporate probabilities for 3 different situations. For example, consider position 9 for the case where the reported movement is 2. It should be clear that after the move we need to incorporate the probability that was at position 7 (9-2). However, there is a small chance that our dog actually moved from either 1 or 3 spaces away due to the sensor noise, so we also need to use positions 6 and 8. How much? Well, we have the probabilities, so we can just multiply and add. It would be 80% of position 7 plus 10% of position 6 and 10% of position 8! Let's try coding that:

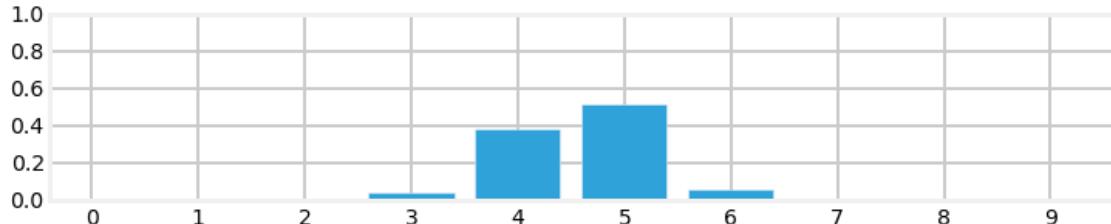
```
In [10]: def predict(belief, move, p_correct, p_under, p_over):
    n = len(belief)
    result = np.zeros(n)
    for i in range(n):
        result[i] = (
            belief[(i-move) % n] * p_correct +
            belief[(i-move-1) % n] * p_over +
            belief[(i-move+1) % n] * p_under)
    belief[:] = result

belief = np.array([0., 0., 0., 1., 0., 0., 0., 0., 0.])
predict(belief, move=2, p_correct=.8, p_under=.1, p_over=.1)
bp.bar_plot(belief)
```



The simple test case that we ran appears to work correctly. We initially believed that the dog was in position 3 with 100% certainty; after the movement update we now give an 80% probability to the dog being in position 5, and a 10% chance to undershooting to position 4, and a 10% chance of overshooting to position 6. Let us look at a case where we have multiple beliefs:

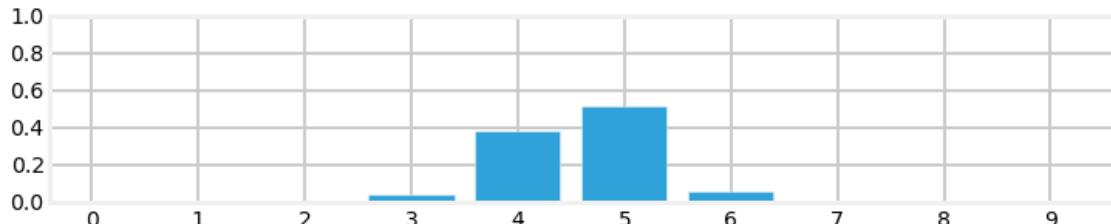
```
In [11]: belief = np.array([0, 0, .4, .6, 0, 0, 0, 0, 0])
predict(belief, move=2, p_correct=.8, p_under=.1, p_over=.1)
bp.bar_plot(belief)
```



Here the results are more complicated, but you should still be able to work it out in your head. The 0.04 is due to the possibility that the 0.4 belief undershot by 1. The 0.38 is due to the following: the 80% chance that we moved 2 positions ( $0.4 \times 0.8$ ) and the 10% chance that we undershot ( $0.6 \times 0.1$ ). Overshooting plays no role here because if we overshot both 0.4 and 0.6 would be past this position. **I strongly suggest working some examples until all of this is very clear, as so much of what follows depends on understanding this step.**

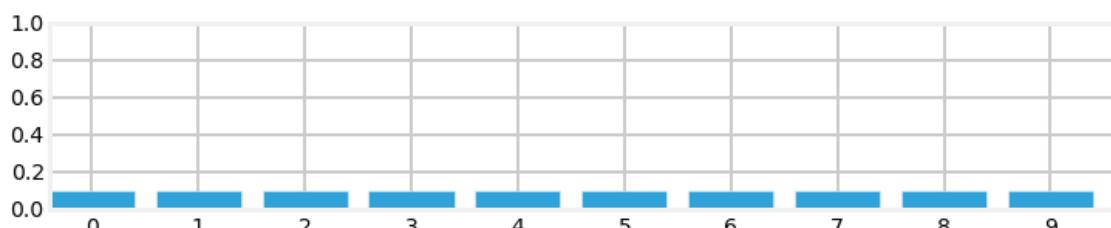
If you look at the probabilities after performing the update you probably feel dismay. In the example above we started with probabilities of 0.4 and 0.6 in two fields; after performing the update the probabilities are not only lowered, but they are strewn out across the map.

In [12]: `bp.bar_plot(belief)`



This is not a coincidence, or the result of a carefully chosen example - it is always true of the evolution (predict step). This is inevitable; if our sensor is noisy we will lose a bit of information on every prediction. Suppose we were to perform the prediction an infinite number of times - what would the result be? If we lose information on every step, we must eventually end up with no information at all, and our probabilities will be equally distributed across the `belief` array. Let's try this with 500 iterations.

In [13]: `belief = np.array([1.0, 0, 0, 0, 0, 0, 0, 0, 0, 0])  
for i in range(500):  
 predict(belief, move=1, p_correct=.8, p_under=.1, p_over=.1)  
bp.bar_plot(belief)`



After 500 iterations we have lost all information, even though we were 100% sure that we started in position 1. Feel free to play with the numbers to see the effect of differing number of updates. For example, after 100 updates we have a small amount of information left.

And, if you are viewing this on the web or in Jupyter Notebook, here is an animation of that output.

## 2.7 Generalizing with Convolution

In the code above we made the assumption that the movement error is only one position. In any real problem it will almost always be possible for the error to be two, three, or more positions.

This is easily solved with `convolution`. Convolution is the mathematical technique we use to modify one function with another function. In our case we are modifying our probability distribution with the error function that represents the probability of a movement error. In fact, the implementation of `predict()` is a convolution, though we did not call it that. Formally, convolution is defined as

$$(f * g)(t) = \int_0^t f(\tau) g(t - \tau) d\tau$$

where  $f * g$  is the notation for convolving  $f$  by  $g$ . It does not mean multiply.

Integrals are for continuous functions, but we are using discrete functions, so we replace the integral with a summation, and the parenthesis with array brackets.

$$(f * g)[t] = \sum_{\tau=0}^t f[\tau] g[t - \tau]$$

If you look at that equation and compare it to the `predict()` function you can see that they are doing the same thing.

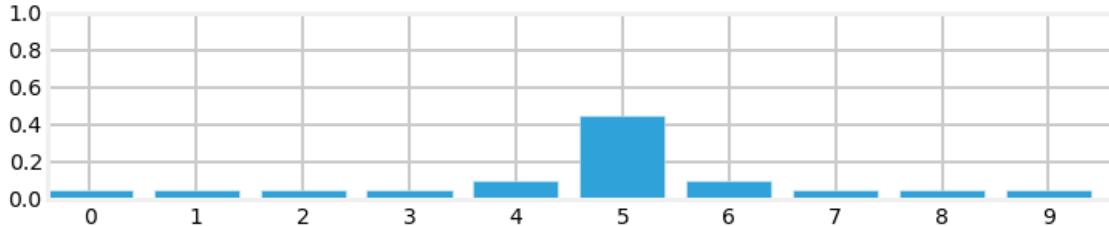
I would love to go on and on about convolution, but we have a filter to implement. Khan Academy [4] has a good mathematical introduction to convolution, and Wikipedia has some excellent animations of convolutions [5]. But the general idea is already clear to you. You slide across an array, multiplying the neighbors of the current cell with the values of a second array. This second array is called the `kernel`. In our example above we used 0.8 for the probability of moving to the correct location, 0.1 for undershooting, and 0.1 for overshooting. We make a kernel of this with the array `[0.1, 0.8, 0.1]`. So all we need to do is write a loop that goes over each element of our array, multiplying by the kernel, and summing the results. To emphasize that the belief is a probability distribution I have named the belief `distribution`.

```
In [14]: def predict(distribution, offset, kernel):
    N = len(distribution)
    kN = len(kernel)
    width = int((kN - 1) / 2)

    result = np.zeros(N)
    for i in range(N):
        for k in range(kN):
            index = (i + (width-k) - offset) % N
            result[i] += distribution[index] * kernel[k]
    distribution[:] = result[:] # update belief
```

We should test that this is correct:

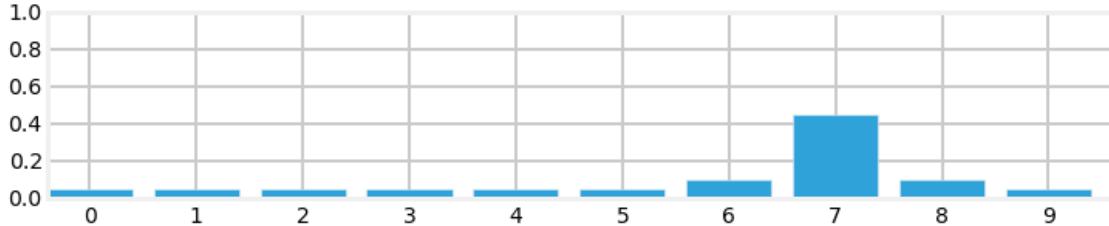
```
In [15]: belief = np.array([.05] * 10)
belief[4] = .55
predict(belief, offset=1, kernel=[.1, .8, .1])
bp.bar_plot(belief)
```



All of the elements are unchanged except the middle ones, which is correct. The value in position 4 should be  $(0.1 \times 0.05) + (0.8 \times 0.05) + (0.1 \times 0.55) = 0.1$ , which it does. Position 5 should be  $(0.1 \times 0.05) + (0.8 \times 0.55) + (0.1 \times 0.05) = 0.45$ , which it does. Finally, position 6 is computed the same as position 4, and it is also correct.

Finally, let's ensure that it shifts the positions correctly for movements greater than one.

```
In [16]: belief = np.array([.05] * 10)
belief[4] = .55
predict(belief, offset=3, kernel=[.1, .8, .1])
bp.bar_plot(belief)
```



The position was correctly shifted by 3 positions, so the code seems to be correct.

It will usually be true that a small error is very likely, and a large error is very unlikely. For example, a kernel might look like

[.02, .1, .2, .36, .2, .1, .02]

But maybe an overshoot is very likely, and an undershoot is impossible. That kernel might be

[.0, .0, .5, .3, .2]

You probably will never be able to compute exact probabilities for the kernel of a real problem. It is usually enough to be approximately right. What you must do is ensure that the terms of the kernel sum to 1. The kernel expresses the probability of any given move, and the sum of any probability distribution must be 1.

An easy way to do this is to make each term proportionally correct and then normalize it to sum to 1.

```
In [17]: kernel = np.array([1, 2, 6, 2, 1], dtype=float)
normalize(kernel)
print(kernel)

[ 0.083  0.167  0.5    0.167  0.083]
```

## 2.8 Integrating Measurements and Movement Updates

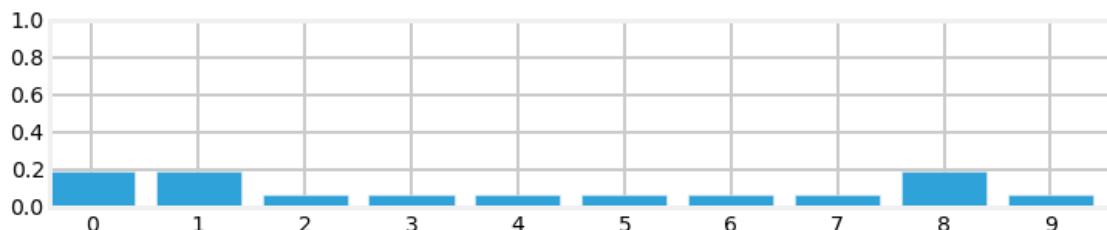
The problem of losing information during a prediction may make it seem as if our system would quickly devolve into no knowledge. However, recall that our process is not an endless series of predictions, but a cycle of predictions followed by updates. The output of the update step, where we measure the current position, is fed into the prediction. The prediction step, with a degraded certainty, is then fed back into the update step where we measure the position again.

Let's think about this intuitively. After the first update->predict round we have degraded the knowledge we gained by the measurement by a small amount. But now we take another measurement. When we try to incorporate that new measurement into our belief, do we become more certain, less certain, or equally certain? Consider a simple case - you are sitting in your office. A co-worker asks another co-worker where you are, and they report "in his office". You keep sitting there while they ask and answer "has he moved"? "No" "Where is he" "In his office". Eventually you get up and move, and lets say the person didn't see you move. At that time the questions will go "Has he moved" "no" (but you have!) "Where is he" "In the kitchen". Wow! At that moment the statement that you haven't moved conflicts strongly with the next measurement that you are in the kitchen. If we were modeling these with probabilities the probability that you are in your office would lower, and the probability that you are in the kitchen would go up a little bit. But now imagine the subsequent conversation: "has he moved" "no" "where is he" "in the kitchen". Pretty quickly the belief that you are in your office would fade away, and the belief that you are in the kitchen would increase to near certainty. The belief that you are in the office will never go to zero, nor will the belief that you are in the kitchen ever go to 1.0 because of the chances of error, but in practice your co-workers would be correct to be quite confident in their system.

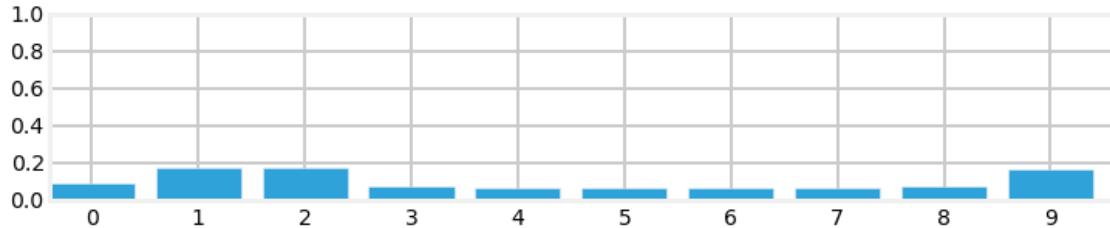
That is what intuition tells us. What does the math tell us?

Well, we have already programmed the update step, and we have programmed the predict step. All we need to do is feed the result of one into the other, and we will have programmed our dog tracker!!! Let's see how it performs. We will input measurements as if the dog started at position 0 and moved right at each update. However, as in a real world application, we will start with no knowledge and assign equal probability to all positions.

```
In [18]: hallway = np.array([1, 1, 0, 0, 0, 0, 0, 0, 1, 0])
belief = np.array(.1 * 10)
update(hallway, belief, z=1, prob_correct=.75)
bp.bar_plot(belief)
```

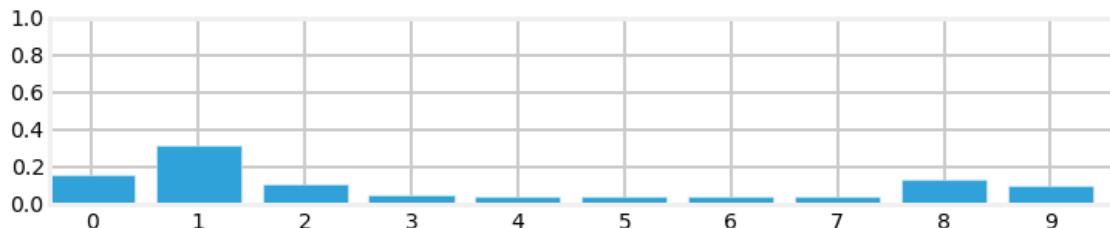


```
In [19]: kernel = (.1, .8, .1)
predict(belief, 1, kernel)
bp.bar_plot(belief)
```



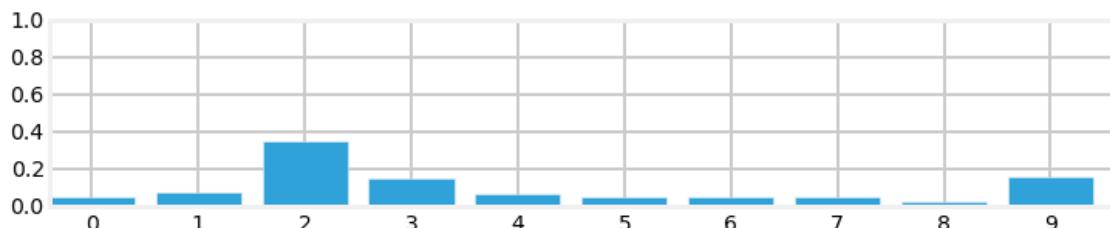
So after the first update we have assigned a high probability to each door position, and a low probability to each wall position. The predict step shifted these probabilities to the right, smearing them about a bit. Now let's look at what happens at the next sense.

```
In [20]: update(hallway, belief, z=1, prob_correct=.75)
bp.bar_plot(belief)
```



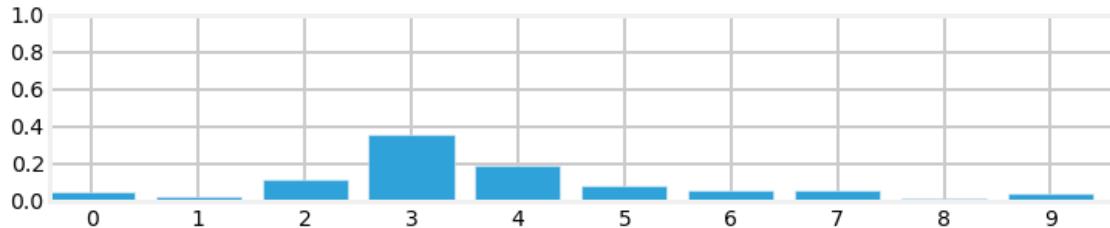
Notice the tall bar at position 1. This corresponds with the (correct) case of starting at position 0, sensing a door, shifting 1 to the right, and sensing another door. No other positions make this set of observations as likely. Now lets add an update and then sense the wall.

```
In [21]: predict(belief, 1, kernel)
update(hallway, belief, z=0, prob_correct=.75)
bp.bar_plot(belief)
```



This is exciting! We have a very prominent bar at position 2 with a value of around 35%. It is over twice the value of any other bar in the plot, and is about 4% larger than our last plot, where the tallest bar was around 31%. Let's see one more sense->update cycle.

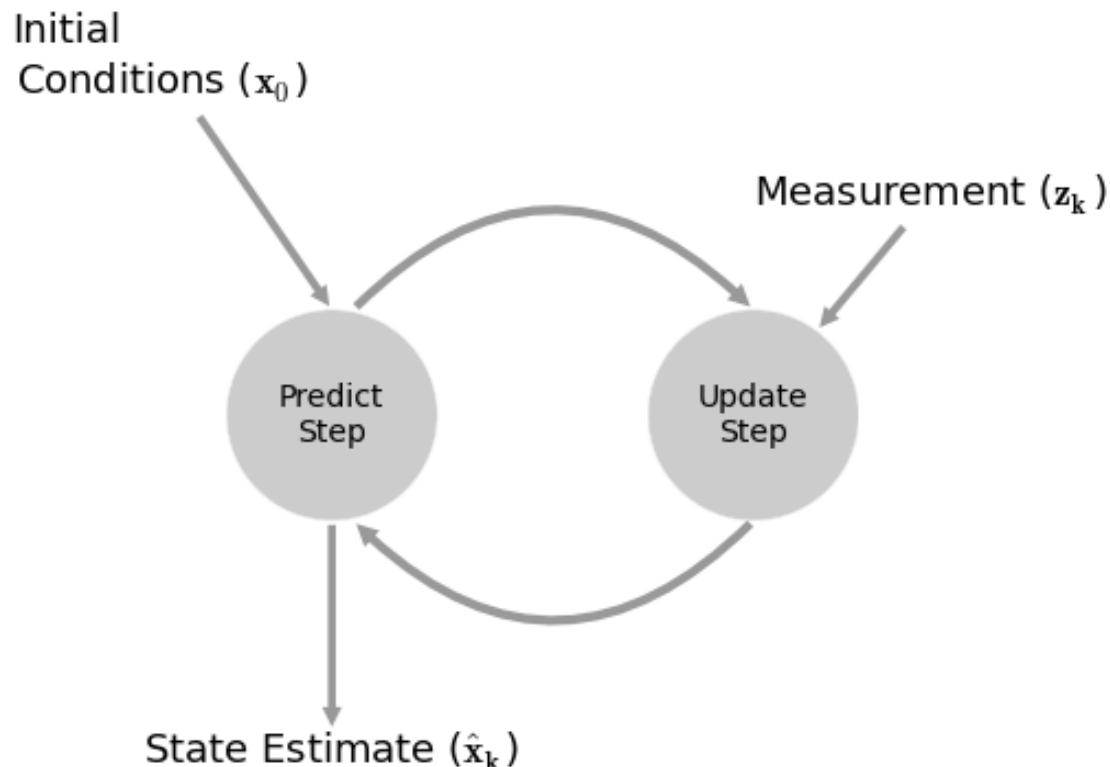
```
In [22]: predict(belief, 1, kernel)
        update(hallway, belief, z=0, prob_correct=.75)
        bp.bar_plot(belief)
```



Here things have degraded a bit due to the long string of wall positions in the map. We cannot be as sure where we are when there is an undifferentiated line of wall positions, so naturally our probabilities spread out a bit.

I spread the computation across several cells, iteratively calling `predict()` and `update()`. This chart from the **g-h Filter** chapter illustrates the algorithm.

```
In [23]: import gh_internal
        gh_internal.create_predict_update_chart()
```



This filter is a form of the g-h filter. Here we are using the percentages for the errors to implicitly compute the  $g$  and  $h$  parameters. We could express the discrete Bayes algorithm as a g-h filter, but that would obscure the logic of this filter.

We can express this in pseudocode.

### Initialization

1. Initialize our belief in the state

### Predict

1. Based on the system behavior, predict state at the next time step
2. Adjust belief to account for the uncertainty in prediction

### Update

1. Get a measurement and associated belief about its accuracy
2. Compute residual between estimated state and measurement
3. Determine whether the measurement matches each state
4. Update state belief if it matches the measurement

When we cover the Kalman filter we will use this exact same algorithm; only the details of the computation will differ.

Finally, for those viewing this in a Jupyter Notebook or on the web, here is an animation of that algorithm.

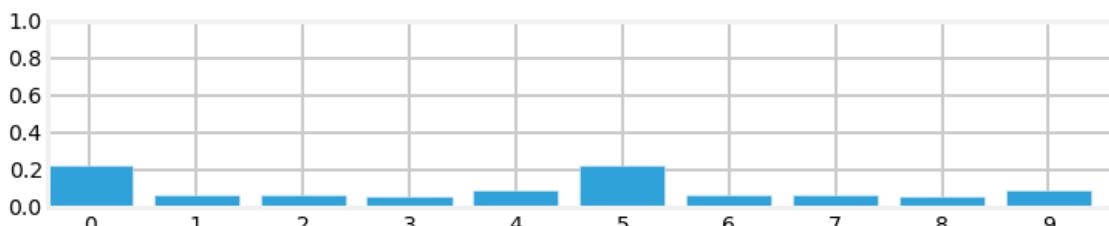
## 2.9 The Effect of Bad Sensor Data

You may be suspicious of the results above because I always passed correct sensor data into the functions. However, we are claiming that this code implements a filter - it should filter out bad sensor measurements. Does it do that?

To make this easy to program and visualize I will change the layout of the hallway to mostly alternating doors and hallways, and run the algorithm on 5 correct measurements:

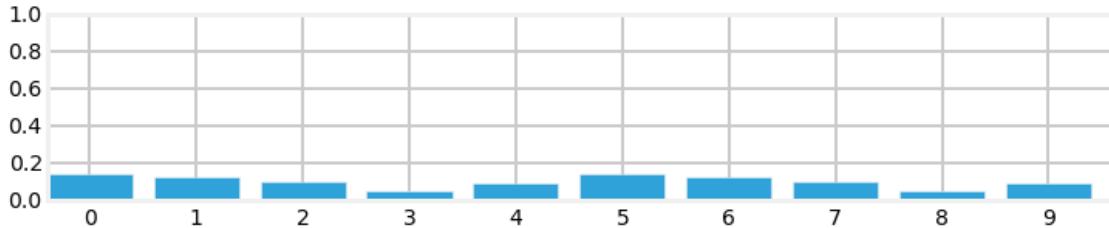
```
In [24]: hallway = [1, 0, 1, 0, 0, 1, 0, 1, 0, 0]
kernel = (.1, .8, .1)
belief = np.array([.1] * 10)
measurements = [1, 0, 1, 0, 0]

for m in measurements:
    update(hallway, belief, z=m, prob_correct=.75)
    predict(belief, 1, kernel)
bp.bar_plot(belief)
```



At this point we have correctly identified the likely cases, we either started at position 0 or 5, because we saw the following sequence of doors and walls 1,0,1,0,0. But now lets inject a bad measurement, and see what happens:

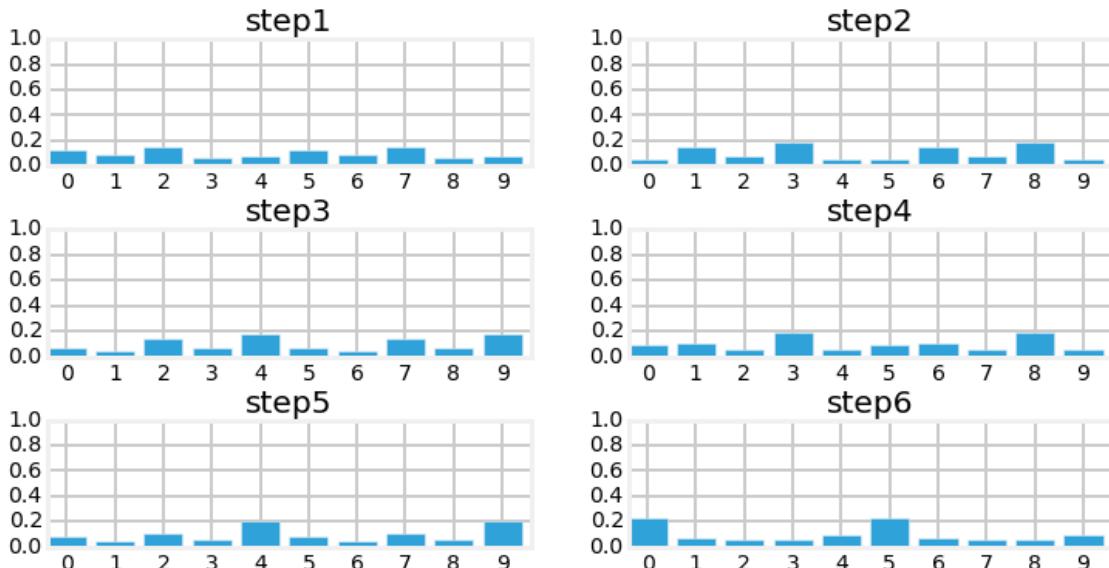
```
In [25]: update(hallway, belief, z=m, prob_correct=.75)
predict(belief, 1, kernel)
bp.bar_plot(belief)
```



That one bad measurement appears to have significantly eroded our knowledge. However, note that our highest probabilities are still at 0 and 5, which is correct. Now let's continue with a series of correct measurements

```
In [26]: with figsize(y=5.5):
    measurements = [0, 1, 0, 1, 0, 0]

    for i, m in enumerate(measurements):
        update(hallway, belief, z=m, prob_correct=.75)
        predict(belief, 1, kernel)
        plt.subplot(3, 2, i+1)
        bp.bar_plot(belief, title='step{}'.format(i+1))
```



As you can see we quickly filtered out the bad sensor reading and converged on the most likely positions for our dog.

## 2.10 Drawbacks and Limitations

Do not be mislead by the simplicity of the examples I chose. This is a robust and complete implementation of a histogram filter, and you may use the code in real world solutions. If you need a multimodal, discrete filter, this filter works.

With that said, while this filter is used in industry, it is not used often because it has several limitations. Getting around those limitations is the motivation behind the chapters in the rest of this book.

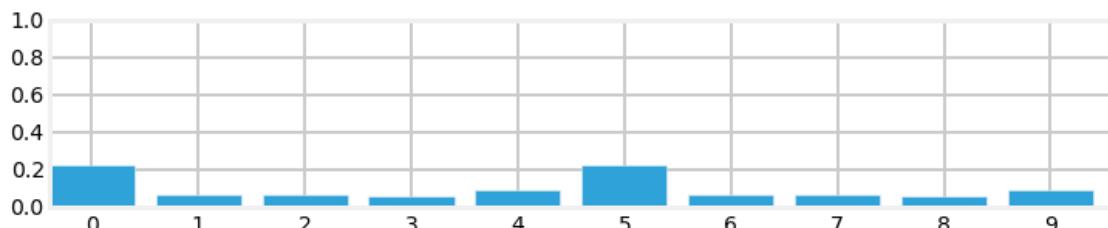
The first problem is scaling. Our dog tracking problem used only one variable, *pos*, to denote the dog's position. Most interesting problems will want to track several things in a large space. Realistically, at a minimum we would want to track our dogs (*x*, *y*) coordinate, and probably his velocity ( $\dot{x}$ ,  $\dot{y}$ ) as well. We have not covered the multidimensional case, but instead of a histogram we use a multidimensional grid to store the probabilities at each discrete location. Each `update()` and `predict()` step requires updating all values in the grid, so a simple four variable problem would require  $O(n^4)$  running time per time step. Realistic filters can have 10 or more variables to track, leading to exorbitant computation requirements.

The second problem is that the histogram is discrete, but we live in a continuous world. The histogram requires that you model the output of your filter as a set of discrete points. In our dog in the hallway example, we used 10 positions, which is obviously far too few positions for anything but a toy problem. For example, for a 100 meter hallway you would need 10,000 positions to model the hallway to 1cm accuracy. So each update and predict operation would entail performing calculations for 10,000 different probabilities. It gets exponentially worse as we add dimensions. If our dog was roaming in a 100x100 m<sup>2</sup> courtyard, we would need 100,000,000 bins ( $10,000^2$ ) to get 1cm accuracy.

A third problem is that the histogram is multimodal. This is not always a problem - an entire class of filters, the particle filters, are multimodal and are often used because of this property. But imagine if the GPS in your car reported to you that it is 40% sure that you are on D street, but 30% sure you are on Willow Avenue. I doubt that you would find that useful. Also, GPSs report their error - they might report that you are at (109.878W, 38.326N) with an error of 9 meters. There is no clear mathematical way to extract error information from a histogram. Heuristics suggest themselves to be sure, but there is no exact determination. You may or may not care about that while driving, but you surely do care if you are trying to send a rocket to Mars or track and hit an oncoming missile.

This difficulty is related to the fact that the filter often does not represent what is physically occurring in the world. Consider this distribution for our dog:

```
In [27]: belief = [0.2245871, 0.06288015, 0.06109133,
                 0.0581008, 0.09334062, 0.2245871,
                 0.06288015, 0.06109133, 0.0581008, 0.09334062]
bp.bar_plot(belief)
```



The largest probabilities are in position 0 and position 5. This does not fit our physical intuition at all. A dog cannot be in two places at once (my dog Simon certainly tries - his food bowl and my lap often have equal allure to him). We would have to use heuristics to decide how to interpret this distribution, and there is usually no satisfactory answer. This is not always a weakness - a considerable amount of literature has been written on Multi-Hypothesis Tracking (MHT). We cannot always distill our knowledge to one conclusion, and MHT uses various techniques to maintain multiple story lines at once, using backtracking schemes to go back in time to correct hypothesis once more information is known. This will be the subject of later chapters. In other cases we truly have a multimodal situation - we may be optically tracking pedestrians on the street and need to represent all of their positions.

In practice it is the exponential increase in computation time that leads to the discrete Bayes filter being the least frequently used of all filters in this book. Many problems are best formulated as discrete or multimodal, but we have other filter choices with better performance. With that said, if I had a small problem that this technique could handle I would choose to use it; it is trivial to implement, debug, and understand, all virtues.

## 2.11 Tracking and Control

So far we have been tracking an object which is moving independently. But consider this very similar problem. I am automating a warehouse and want to use robots to collect all of the items for a customer's order. Perhaps the easiest way to do this is to have the robots travel on a train track. I want to be able to send the robot a destination and have it go correctly to that point. But train tracks and robot motors are imperfect. Wheel slippage and imperfect motors means that the robot is unlikely to travel to exactly the position you command.

So, we add sensors. We can add some sort of device to help the robot determine it's position. Perhaps we mount magnets on the track every few feet, and use a Hall sensor to count how many magnets are passed. If we have counted 10 magnets then the robot should be at the 10th magnet. Of course it is possible to either miss a magnet or to count it twice, so we have to accommodate some degree of error. In any case, it should be clear that we can use the code in the previous section to track our robot since the magnet counting is very similar to doorway sensing.

But we are not done. A key lesson from the g-h filters chapter is to never throw information away. If you have information you should use it to improve your estimate. What information are we leaving out? Well, we know what our destination is, and we know what control inputs we are feeding to the wheels of the robot at each moment in time. For example, let's say that once a second we send a movement command to the robot - move left 1 unit, move right 1 unit, or stand still. This is obviously a simplification because I am not taking acceleration into account, but I am not trying to teach control theory. If I send the command 'move left 1 unit' I expect that in one second from now the robot will be 1 unit to the left of where it is now. But, wheels and motors are imperfect. We will assume that it never makes a mistake and goes right when told to go left, so the errors should be relatively small. Thus the robot might end up 0.9 units away, or maybe 1.2 units away.

Now the entire solution is clear. For the dog which was moving independently we assumed that he kept moving in whatever direction he was previously moving. That is a dubious assumption for my dog! Robots are far more predictable. Instead of making a dubious prediction based on assumption of behavior we will feed in the command that we sent to the robot! In other words, when we call `predict()` we will pass in the commanded movement that we gave the robot along with a kernel that describes the likelihood of that movement.

### 2.11.1 Simulating the Train Behavior

To fully implement this filter we need to simulate an imperfect train. When we command it to move it will sometimes make a small mistake, and its sensor will sometimes return the incorrect value.

In [28]:

```
class Train(object):
```

```
    def __init__(self, track_len, kernel=[1.], sensor_accuracy=.9):
        self.track_len = track_len
        self.pos = 0
        self.kernel = kernel
        self.sensor_accuracy = sensor_accuracy

    def move(self, distance=1):
        """ move in the specified direction
        with some small chance of error"""

        self.pos += distance
        # insert random movement error according to kernel
        r = random.random()
        s = 0
        offset = -(len(self.kernel) - 1) / 2
        for k in self.kernel:
            s += k
            if r <= s:
                break
            offset += 1
        self.pos = int((self.pos + offset) % self.track_len)
        return self.pos

    def sense(self):
        pos = self.pos
        # insert random sensor error
        r = random.random()
        if r > self.sensor_accuracy:
            if random.random() > 0.5:
                pos += 1
            else:
                pos -= 1
        return pos
```

With that we are ready to write the simulation. We will put it in a function so that we can run it with different assumptions. I will assume that the robot always starts at the beginning of the track. The track is implemented as being 10 units long, but think of it as a track of length, say 10,000, with the magnet pattern repeated every 10 units. 10 makes it easier to plot and inspect.

In [29]:

```
def simulate(iterations, kernel, sensor_accuracy,
             move_distance, do_print=True):
    track = np.array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])

    belief = np.array([0.01] * 10)
    belief[0] = .9
    normalize(belief)
```

```

robot = Train(len(track), kernel, sensor_accuracy)
for i in range(iterations):
    robot.move(distance=move_distance)
    m = robot.sense()
    if do_print:
        print('' time {}: pos {}, sensed {}, '''''
              '''at position {}'''.format(
                  i, robot.pos, m, track[robot.pos]))

    update(track, belief, m, sensor_accuracy)
    index = np.argmax(belief)
    if do_print:
        print('' predicted position is {}'''''
              ''' with confidence {:.4f}%:''''.format(
                  index, belief[index]*100))
    if i < iterations - 1:
        predict(belief, move_distance, kernel)

bp.bar_plot(belief)
if do_print:
    print()
    print('final position is', robot.pos)
    iindex = np.argmax(belief)
    print('predicted position is {} with '''''
          '''confidence {:.4f}%:''''.format(
              index, belief[index]*100))

```

Read the code and make sure you understand it. Now let's do a run with no sensor or movement error. If the code is correct it should be able to locate the robot with no error. The output is a bit tedious to read, but if you are at all unsure of how the update/predict cycle works make sure you read through it carefully to solidify your understanding.

In [30]: `import random`

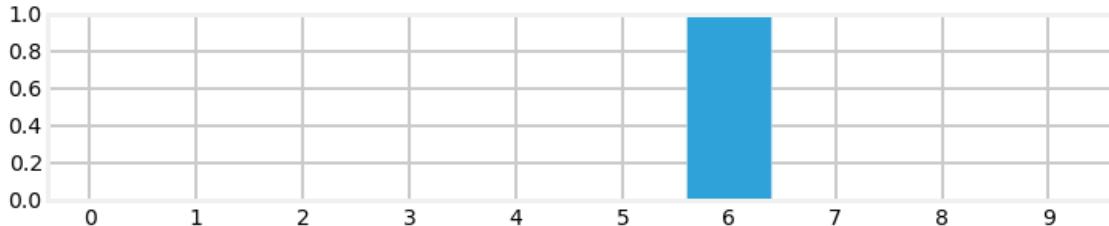
```

random.seed(3)
np.set_printoptions(precision=2, suppress=True, linewidth=60)
simulate(4, kernel=[1.], sensor_accuracy=.999,
         move_distance=4, do_print=True)

time 0: pos 4, sensed 4, at position 4
    predicted position is 4 with confidence 91.0665%:
time 1: pos 8, sensed 8, at position 8
    predicted position is 8 with confidence 99.9902%:
time 2: pos 2, sensed 2, at position 2
    predicted position is 2 with confidence 100.0000%:
time 3: pos 6, sensed 6, at position 6
    predicted position is 6 with confidence 100.0000%:

final position is 6
predicted position is 6 with confidence 100.0000%:

```

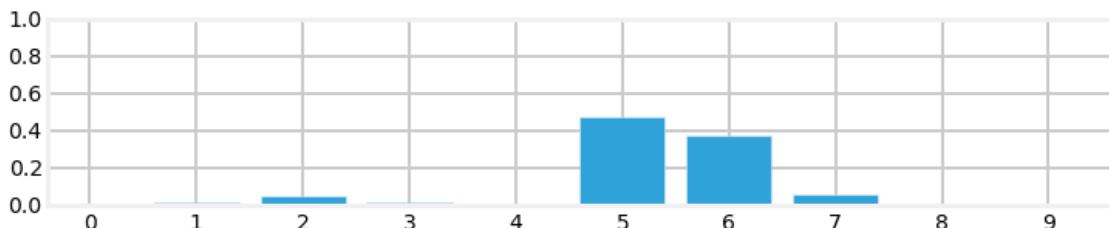


We can see that the code was able to perfectly track the robot so we should feel reasonably confident that the code is working. Now let's see how it fairs with some errors.

```
In [31]: random.seed(3)
        simulate(4, kernel=[.1, .8, .1], sensor_accuracy=.9,
                  move_distance=4, do_print=True)

time 0: pos 4, sensed 4, at position 4
          predicted position is 0 with confidence 84.1121%:
time 1: pos 8, sensed 8, at position 8
          predicted position is 4 with confidence 43.4416%:
time 2: pos 2, sensed 2, at position 2
          predicted position is 2 with confidence 81.0925%:
time 3: pos 5, sensed 5, at position 5
          predicted position is 5 with confidence 47.8278%:

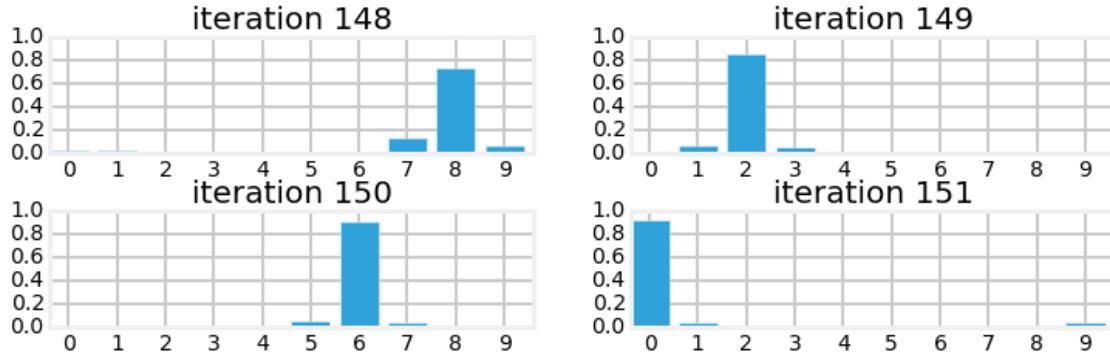
final position is 5
predicted position is 5 with confidence 47.8278%:
```



Here we see that there was a sense error at time 1, but we are still quite confident in our position.

Now lets run a very long simulation and see how the filter responds to errors.

```
In [32]: with figsize(y=5):
        for i in range(4):
            random.seed(3)
            plt.subplot(321+i)
            simulate(148+i, kernel=[.1, .8, .1],
                      sensor_accuracy=.8,
                      move_distance=4, do_print=False)
            plt.title('iteration {}'.format(148+i))
```



We can see that there was a problem on iteration 149 as the confidence degrades. But within a few iterations the filter is able to correct itself and regain confidence in the estimated position.

## 2.12 Bayes Theorem

We developed the math in this chapter merely by reasoning about the information we have at each moment. In the process we discovered Bayes Theorem. We will go into the specifics of the math of Bayes theorem later in the book. For now we will take a more intuitive approach. Recall from the preface that Bayes theorem tells us how to compute the probability of an event given previous information. That is exactly what we have been doing in this chapter. With luck our code should match the Bayes Theorem equation!

We implemented the `update()` function with this probability calculation:

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{normalization}}$$

To review, the prior is the probability of something happening before we include the probability of the measurement (the likelihood) and the posterior is the probability we compute after incorporating the information from the measurement.

Bayes theorem is

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

If you are not familiar with this notation, let's review.  $P(A)$  means the probability of event  $A$ . If  $A$  is the event of a fair coin landing heads, then  $P(A) = 0.5$ .

$P(A|B)$  is called a conditional probability. That is, it represents the probability of  $A$  happening if  $B$  happened. For example, it is more likely to rain today if it also rained yesterday because rain systems tend to last more than one day. We'd write the probability of it raining today given that it rained yesterday as  $P(\text{rain}_{\text{today}}|\text{rain}_{\text{yesterday}})$ .

In the Bayes theorem equation above  $B$  is the evidence,  $P(A)$  is the prior,  $P(B|A)$  is the likelihood, and  $P(A|B)$  is the posterior. By substituting the mathematical terms with the corresponding words you can see that Bayes theorem matches out update equation. Let's rewrite the equation in terms of our problem. We will use  $x_i$  for the position at  $i$ , and  $Z$  for the measurement. Hence, we want to know  $P(x_i|Z)$ , that is, the probability of the dog being at  $x_i$  given the measurement  $Z$ .

So, let's plug that into the equation and solve it.

$$P(x_i|Z) = \frac{P(Z|x_i)P(x_i)}{P(Z)}$$

That looks ugly, but it is actually quite simple. Let's figure out what each term on the right means. First is  $P(Z|x_i)$ . This is the probability for the measurement at every cell  $x_i$ .  $P(x_i)$  is the prior - our belief before incorporating the measurements. We multiply those together. This is just the unnormalized multiplication in the `update()` function, where `belief` is our prior  $P(x_i)$ .

```
for i, val in enumerate(map_):
    if val == z:
        belief[i] *= correct_scale
    else:
        belief[i] *= 1.
```

I added the `else` here, which has no mathematical effect, to point out that every element in  $x$  (called `belief` in the code) is multiplied by a probability. You may object that I am multiplying by a scale factor, which I am, but this scale factor is derived from the probability of the measurement being correct vs the probability being incorrect.

The last term to consider is the denominator  $P(Z)$ . This is the probability of getting the measurement  $Z$  without taking the location into account. We compute that by taking the sum of  $x$ , or `sum(belief)` in the code. That is how we compute the normalization! So, the `update()` function is doing nothing more than computing Bayes theorem.

The literature often gives you these equations in the form of integrals. After all, an integral is just a sum over a continuous function. So, you might see Bayes' theorem written as

$$P(A|B) = \frac{P(B|A)P(A)}{\int P(B|A)P(B)dy}.$$

In practice the denominator can be fiendishly difficult to solve analytically (a recent opinion piece for the Royal Statistical Society called it a “dog’s breakfast” [8], and filtering textbooks are filled with integral laden equations which you cannot be expected to solve. We will learn more techniques to handle this in the **Particle Filters** chapter. Until then, recognize that in practice it is just a normalization term over which we can sum. What I’m trying to say is that when you are faced with a page of integrals, just think of them of sums, and relate them back to this chapter, and often the difficulties will fade. Ask yourself “why are we summing these values”, and “why am I dividing by this term”. Surprisingly often the answer is readily apparent.

## 2.13 Total Probability Theorem

We know now the formal mathematics behind the `update()` function; what about the `predict()` function? `predict()` implements the total probability theorem. Let's recall what `predict()` computed. It computed the probability of being at any given position given the probability of all the possible movement events. Let's express that as an equation. The probability of being at any position  $i$  at time  $t$  can be written as  $P(X_i^t)$ . We computed that as the sum of the prior at time  $t-1$   $P(X_j^{t-1})$  multiplied by the probability of moving from cell  $x_j$  to  $x_i$ . That is

$$P(X_i^t) = \sum_j P(X_j^{t-1})P(x_i|x_j)$$

That equation is called the total probability theorem. Quoting from Wikipedia [6] “It expresses the total probability of an outcome which can be realized via several distinct events”. I could have given you that

equation and implemented `predict()`, but your chances of understanding why the equation works would be slim. As a reminder, here is the code that computes this equation

```
for i in range(N):
    for k in range (kN):
        index = (i + (width-k) - offset) % N
        result[i] += prob_dist[index] * kernel[k]
```

## 2.14 Summary

The code is very short, but the result is huge! We have implemented a form of a Bayesian filter. We have learned how to start with no information and derive information from noisy sensors. Even though the sensors in this chapter are very noisy (most sensors are more than 80% accurate, for example) we quickly converge on the most likely position for our dog. We have learned how the predict step always degrades our knowledge, but the addition of another measurement, even when it might have noise in it, improves our knowledge, allowing us to converge on the most likely result.

If you followed the math carefully you will realize that all of this math is exact. The bar charts that we are displaying are not an estimate or guess - they are mathematically exact results that exactly represent our knowledge. The knowledge is probabilistic, to be sure, but it is exact, and correct.

Furthermore, through basic reasoning we were able to discover two extremely important theorems: Bayes theorem and the total probability theorem. I hope you spent time on those sections as in almost any other source they will express filtering algorithms in terms of these two theorems. It will be your job to understand what these equations mean and how to turn them into code.

This book is mostly about the Kalman filter. In the g-h filter chapter I told you that the Kalman filter is a type of g-h filter. It is also a type of Bayesian filter. It also uses Bayes theorem and the total probability theorem to filter data, although with a different set of assumptions and conditions than used in this chapter.

The discrete Bayes filter allows us to filter sensors and track an object, but we are a long way from tracking an airplane or a car. This code only handles the 1 dimensional case, whereas cars and planes operate in 2 or 3 dimensions. Also, our position vector is multimodal. It expresses multiple beliefs at once. Imagine if your GPS told you “it’s 20% likely that you are here, but 10% likely that you are on this other road, and 5% likely that you are at one of 14 other locations”. That would not be very useful information. Also, the data is discrete. We split an area into 10 (or whatever) different locations, whereas in most real world applications we want to work with continuous data. We want to be able to represent moving 1 km, 1 meter, 1 mm, or any arbitrary amount, such as 2.347 cm.

Finally, the bar charts may strike you as being a bit less certain than we would want. A 25% certainty may not give you a lot of confidence in the answer. Of course, what is important here is the ratio of this probability to the other probabilities in your vector. If the next largest bar is 23% then we are not very knowledgeable about our position, whereas if the next largest is 3% we are in fact quite certain. But this is not clear or intuitive. However, there is an extremely important insight that Kalman filters implement that will significantly improve our accuracy from the same data.

**If you can understand this chapter you will be able to understand and implement Kalman filters.** I cannot stress this enough. If anything is murky, go back and reread this chapter and play with the code. The rest of this book will build on the algorithms that we use here. If you don’t intuitively understand why this filter works, and can’t at least work through the math, you will have little success with the rest of the material. However, if you grasp the fundamental insight - multiplying probabilities when we measure, and shifting probabilities when we update leads to a converging solution - then you understand everything important you need for the Kalman filter.

## 2.15 References

- [1] D. Fox, W. Burgard, and S. Thrun. “Monte carlo localization: Efficient position estimation for mobile robots.” In Journal of Artificial Intelligence Research, 1999.

<http://www.cs.cmu.edu/afs/cs/project/jair/pub/volume11/fox99a-html/jair-localize.html>

- [2] Dieter Fox, et. al. “Bayesian Filters for Location Estimation”. In IEEE Pervasive Computing, September 2003.

<http://swarmlab.unimaas.nl/wp-content/uploads/2012/07/fox2003bayesian.pdf>

- [3] Sebastian Thrun. “Artificial Intelligence for Robotics”.

<https://www.udacity.com/course/cs373>

- [4] Khan Academy. “Introduction to the Convolution”

<https://www.khanacademy.org/math/differential-equations/laplace-transform/convolution-integral/v/introduction-to-the-convolution>

- [5] Wikipedia. “Convolution”

<http://en.wikipedia.org/wiki/Convolution>

- [6] Wikipedia. “Law of total probability”

[http://en.wikipedia.org/wiki/Law\\_of\\_total\\_probability](http://en.wikipedia.org/wiki/Law_of_total_probability)

- [7] Wikipedia. “Time Evolution”

[https://en.wikipedia.org/wiki/Time\\_evolution](https://en.wikipedia.org/wiki/Time_evolution)

- [8] We need to rethink how we teach statistics from the ground up

<http://www.statslife.org.uk/opinion/2405-we-need-to-rethink-how-we-teach-statistics-from-the-ground-up>



# Chapter 3

## Gaussian Probabilities

### 3.1 Introduction

The last chapter ended by discussing some of the drawbacks of the Discrete Bayesian filter. For many tracking and filtering problems our desire is to have a filter that is unimodal and continuous. That is, we want to model our system using floating point math (continuous) and to have only one belief represented (unimodal). For example, we want to say an aircraft is at (12.34, -95.54, 2389.5) where that is latitude, longitude, and altitude. We do not want our filter to tell us “it might be at (1.65, -78.01, 2100.45) or it might be at (34.36, -98.23, 2543.79).” That doesn’t match our physical intuition of how the world works, and as we discussed, it can be prohibitively expensive to compute the multimodal case. And, of course, multiple position estimates makes navigating impossible.

So we desire a unimodal, continuous way to represent probabilities that models how the real world works, and that is computationally efficient to calculate. As you might guess from the chapter name, Gaussian distributions provide all of these features.

### 3.2 Mean, Variance, and Standard Deviations

#### 3.2.1 Random Variables

Each time you roll a die the outcome will be between 1 and 6. If we rolled a fair die a million times we’d expect to get 1 1/6 of the time. Thus we say the probability, or odds of the outcome 1 is 1/6. Likewise, if I asked you the chance of 1 being the result of the next roll you’d reply 1/6.

This combination of values and associated probabilities is called a random variable. Random does not mean the process is nondeterministic, only that we lack information. The result of a die toss is deterministic, but we lack enough information to compute the result. We don’t know what will happen, except probabilistically.

While we are defining things, the range of values is called the sample space. For a die the sample space is {1, 2, 3, 4, 5, 6}. For a coin the sample space is {H, T}. Space is a mathematical term which means a set with structure. The sample space for the die is a subset of the natural numbers in the range of 1 to 6.

Another example of a random variable is the heights of students in a university. Here the sample space is a range of values in the real numbers between two limits defined by biology.

Random variables such as coin tosses and die rolls are discrete random variables. This means their sample space is represented by either a finite number of values or a countably infinite number of values such as the natural numbers. Heights of humans are called continuous random variables since they can take on any real value between two limits.

Do not confuse the measurement of the random variable with the actual value. If we can only measure the height of a person to 0.1 meters we would only record values from 0.1, 0.2, 0.3...2.7, yielding 27 discrete choices. Nonetheless a person's height can vary between any arbitrary real value between those ranges, and so height is a continuous random variable.

In statistics capital letters are used for random variables, usually from the latter half of the alphabet. So, we might say that  $X$  is the random variable representing the die toss, or  $Y$  are the heights of the students in the freshmen poetry class. In later chapters we will follow the convention of using lower case for vectors, and upper case for matrices. Unfortunately these conventions clash, and you will have to determine which an author is using from context.

### 3.3 Probability Distribution

The probability distribution gives the probability for the random variable to take any value in a sample space. For example, for a fair six sided die we might say:

Value	Probability
1	1/6
2	1/6
3	1/6
4	1/6
5	1/6
6	1/6

Some sources call this the probability function. Using ordinary function notation, we would write:

$$P(X=4) = f(4) = \frac{1}{6}$$

This states that the probability of the die landing on 4 is  $\frac{1}{6}$ .  $P(X=x_k)$  is notation for “the probability of  $X$  being  $x_k$ . Some texts use  $Pr$  or  $Prob$  instead of  $P$ .

Another example is a fair coin. It has the sample space  $\{H, T\}$ . The coin is fair, so the probability for heads ( $H$ ) is 50%, and the probability for tails ( $T$ ) is 50%. We write this as

$$P(X=H) = 0.5P(X=T) = 0.5$$

Sample spaces are not unique. One sample space for a die is  $\{1, 2, 3, 4, 5, 6\}$ . Another valid sample space would be  $\{\text{even}, \text{odd}\}$ . Another might be  $\{\text{dots in all corners, not dots in all corners}\}$ . A sample space is valid so long as it covers all possibilities, and any single event is described by only one element.  $\{\text{even}, 1, 3, 4, 5\}$  is not a valid sample space for a die since a value of 4 is matched both by ‘even’ and ‘4’.

The probabilities for all values of a discrete random value is known as the discrete probability distribution and the probabilities for all values of a continuous random value is known as the continuous probability distribution.

To be a probability distribution the probability of each value  $x_i$  must be  $x_i \geq 0$ , since no probability can be less than zero. Secondly, the sum of the probabilities for all values must equal one. This should be intuitively clear for a coin toss: if the odds of getting heads is 70%, then the odds of getting tails must be 30%. We formalize this requirement as

$$\sum_u P(X=u) = 1$$

for discrete distributions, and as

$$\int P(X=u) = 1$$

for continuous distributions.

### 3.3.1 The Mean, Median, and Mode of a Random Variable

Given a set of data we often want to know a representative or average value for that set. There are many measures for this, and the concept is called a measure of central tendency. For example we will want to know the average height of the students. We all know how to find the average, but let me belabor the point so I can introduce more formal notation and terminology. Another word for average is the mean. We compute the mean by summing the values and dividing by the number of values. If the heights of the students in meters is

$$X = \{1.8, 2.0, 1.7, 1.9, 1.6\}$$

we compute the mean as

$$\mu = \frac{1.8 + 2.0 + 1.7 + 1.9 + 1.6}{5} = 1.8$$

It is traditional to use the symbol  $\mu$  (mu) to denote the mean.

We can formalize this computation with the equation

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

NumPy provides `numpy.mean()` for computing the mean.

```
In [2]: import numpy as np
        x = [1.85, 2.0, 1.7, 1.9, 1.6]
        print(np.mean(x))
```

1.81

The mode of a set of numbers is the number that occurs most often. If only one number occurs most often we say it is a unimodal set, and if two or more numbers occur the most with equal frequency than the set is multimodal. For example the set  $\{1, 2, 2, 2, 3, 4, 4, 4\}$  has modes 2 and 4, which is multimodal, and the set  $\{5, 7, 7, 13\}$  has the mode 7, and so is unimodal. We will not be computing the mode in this manner in this book, but we do use the concepts of unimodal and multimodal in a more general sense. For example, in the **Discrete Bayes** chapter we talked about our belief in the dog's position as a multimodal distribution because we assigned different probabilities to different positions.

Finally, the median of a set of numbers is the middle point of the set so that half the values are below the median and half are above the median. Here, above and below is in relation to the set being sorted. If the set contains an even number of values then the two middle numbers are averaged together.

Numpy provides `numpy.median()` to compute the median. As you can see the median of  $\{1.85, 2.0, 1.7, 1.9, 1.6\}$  is 1.85, because 1.85 is the third element of this set after being sorted.

In [3]: `print(np.median(x))`

1.85

## 3.4 Expected Value of a Random Variable

The expected value of a random variable is the average value it would have if we took an infinite number of samples of it and then averaged those samples together. Let's say we have  $x = [1, 3, 5]$  and each value is equally probable. What would we expect  $x$  to have, on average?

It would be the average of 1, 3, and 5, of course, which is 3. That should make sense; we would expect equal numbers of 1, 3, and 5 to occur, so  $(1 + 3 + 5)/3 = 3$  is clearly the average of that infinite series of samples. In other words, here the expected value is just the mean of the sample space.

Now suppose that each value has a different probability of happening. Say 1 has an 80% chance of occurring, 3 has an 15% chance, and 5 has only a 5% chance. In this case we compute the expected value by multiplying each value of  $x$  by the percent chance of it occurring, and summing the result. So for this case we could compute

$$\mathbb{E}[X] = (1)(0.8) + (3)(0.15) + (5)(0.05) = 1.5$$

Here I have introduced the notation  $\mathbb{E}[X]$  for the expected value of  $x$ . Some texts use  $E(x)$ . The value 1.5 for  $x$  makes intuitive sense because x is far more like to be 1 than 3 or 5, and 3 is more likely than 5 as well.

We can formalize this by letting  $x_i$  be the  $i^{th}$  value of  $X$ , and  $p_i$  be the probability of its occurrence. This gives us

$$\mathbb{E}[X] = \sum_{i=1}^n p_i x_i$$

A trivial bit of algebra shows that if the probabilities are all equal, the expected value is the same as the mean:

$$\mathbb{E}[X] = \sum_{i=1}^n p_i x_i = \sum_{i=1}^n \frac{1}{n} x_i = \mu_x$$

If  $x$  is continuous we substitute the sum for an integral, like so

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} x f(x)$$

where  $f(x)$  is the probability distribution function of  $x$ . We won't be using this equation yet, but we will be using it in the next chapter.

### 3.4.1 Variance of a Random Variable

The computation above tells us the average height of the students, but it doesn't tell us everything we might want to know. For example, suppose we have three classes of students, which we label  $X$ ,  $Y$ , and  $Z$ , with these heights:

In [4]: `X = [1.8, 2.0, 1.7, 1.9, 1.6]`  
`Y = [2.2, 1.5, 2.3, 1.7, 1.3]`  
`Z = [1.8, 1.8, 1.8, 1.8, 1.8]`

Using NumPy we see that the mean height of each class is the same.

```
In [5]: print(np.mean(X))
      print(np.mean(Y))
      print(np.mean(Z))
```

```
1.8
1.8
1.8
```

The mean of each class is 1.8 meters, but notice that there is a much greater amount of variation in the heights in the second class than in the first class, and that there is no variation at all in the third class.

So the mean tells us something about the data, but it does not tell the whole story. We want to be able to specify how much variation there is between the heights of the students. You can imagine a number of reasons for this. Perhaps a school district needs to order 5,000 desks, and they want to be sure they buy sizes that accommodate the range of heights of the students.

Statistics has formalized this concept of measuring variation into the notion of standard deviation and variance. The equation for computing the variance is

$$VAR(X) = E[(X - \mu)^2]$$

Ignoring the squared terms for a moment, you can see that the variance is the expected value for how much the sample space ( $X$ ) varies from the mean (squared, of course). We have the formula for the expected value  $E[X] = \sum_{i=1}^n p_i x_i$ , and we will assume that any height is equally probable, so we can substitute that into the equation above to get

$$VAR(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

Let's compute the variance of the three classes to see what values we get and to become familiar with this concept.

The mean of  $X$  is 1.8 ( $\mu_x = 1.8$ ) so we compute

$$\begin{aligned} VAR(X) &= \frac{(1.8 - 1.8)^2 + (2 - 1.8)^2 + (1.7 - 1.8)^2 + (1.9 - 1.8)^2 + (1.6 - 1.8)^2}{5} \\ &= \frac{0 + 0.04 + 0.01 + 0.01 + 0.04}{5} \\ VAR(X) &= 0.02 \text{ m}^2 \end{aligned}$$

NumPy provides the function `var()` to compute the variance:

```
In [6]: print(np.var(X), "meters squared")
```

```
0.02 meters squared
```

This is perhaps a bit hard to interpret. Heights are in meters, yet the variance is meters squared. Thus we have a more commonly used measure, the standard deviation, which is defined as the square root of the variance:

$$\sigma = \sqrt{VAR(X)} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

It is typical to use  $\sigma$  for the standard deviation and  $\sigma^2$  for the variance. In most of this book I will be using  $\sigma^2$  instead of  $VAR(X)$  for the variance; they symbolize the same thing.

For the first class we compute the standard deviation with

$$\begin{aligned}\sigma_x &= \sqrt{\frac{(1.8 - 1.8)^2 + (2 - 1.8)^2 + (1.7 - 1.8)^2 + (1.9 - 1.8)^2 + (1.6 - 1.8)^2}{5}} \\ &= \sqrt{\frac{0 + 0.04 + 0.01 + 0.01 + 0.04}{5}} \\ \sigma_x &= 0.1414\end{aligned}$$

We can verify this computation with the NumPy method `numpy.std()` which computes the standard deviation. ‘std’ is a common abbreviation for standard deviation.

```
In [7]: print('std {:.4f}'.format(np.std(X)))
print('var {:.4f}'.format(np.std(X)**2))

std 0.1414
var 0.0200
```

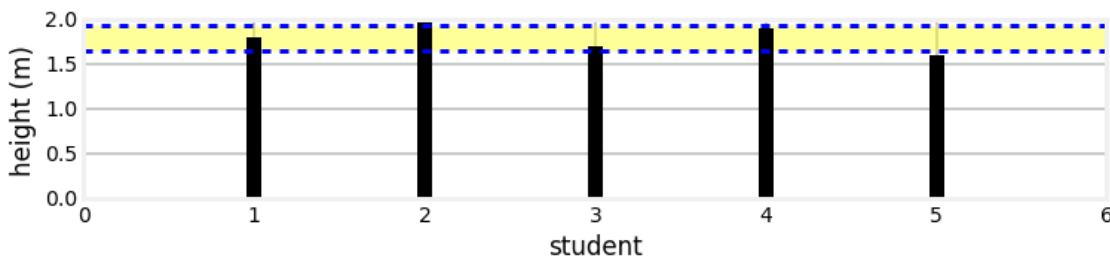
And, of course,  $0.1414^2 = 0.02$ , which agrees with our earlier computation of the variance.

What does the standard deviation signify? It tells us how much the heights vary amongst themselves. “How much” is not a mathematical term. We will be able to define it much more precisely once we introduce the concept of a Gaussian in the next section. For now I’ll say that for many things 68% of all values lie within one standard deviation of the mean. In other words we can conclude that for a random class 68% of the students will have heights between 1.66 (1.8-0.1414) meters and 1.94 (1.8+0.1414) meters.

We can view this in a plot:

```
In [8]: from book_format import set_figsize, figsize
from gaussian_internal import plot_height_std
import matplotlib.pyplot as plt

with figsize(y=2):
    plot_height_std(X)
```



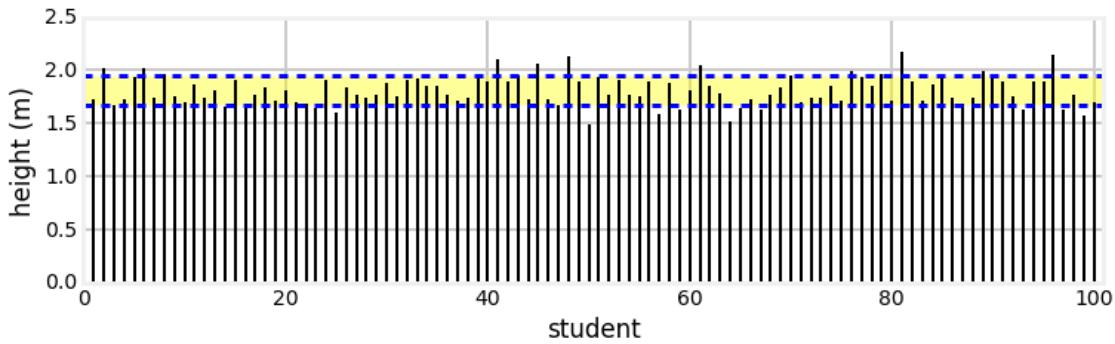
For only 5 students we obviously will not get exactly 68% within one standard deviation. We do see that 3 out of 5 students are within  $1\sigma$ , or 60%, which is as close as you can get to 68% with only 5 samples. I haven't yet introduced enough math or Python for you to fully understand the next bit of code, but let's look at the results for a class with 100 students.

We write one standard deviation as  $1\sigma$ , which is pronounced “one standard deviation”, not “one sigma”. Two standard deviations is  $2\sigma$ , and so on.

```
In [9]: from numpy.random import randn
data = [1.8 + .1414*randn() for i in range(100)]

with figsize(y=3.):
    plot_height_std(data, lw=2)

print('mean = {:.3f}'.format(np.mean(data)))
print('std  = {:.3f}'.format(np.std(data)))
```



```
mean = 1.809
std  = 0.137
```

We can see by eye that roughly 68% of the heights lie within  $1\sigma$  of the mean 1.8.

We'll discuss this in greater depth soon. For now let's compute the standard deviation for

$$Y = [2.2, 1.5, 2.3, 1.7, 1.3]$$

The mean of  $Y$  is  $\mu = 1.8$  m, so

$$\begin{aligned}\sigma_y &= \sqrt{\frac{(2.2 - 1.8)^2 + (1.5 - 1.8)^2 + (2.3 - 1.8)^2 + (1.7 - 1.8)^2 + (1.3 - 1.8)^2}{5}} \\ &= \sqrt{0.152} = 0.39 \text{ m}\end{aligned}$$

We will verify that with NumPy with

```
In [10]: print('std of Y is {:.4f} m'.format(np.std(Y)))
```

```
std of Y is 0.3899 m
```

This corresponds with what we would expect. There is more variation in the heights for  $Y$ , and the standard deviation is larger.

Finally, let's compute the standard deviation for  $Z$ . There is no variation in the values, so we would expect the standard deviation to be zero. We show this to be true with

$$\begin{aligned}\sigma_z &= \sqrt{\frac{(1.8 - 1.8)^2 + (1.8 - 1.8)^2 + (1.8 - 1.8)^2 + (1.8 - 1.8)^2 + (1.8 - 1.8)^2}{5}} \\ &= \sqrt{\frac{0 + 0 + 0 + 0 + 0}{5}} \\ \sigma_z &= 0.0 \text{ m}\end{aligned}$$

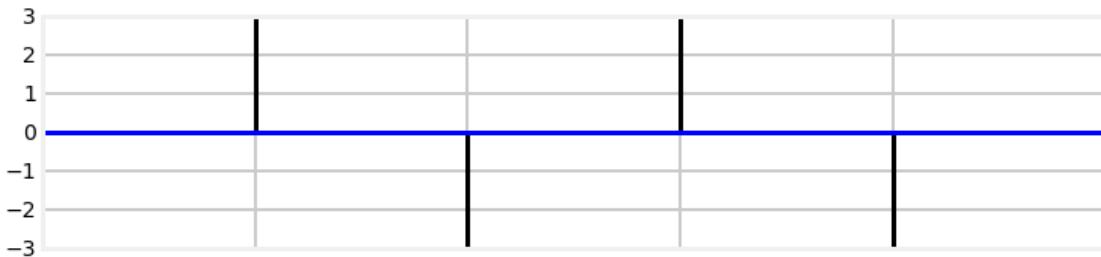
In [11]: `print(np.std(Z))`

0.0

### 3.4.2 Why the Square of the Differences

Why are we taking the square of the differences for the variance? I could go into a lot of math, but let's look at this in a simple way. Here is a chart of the values of  $X$  plotted against the mean for  $X = [3, -3, 3, -3]$

```
In [12]: with figsize(y=2.5):
    X = [3, -3, 3, -3]
    mean = np.average(X)
    for i in range(len(X)):
        plt.plot([i, i], [mean, X[i]], color='k')
    plt.axhline(mean)
    plt.xlim(-1, len(X))
    plt.tick_params(axis='x', labelbottom='off')
```



If we didn't take the square of the differences the signs would cancel everything out:

$$\frac{(3 - 0) + (-3 - 0) + (3 - 0) + (-3 - 0)}{4} = 0$$

This is clearly incorrect, as there is more than 0 variance in the data.

Maybe we can use the absolute value? We can see by inspection that the result is  $12/4 = 3$  which is certainly correct - each value varies by 3 from the mean. But what if have  $Y = [6, -2, -3, 1]$ ? In this case we get  $12/4 = 3$ .  $Y$  is clearly more spread out than  $X$ , but the computation yields the same variance. If we use the correct formula we get a variance of 3.5 for  $Y$ , which reflects its larger variation.

This is not a proof of correctness. Indeed, Carl Friedrich Gauss, the inventor of the technique, recognized that it is somewhat arbitrary. If there are outliers then squaring the difference gives disproportionate weight to that term. For example, let's see what happens if we have  $X = [1, -1, 1, -2, 3, 2, 100]$ .

```
In [13]: X = [1, -1, 1, -2, 3, 2, 100]
print('Variance of X = {:.2f}'.format(np.var(X)))
```

Variance of X = 1210.69

Is this “correct”? You tell me. Without the outlier of 100 we get  $\sigma^2 = 2.89$ , which accurately reflects how  $X$  is varying absent the outlier. The one outlier swamps the computation. I will not continue down this path; if you are interested you might want to look at the work that James Berger has done on this problem, in a field called Bayesian robustness, or the excellent publications on robust statistics by Peter J. Huber [3].

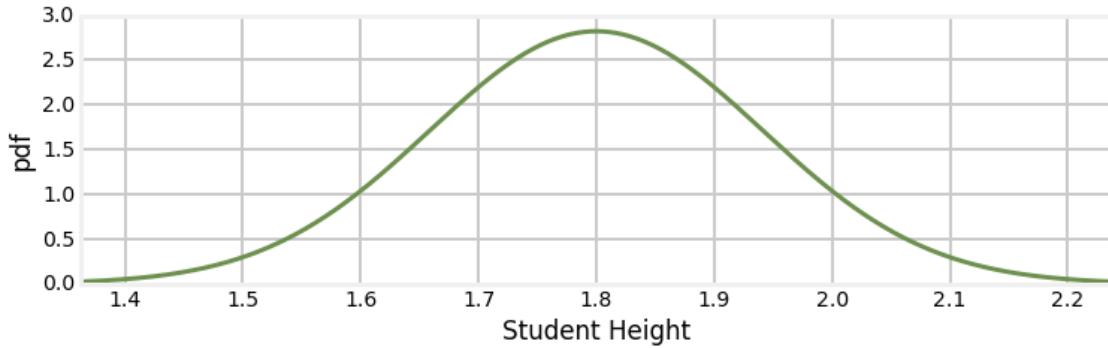
## 3.5 Gaussians

We are now ready to learn about Gaussians. Let's remind ourselves of the motivation for this chapter.

We desire a unimodal, continuous way to represent probabilities that models how the real world works, and that is computationally efficient to calculate.

Let's look at a graph of a Gaussian distribution to get a sense of what we are talking about.

```
In [14]: from filterpy.stats import plot_gaussian_pdf
plot_gaussian_pdf(mean=1.8, variance=0.1414**2,
                    xlabel='Student Height', ylabel='pdf');
```



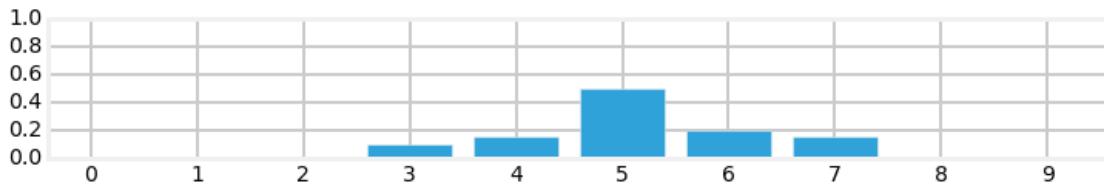
I explain how to plot Gaussians, and much more, in the Notebook Computing\_and\_Plotting\_PDFs in the Supporting\_Notebooks folder. You can read it online [here](#)[1]

This may be recognizable to you as a ‘bell curve’. This curve is ubiquitous because under real world conditions many observations are distributed in such a manner. In fact, this is the curve for the student heights given earlier. I will not use the term ‘bell curve’ to refer to a Gaussian because many probability distributions have a similar bell curve shape. Non-mathematical sources might not be so precise, so be judicious in what you conclude when you see the term used without definition.

This curve is not unique to heights - a vast amount of natural phenomena exhibits this sort of distribution, including the sensors that we use in filtering problems. As we will see, it also has all the attributes that we are looking for - it represents a unimodal belief or value as a probability, it is continuous, and it is computationally efficient. We will soon discover that it also has other desirable qualities which we may not realize we desire.

To further motivate you, recall the shapes of the probability distributions in the [Discrete Bayes](#) chapter. They were not perfect Gaussian curves, but they were similar, as in the plot below. We will be using Gaussians to replace the discrete probabilities used in that chapter!

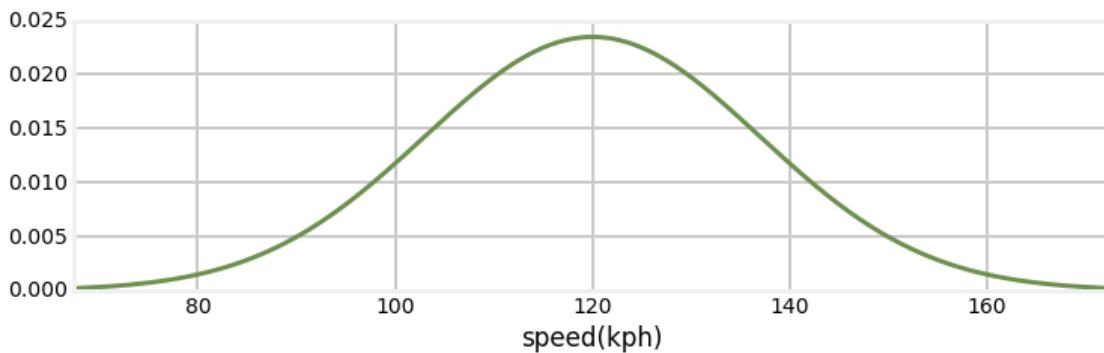
```
In [15]: import book_plots
belief = [0., 0., 0., 0.1, 0.15, 0.5, 0.2, .15, 0, 0]
with figsize(y=1.5):
    book_plots.bar_plot(belief)
```



## 3.6 Nomenclature

A bit of nomenclature before we continue - this chart depicts the [probability density](#) of a random variable having any value between  $(-\infty, \infty)$ . What does that mean? Imagine we take an infinite number of infinitely precise measurements of the speed of automobiles on a section of highway. We could then plot the results by showing the relative number of cars going at any given speed. If the average was 120 kph, it might look like this:

```
In [16]: with figsize(y=3.):
    plot_gaussian_pdf(mean=120, variance=17**2, xlabel='speed(kph)')
```



The y-axis depicts the probability density - the relative amount of cars that are going the speed at the corresponding x-axis.

You may object that human heights or automobile speeds cannot be less than zero, let alone  $-\infty$  or  $+\infty$ . This is true, but this is a common limitation of mathematical modeling. “The map is not the territory” is a common expression, and it is true for Bayesian filtering and statistics. The Gaussian distribution above models the distribution of the measured automobile speeds, but being a model it is necessarily imperfect. The difference between model and reality will come up again and again in these filters. Gaussians are used in many branches of mathematics, not because they perfectly model reality, but because they are easier to use than any other relatively accurate choice. However, even in this book Gaussians will fail to model reality, forcing us to use computationally expensive alternative.

You will see these distributions called Gaussian distributions or normal distributions. Gaussian and normal both mean the same thing in this context, and are used interchangeably. I will use both throughout this book as different sources will use either term, and I want you to be used to seeing both. Finally, as in this paragraph, it is typical to shorten the name and talk about a Gaussian or normal - these are both typical shortcut names for the Gaussian distribution.

## 3.7 Gaussian Distributions

So let us explore how Gaussians work. A Gaussian is a continuous probability distribution that is completely described with two parameters, the mean ( $\mu$ ) and the variance ( $\sigma^2$ ). It is defined as:

$$f(x, \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right]$$

$\exp[x]$  is notation for  $e^x$ ; I avoid using superscripts in print so that the fonts are larger and more readable. Don’t be dissuaded by the equation if you haven’t seen it before; you will not need to memorize or manipulate it. The computation of this function is stored in `stats.py` with the function `gaussian(x, mean, var)`.

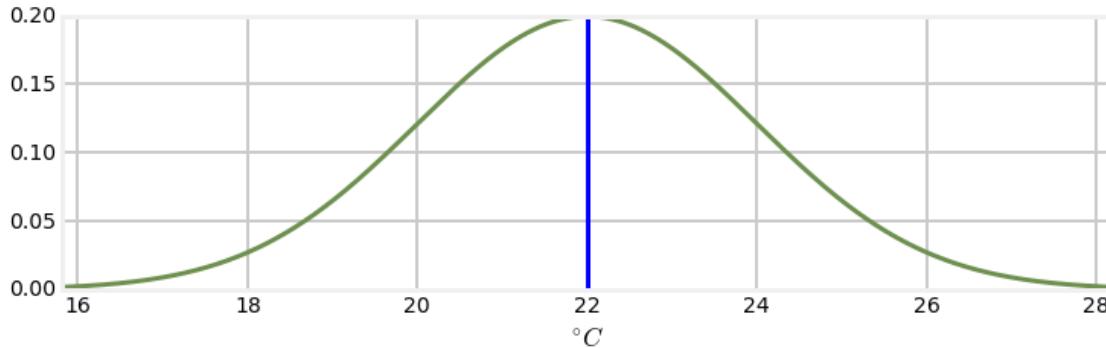
**Optional:** Let’s remind ourselves how to look at a function stored in a file by using the `%load` magic. If you type `%load -s gaussian stats.py` into a code cell and then press CTRL-Enter, the notebook will create a new input cell and load the function into it.

```
%load -s gaussian stats.py

def gaussian(x, mean, var):
    """returns normal distribution for x given a
    gaussian with the specified mean and variance.
    """
    return (np.exp((-0.5*(np.asarray(x)-mean)**2)/var) /
           math.sqrt(2*math.pi*var))
```

We will plot a Gaussian with a mean of 22 ( $\mu = 22$ ), with a variance of 4 ( $\sigma^2 = 4$ ), and then discuss what this means.

```
In [17]: from filterpy.stats import gaussian, norm_cdf
plot_gaussian_pdf(22, 4, mean_line=True, xlabel='$\circ$');
```



So what does this curve mean? Assume we have a thermometer which reads  $22^{\circ}\text{C}$ . No thermometer is perfectly accurate, and so we normally expect that thermometer will read slightly plus or minus that temperature each time we read it. However, a theorem called Central Limit Theorem states that if we make many measurements that the measurements will be normally distributed. So, when we look at this chart we can “sort of” think of it as representing the probability of the thermometer reading a particular value given the actual temperature of  $22^{\circ}\text{C}$ .

Recall that a Gaussian distribution is continuous. Think of an infinitely long straight line - what is the probability that a point you pick randomly is at 2. Clearly 0%, as there is an infinite number of choices to choose from. The same is true for normal distributions; in the graph above the probability of being exactly  $22^{\circ}\text{C}$  is 0% because there are an infinite number of values the reading can take.

So what then is this curve? It is something we call the probability density function. The area under the curve at any region gives you the probability of those values. So, for example, if you compute the area under the curve between 20 and 22 the resulting area will be the probability of the temperature reading being between those two temperatures.

We can think of this in Bayesian terms or frequentist terms. As a Bayesian, if the thermometer reads exactly  $22^{\circ}\text{C}$ , then our belief is described by the curve - our belief that the actual (system) temperature is near 22 is very high, and our belief that the actual temperature is near 18 is very low. As a frequentist we would say that if we took 1 billion temperature measurements of a system at exactly  $22^{\circ}\text{C}$ , then a histogram of the measurements would look like this curve.

So how do you compute the probability, or area under the curve? Well, you integrate the equation for the Gaussian

$$\int_{x_0}^{x_1} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu)^2/\sigma^2} dx$$

I wrote `filterpy.stats.norm_cdf` which computes the integral for you. So, for example, we can compute

```
In [18]: print('Probability of range 21.5 to 22.5 is {:.2f}%'.format(
    norm_cdf((21.5, 22.5), 22, 4)*100))
print('Probability of range 23.5 to 24.5 is {:.2f}%'.format(
    norm_cdf((23.5, 24.5), 22, 4)*100))
```

```
Probability of range 21.5 to 22.5 is 19.74%
Probability of range 23.5 to 24.5 is 12.10%
```

So the mean ( $\mu$ ) is what it sounds like - the average of all possible probabilities. Because of the symmetric shape of the curve it is also the tallest part of the curve. The thermometer reads  $22^{\circ}\text{C}$ , so that is what we used for the mean.

The notation for a normal distribution for a random variable  $X$  is  $X \sim \mathcal{N}(\mu, \sigma^2)$  where  $\sim$  means distributed according to. This means I can express the temperature reading of our thermometer as

$$\text{temp} \sim \mathcal{N}(22, 4)$$

This is an extremely important result. Gaussians allow me to capture an infinite number of possible values with only two numbers! With the values  $\mu = 22$  and  $\sigma^2 = 4$  I can compute the distribution of measurements for over any range.

## 3.8 The Variance and Belief

Since this is a probability density distribution it is required that the area under the curve always equals one. This should be intuitively clear - the area under the curve represents all possible outcomes, something happened, and the probability of something happening is one, so the density must sum to one. We can prove this ourselves with a bit of code. (If you are mathematically inclined, integrate the Gaussian equation from  $-\infty$  to  $\infty$ )

```
In [19]: print(norm_cdf((-1e8, 1e8), mu=0, var=4))
```

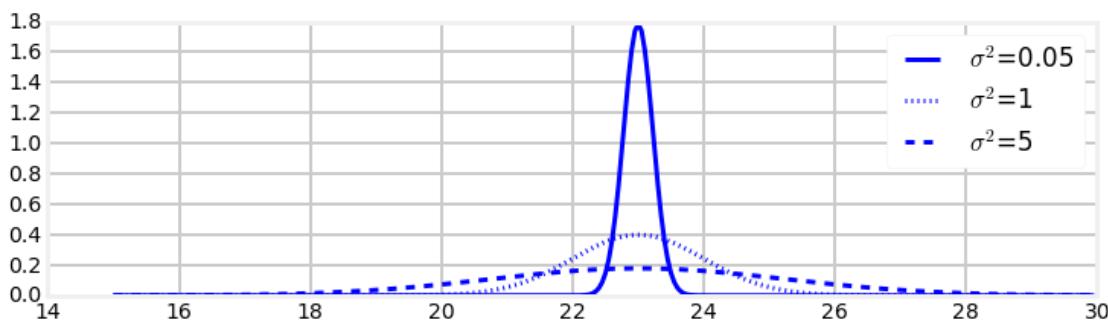
```
1.0
```

This leads to an important insight. If the variance is small the curve will be narrow. this is because the variance is a measure of how much the samples vary from the mean. To keep the area equal to 1, the curve must also be tall. On the other hand if the variance is large the curve will be wide, and thus it will also have to be short to make the area equal to 1.

Let's look at that graphically:

```
In [20]: import numpy as np
import matplotlib.pyplot as plt

xs = np.arange(15, 30, 0.05)
plt.plot(xs, gaussian(xs, 23, 0.05), label='$\sigma^2=0.05$', c='b')
plt.plot(xs, gaussian(xs, 23, 1), label='$\sigma^2=1$', ls=':', c='b')
plt.plot(xs, gaussian(xs, 23, 5), label='$\sigma^2=5$', ls='--', c='b')
plt.legend();
```



So what is this telling us? The Gaussian with  $\sigma^2 = 0.05$  is very narrow. It is saying that we believe  $x = 23$ , and that we are very sure about that. In contrast, the Gaussian with  $\sigma^2 = 5$  also believes that  $x = 23$ , but we are much less sure about that. Our belief that  $x = 23$  is lower, and so our belief about the likely possible values for  $x$  is spread out - we think it is quite likely that  $x = 20$  or  $x = 26$ , for example.  $\sigma^2 = 0.05$  has almost completely eliminated 22 or 24 as possible values, whereas  $\sigma^2 = 5$  considers them nearly as likely as 23.

If we think back to the thermometer, we can consider these three curves as representing the readings from three different thermometers. The curve for  $\sigma^2 = 0.05$  represents a very accurate thermometer, and curve for  $\sigma^2 = 5$  represents a fairly inaccurate one. Note the very powerful property the Gaussian distribution affords us - we can entirely represent both the reading and the error of a thermometer with only two numbers - the mean and the variance.

An equivalent formulation for a Gaussian is  $\mathcal{N}(\mu, 1/\tau)$  where  $\mu$  is the mean and  $\tau$  the precision.  $1/\tau = \sigma^2$ ; it is the reciprocal of the variance. While we do not use this formulation in this book, it underscores that the variance is a measure of how precise our data is. A small variance yields large precision - our measurement is very precise. Conversely, a large variance yields low precision - our belief is spread out across a large area. You should become comfortable with thinking about Gaussians in these equivalent forms. In Bayesian terms Gaussians reflect our belief about a measurement, they express the precision of the measurement, and they express how much variance there is in the measurements. These are all different ways of stating the same fact.

I'm getting ahead of myself, but in the next chapters we will use Gaussians to express our belief in things like the estimated position of the object we are tracking, or the accuracy of the sensors we are using.

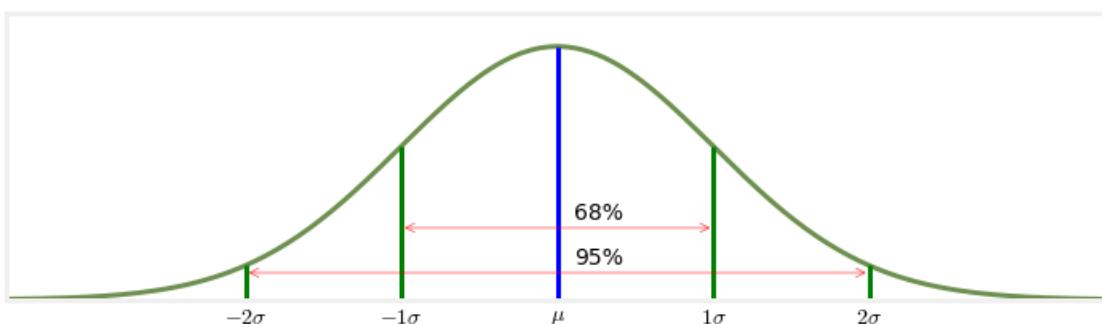
### 3.9 The 68-95-99.7 Rule

It is worth spending a few words on standard deviation now. The standard deviation is a measure of how much variation from the mean exists. For Gaussian distributions, 68% of all the data falls within one standard deviation ( $1\sigma$ ) of the mean, 95% falls within two standard deviations ( $2\sigma$ ), and 99.7% within three ( $3\sigma$ ). This is often called the 68-95-99.7 rule. So if you were told that the average test score in a class was 71 with a standard deviation of 9.4, you could conclude that 95% of the students received a score between 52.2 and 89.8 if the distribution is normal (that is calculated with  $71 \pm (2 * 9.4)$ ).

Finally, these are not arbitrary numbers. If the Gaussian for our position is  $\mu = 22$  meters, then the standard deviation also has units meters. Thus  $\sigma = 0.2$  implies that 68% of the measurements range from 21.8 to 22.2 meters. Variance is the standard deviation squared, so  $\sigma^2 = .04$  meters<sup>2</sup>.

The following graph depicts the relationship between the standard deviation and the normal distribution.

```
In [21]: from gaussian_internal import display_stddev_plot
with figsize(y=3):
    display_stddev_plot()
```



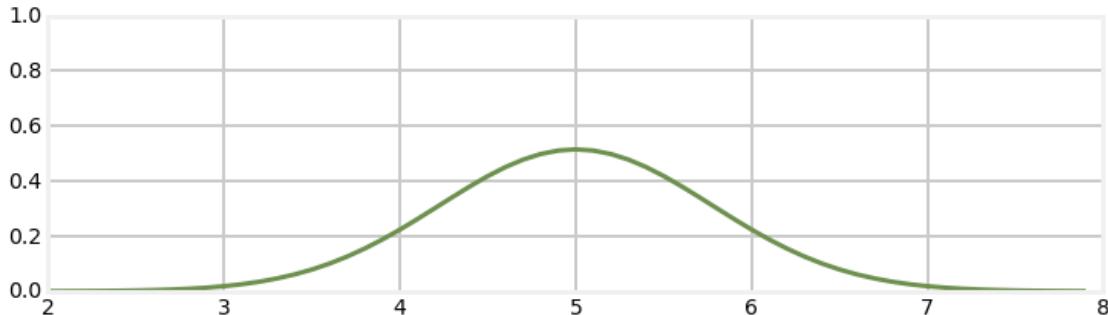
## 3.10 Interactive Gaussians

For those that are reading this in a Jupyter Notebook, here is an interactive version of the Gaussian plots. Use the sliders to modify  $\mu$  and  $\sigma^2$ . Adjusting  $\mu$  will move the graph to the left and right because you are adjusting the mean, and adjusting  $\sigma^2$  will make the bell curve thicker and thinner.

```
In [22]: import math
from IPython.html.widgets import interact, interactive, fixed

set_figsize(y=3)
def plt_g(mu,variance):
    xs = np.arange(2, 8, 0.1)
    ys = gaussian(xs, mu, variance)
    plt.plot(xs, ys)
    plt.ylim((0, 1))

interact (plt_g, mu=(0., 10), variance = (.2, 1.));
```



Finally, if you are reading this online, here is an animation of a Gaussian. First, the mean is shifted to the right. Then the mean is centered at  $\mu = 5$  and the variance is modified.

## 3.11 Computational Properties of the Gaussian

Recall how our discrete Bayesian filter worked. We had a vector implemented as a NumPy array representing our belief at a certain moment in time. When we integrated another measurement into our belief using the `update()` function we had to multiply probabilities together, and when we performed the motion step using the `predict()` function we had to shift and add probabilities. I've promised you that the Kalman filter uses essentially the same process, and that it uses Gaussians instead of histograms, so you might reasonable expect that we will be multiplying, adding, and shifting Gaussians in the Kalman filter.

A typical textbook would directly launch into a multi-page proof of the behavior of Gaussians under these operations, but I don't see the value in that right now. I think the math will be much more intuitive and clear if we just start developing a Kalman filter using Gaussians. I will provide the equations for multiplying and shifting Gaussians at the appropriate time. You will then be able to develop a physical intuition for what these operations do, rather than be forced to digest a lot of fairly abstract math.

The key point, which I will only assert for now, is that all the operations are very simple, and that they preserve the properties of the Gaussian. This is somewhat remarkable, in that the Gaussian is a nonlinear function, and typically if you multiply a nonlinear equation with itself you end up with a different equation. For example, the shape of `sin(x)sin(x)` is very different from `sin(x)`. But the result of multiplying two Gaussians is yet another Gaussian. This is a fundamental property, and a key reason why Kalman filters are computationally feasible. Said another way, Kalman filters use Gaussians because they are so computationally nice.

## 3.12 Computing Probabilities with `scipy.stats`

In this chapter I have used custom code from FilterPy for computing Gaussians, plotting, and so on. I chose to do that to give you a chance to look at the code and see how these functions are implemented. However, Python comes with “batteries included” as the saying goes, and it comes with a wide range of statistics functions in the module `scipy.stats`. I find the performance of some of the functions rather slow (the `scipy.stats` documentation contains a warning to this effect), but this is offset by the fact that this is standard code available to everyone, and it is well tested. So let’s walk through how to use `scipy.stats` to compute statistics and probabilities.

The `scipy.stats` module contains a number of objects which you can use to compute attributes of various probability distributions. The full documentation for this module is here: <http://docs.scipy.org/doc/scipy/reference/stats.html>. However, we will focus on the `norm` variable, which implements the normal distribution. Let’s look at some code that uses `scipy.stats.norm` to compute a Gaussian, and compare its value to the value returned by the `gaussian()` function from FilterPy.

```
In [23]: from scipy.stats import norm
import filterpy.stats
print(norm(2, 3).pdf(1.5))
print(filterpy.stats.gaussian(x=1.5, mean=2, var=3*3))

0.131146572034
0.131146572034
```

The call `norm(2, 3)` creates what scipy calls a ‘frozen’ distribution - it creates and returns an object with a mean of 2 and a standard deviation of 3. You can then use this object multiple times to get the probability density of various values, like so:

```
In [24]: n23 = norm(2, 3)
print('pdf of 1.5 is      %.4f' % n23.pdf(1.5))
print('pdf of 2.5 is also %.4f' % n23.pdf(2.5))
print('pdf of 2 is        %.4f' % n23.pdf(2))

pdf of 1.5 is      0.1311
pdf of 2.5 is also 0.1311
pdf of 2 is        0.1330
```

If we look at the documentation for `scipy.stats.norm` here[2] we see that there are many other functions that `norm` provides.

For example, we can generate  $n$  samples from the distribution with the `rvs()` function.

```
In [25]: np.set_printoptions(precision=3, linewidth=50)
print(n23.rvs(size=15))
```

```
[ 0.409  2.668  6.657  2.062  4.686  4.65  -2.214
 -2.35   6.065  1.219 -0.904  5.647  8.615 -4.476
 0.463]
```

We can get the cumulative distribution function (CDF), which is the probability that a randomly drawn value from the distribution is less than or equal to  $x$ .

```
In [26]: # probability that a random value is less than the mean 2
print(n23.cdf(2))
```

0.5

We can get various properties of the distribution:

```
In [27]: print('variance is', n23.var())
print('standard deviation is', n23.std())
print('mean is', n23.mean())
```

```
variance is 9.0
standard deviation is 3.0
mean is 2.0
```

There are many other functions available, and if you are interested I urge you to peruse the documentation. Sometimes the documentation is terse, but with a bit of searching you can find out what a function does and some examples of how to use it. Most of this functionality is not of immediate interest to the book, so I will leave the topic in your hands to explore. The SciPy tutorial [3] is quite approachable, and I suggest starting there.

## 3.13 Fat Tails

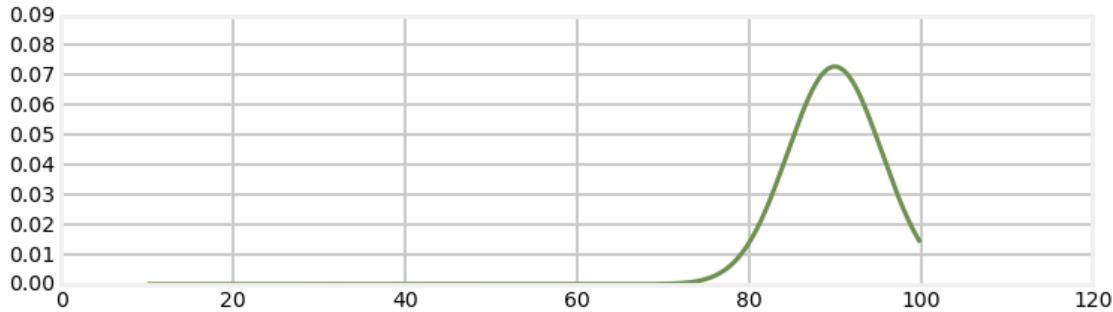
Earlier I spoke very briefly about the central limit theorem, which states that under certain conditions the arithmetic sum of any independent random variable will be normally distributed, regardless of how the random variables are distributed. This is extremely important for (at least) two reasons. First, nature is full of distributions which are not normal, but when we apply the central limit theorem over large populations we end up with normal distributions. Second, Gaussians are mathematically tractable. We will see this more as we develop the Kalman filter theory, but there are very nice closed form solutions for operations on Gaussians that allow us to use them analytically.

However, a key part of the proof is “under certain conditions”. These conditions often do not hold for the physical world. The resulting distributions are called fat tailed. Tails is a colloquial term for the far left and right side parts of the curve where the probability density is close to zero.

Let’s consider a trivial example. We think of things like test scores as being normally distributed. If you have ever had a professor ‘grade on a curve’ you have been subject to this assumption. But of course test scores cannot follow a normal distribution. This is because the distribution assigns a nonzero probability distribution for any value, no matter how far from the mean. So, for example, say your mean is 90 and the standard deviation is 13. The normal distribution assumes that there is a large chance of somebody getting a 90, and a small chance of somebody getting a 40. However, it also implies that there is a tiny chance of somebody getting a grade of -10, or 150. It assigns an infinitesimal chance of getting a score of  $-10^{300}$  or  $10^{32986}$ . The tails of a Gaussian distribution are infinitely long.

But for a test we know this is not true. Ignoring extra credit, you cannot get less than 0, or more than 100. Let’s plot this range of values using a normal distribution.

```
In [28]: xs = np.arange(10,100, 0.05)
ys = [gaussian(x, 90, 30) for x in xs]
plt.plot(xs, ys, label='var=0.2')
plt.xlim((0,120))
plt.ylim(0, 0.09);
```



The area under the curve cannot equal 1, so it is not a probability distribution. What actually happens is that more students than predicted by a normal distribution get scores nearer the upper end of the range (for example), and that tail becomes ‘fat’. Also, the test is probably not able to perfectly distinguish incredibly minute differences in skill in the students, so the distribution to the left of the mean is also probably a bit bunched up in places. The resulting distribution is called a fat tail distribution.

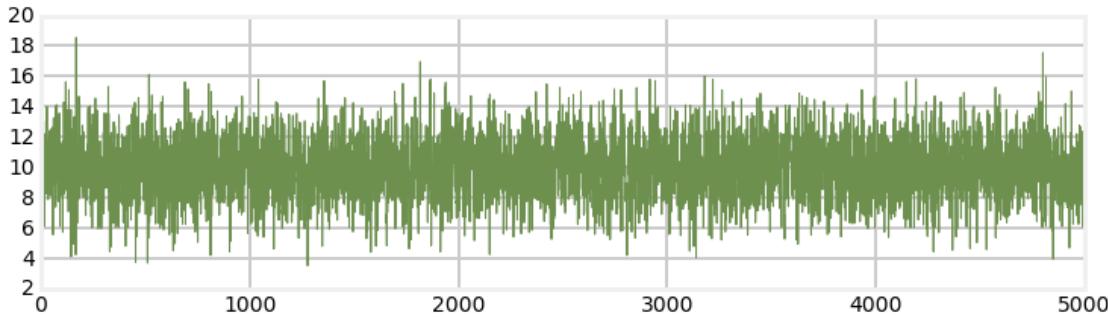
Kalman filters use sensors to measure the world. The errors in sensor’s measurements are rarely truly Gaussian. It is far too early to be talking about the difficulties that this presents to the Kalman filter designer. It is worth keeping in the back of your mind the fact that the Kalman filter math is based on an idealized model of the world. For now I will present a bit of code that I will be using later in the book to form fat tail distributions to simulate various processes and sensors. This distribution is called the Student’s t distribution.

Let’s say I want to model a sensor that has some white noise in the output. For simplicity, let’s say the signal is a constant 10, and the standard deviation of the noise is 2. We can use the function `numpy.random.randn()` to get a random number with a mean of 0 and a standard deviation of 1. So I could simulate this sensor with

```
In [29]: from numpy.random import randn
def sense():
    return 10 + randn()*2
```

Let’s plot that signal and see what it looks like.

```
In [30]: zs = [sense() for i in range(5000)]
plt.plot(zs, lw=1);
```



That looks like I would expect. The signal is centered around 10. A standard deviation of 2 means that 68% of the measurements will be within  $\pm 2$  of 10, and 99% will be within  $\pm 6$  of 10, and that looks like what is happening.

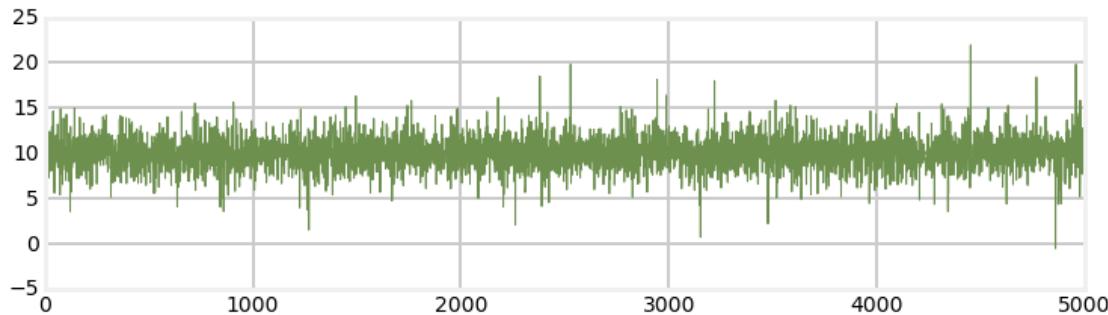
Now let's look at a fat tailed distribution generated with the Student's T distribution. I will not go into the math, but just give you the source code for it and then plot a distribution using it.

```
In [31]: import random
import math

def rand_student_t(df, mu=0, std=1):
    """Return random number distributed by Student's t
    distribution with 'df' degrees of freedom with the
    specified mean and standard deviation.
    """
    x = random.gauss(0, std)
    y = 2.0*random.gammavariate(0.5*df, 2.0)
    return x / (math.sqrt(y / df)) + mu
```

```
In [32]: def sense_t():
    return 10 + rand_student_t(7)*2

zs = [sense_t() for i in range(5000)]
plt.plot(zs, lw=1);
```



We can see from the plot that while the output is similar to the normal distribution there are outliers that go far more than 3 standard deviations from the mean (7 to 13). This is what causes the 'fat tail'.

It is unlikely that the Student's T distribution is an accurate model of how your sensor (say, a GPS or Doppler) performs, and this is not a book on how to model physical systems. However, it does produce reasonable data to test your filter's performance when presented with real world noise. We will be using distributions like these throughout the rest of the book in our simulations and tests.

This is not an idle concern. The Kalman filter equations assume the noise is normally distributed, and perform sub-optimally if this is not true. Designers for mission critical filters, such as the filters on spacecraft, need to master a lot of theory and empirical knowledge about the performance of the sensors on their spacecraft.

The code for `rand_student_t` is included in `filterpy.stats`. You may use it with

```
from filterpy.stats import rand_student_t
```

## 3.14 Summary and Key Points

This chapter is a poor introduction to statistics in general. I've only covered the concepts that needed to use Gaussians in the remainder of the book, no more. What I've covered will not get you very far if you intend to read the Kalman filter literature. If this is a new topic to you I suggest reading a statistics textbook. I've always liked the Schaum series for self study, and Alan Downey's [Think Stats](#) [5] is also very good.

The following points **must** be understood by you before we continue:

- Normals express a continuous probability distribution
- They are completely described by two parameters: the mean ( $\mu$ ) and variance ( $\sigma^2$ )
- $\mu$  is the average of all possible values
- The variance  $\sigma^2$  represents how much our measurements vary from the mean
- The standard deviation ( $\sigma$ ) is the square root of the variance ( $\sigma^2$ )
- Many things in nature approximate a normal distribution

## 3.15 References

- [1] [https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python/blob/master/Supporting\\_Notebooks/Computing\\_and\\_Probability.ipynb](https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python/blob/master/Supporting_Notebooks/Computing_and_Probability.ipynb)
- [2] <http://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.norm.html>
- [3] <http://docs.scipy.org/doc/scipy/reference/tutorial/stats.html>
- [4] Huber, Peter J. [Robust Statistical Procedures](#), Second Edition. Society for Industrial and Applied Mathematics, 1996.
- [5] Downey, Alan. [Think Stats](#), Second Edition. O'Reilly Media.  
<https://github.com/AllenDowney/ThinkStats2>  
<http://greenteapress.com/thinkstats/>

# Chapter 4

## One Dimensional Kalman Filters

### 4.1 One Dimensional Kalman Filters

Now that we understand the discrete Bayes filter and Gaussians we are prepared to implement a Kalman filter. We will do this exactly as we did the discrete Bayes filter - rather than starting with equations we will develop the code step by step based on reasoning about the problem. “One dimensional” just means that the filter only tracks one variable, such as position on the x-axis. In the subsequent chapters we will learn a more general form that can track many variables simultaneously, such as position, velocity, and accelerations in x, y, z.

### 4.2 Tracking A Dog

As in the **Discrete Bayes Filter** chapter we will be tracking a moving object in a long hallway at work, such as a dog or robot. Assume that in our latest hackathon someone created an RFID tracker that provides a reasonably accurate position for our dog. Suppose the hallway is 100 meters long. The sensor returns the distance of the dog from the left end of the hallway in meters. So, 23.4 would mean the dog is 23.4 meters from the left end of the hallway.

Naturally, the sensor is not perfect. A reading of 23.4 could correspond to a real position of 23.7, or 23.0. However, it is very unlikely to correspond to a real position of say 47.6. Testing during the hackathon confirmed this result - the sensor is ‘reasonably’ accurate, and while it had errors, the errors are small. Furthermore, the errors seemed to be evenly distributed on both sides of the measurement; a true position of 23 m would be equally likely to be measured as 22.9 as 23.1.

Before we try to implement the filter, let’s start by implementing a simulation of the dog’s movement and of the sensor. After all, the filter needs to filter data collected from the dog, so we need to understand this movement and sensor behavior to have any hope of filtering it.

#### 4.2.1 Simulating the Sensor Measurement

Accurately simulating an RFID tracker is beyond the scope of this book, so we will use a simple model. We said that the sensor is equally likely to give a measurement that is too large as too small. Furthermore, it is more likely to be a small error than a large one. In the **Discrete Bayes** chapter we modelled this error with a probability. But now we understand Gaussians so let’s use them instead. We can model the error in a system with the Gaussian  $\mathcal{N}(\mu, \sigma^2)$ , where  $\mu$  is the state (such as position or velocity) and  $\sigma^2$  is the variance in that state. So instead of saying our measurement is 22.9, we would say our measurement is 22.9 with a variance of 0.35, or  $\mathcal{N}(22.9, 0.35)$ . The variance accounts for the error in the sensor.

This is a simple model, but it works very well in practice. Error distributions of many sensors approximate a normal distribution. Of course, if the sensor errors are not well modelled by Gaussians your simulations of the filter performance are likely to mislead you.

In the literature you will often see this characterized with an equation like:

$$z = h(x) + \epsilon_z$$

Here  $h(x)$  is the measurement model - a function which describes how measurements are made from the current system state. In other words,  $h(x)$  is the measurement we would get for the state  $x$  if the sensor was perfect.  $\epsilon_z$  is the error that is in the measurement, and hence  $z$  is the actual measurement that results. This is not meant to imply that we can somehow measure or know the exact value of  $\epsilon_z$ , it is a mathematical definition that can be restated as

$$\epsilon_z = z - h(x)$$

The utility of this equation will become apparent later in the book. For now recognize that we will be modelling  $\epsilon_z$  with a Gaussian. If  $\epsilon_z$  is close to being Gaussian our model will be accurate, and if it is not our model may not work well.

### 4.2.2 Simulating the Dog's movement

Let's assume that the dog moves at a constant velocity. We can use Newton's equation to compute the new position based on the current velocity and previous position like so:

$$x_k = v_k \Delta t + x_{k-1}$$

This is a model of the dog's behavior, and we call it the system model or process model. However, clearly it is not perfect. The dog might speed up or slow down due to hills, winds, or whimsy. We could add those things to the model if we knew about them, but there will always be more things that we don't know. This is not due to the tracked object being a living creature. The same holds if we are tracking a missile, aircraft, or paint nozzle on a factory robotic painter. There will always be unpredictable errors in physical systems. We can model this mathematically by saying the dog's actual position is the prediction that comes from our process model plus what we call the process noise. In the literature this is usually written something like this:

$$x = f(x) + \epsilon$$

Here  $f(x)$  is our process model and  $\epsilon$  is some unknown error between the process model's prediction of position and the actual position. Clearly  $\epsilon$  will change over time, and as with the measurement equation in the last section we are not implying that  $\epsilon$  can be derived analytically.

### 4.2.3 Implementation in Python

We will want our implementation to correctly model the noise both in the measurement and the process model. You may recall from the **Gaussians** chapter that we can use `numpy.random.randn()` to generate a random number with a mean of zero and a standard deviation of one. Thus, if we want a random number with a standard deviation of 0.5 we'd multiply the value returned by `randn()` by 0.5.

Our model assumes that the velocity of our dog is constant, but things like hills and obstacles make that very unlikely to be true. We can simulate this by adding a randomly generated value to the velocity each time we compute the dog's position.

```
noisy_vel = vel + randn()*velocity_std
x = x + noisy_vel*dt
```

To model the sensor we need to take the current position of the dog and then add some noise to it proportional to the Gaussian of the measurement noise. This is simply

```
z = x + randn()*measurement_std
```

I don't like the bookkeeping required to keep track of the various values for position, velocity, and standard deviations so I will implement the simulation as a class.

```
In [2]: import math
import matplotlib.pyplot as plt
from numpy.random import randn

class DogSimulation(object):
    def __init__(self, x0=0, velocity=1,
                 measurement_var=0.0,
                 process_var=0.0):
        """
        x0 : initial position
        velocity: (+=right, -=left)
        measurement_var: variance in measurement m^2
        process_var: variance in process (m/s)^2
        """
        self.x = x0
        self.velocity = velocity
        self.meas_noise = math.sqrt(measurement_var)
        self.process_noise = math.sqrt(process_var)

    def move(self, dt=1.0):
        """
        Compute new position of the dog in dt seconds.
        """
        dx = self.velocity + randn()*self.process_noise
        self.x += dx * dt

    def sense_position(self):
        """
        Returns measurement of new position in meters.
        """
        measurement = self.x + randn()*self.meas_noise
        return measurement

    def move_and_sense(self):
        """
        Move dog, and return measurement of new position in meters"""
        self.move()
        return self.sense_position()
```

The constructor `__init__()` initializes the `DogSimulation` class with an initial position `x0`, velocity `vel`, and the variance in the measurement and noise. The `move()` method moves the dog based on its velocity and the process noise. The `sense_position()` method computes and returns a measurement of the dog's position based on the current position and the sensor noise. `move_and_sense()` first performs `move()`, then `sense_position()`.

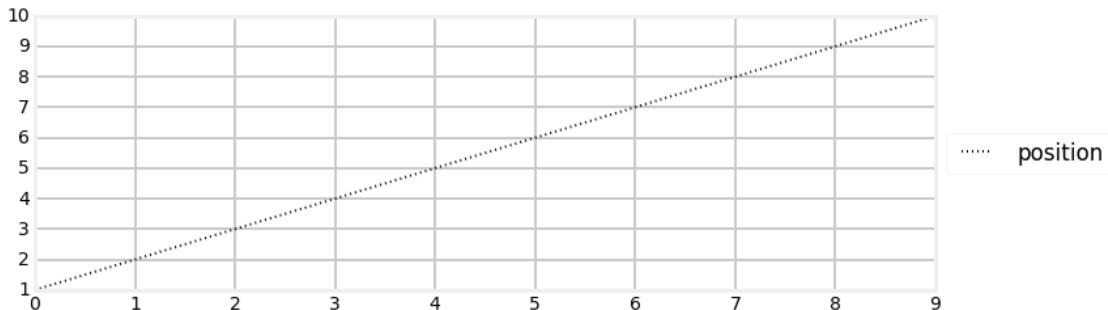
We need to convert the variances that are passed into `__init__()` into standard deviations because `randn()` is scaled by the standard deviation. Variance is the square of the standard deviation, so we take the square root of the variance.

Now let's look at the output of the `DogSimulation` class. We will start by setting the variances to zero. Zero variance means there is no noise in the signal or in the process model, so the results should be a straight line. I use code from `book_plots` to plot the tracks without a lot of superfluous code distracting you from the core ideas.

```
In [3]: import matplotlib.pyplot as plt
        import book_plots as bp
        import numpy as np

        dog = DogSimulation(measurement_var=0.0)
        N = 10
        xs = np.zeros(N)
        for i in range(N):
            dog.move()
            xs[i] = dog.sense_position()
            print("%.1f" % xs[i], end=' ')
        bp.plot_track(xs, label='position')
        bp.show_legend();
```

1.0 2.0 3.0 4.0 5.0 6.0 7.0 8.0 9.0 10.0

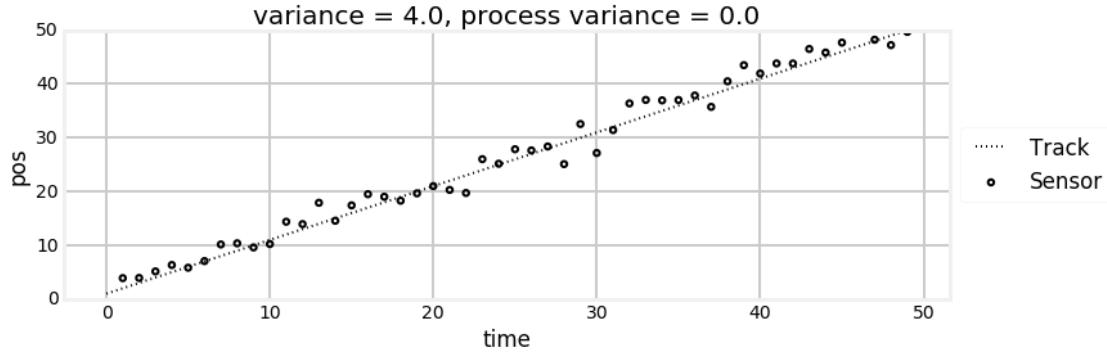


The constructor initialized the dog at position 0 with a velocity of 1 m/s (move 1.0 m to the right). So we would expect to see an output of 1..10, and indeed that is what we see.

Now let's inject some noise in the signal.

```
In [4]: from kf_internal import plot_dog_track
def test_sensor(measurement_var, process_var=0.0):
    dog = DogSimulation(0., 1., measurement_var, process_var)
    N = 50
    xs = np.zeros(N)
    for i in range(N):
        dog.move()
        xs[i] = dog.sense_position()

    plot_dog_track(xs, measurement_var, process_var)
test_sensor(measurement_var=4.0)
```



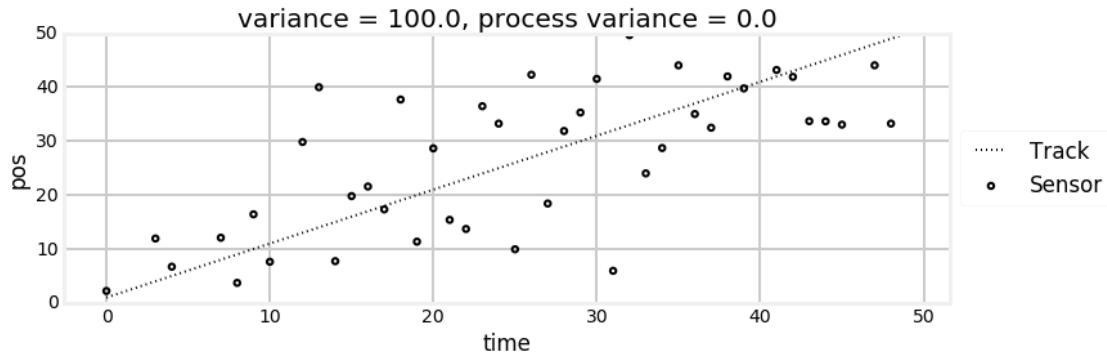
NumPy uses a random number generator to generate the normal distribution samples. The numbers I see as I write this are unlikely to be the ones that you see. If you run the cell above multiple times, you should get a slightly different result each time. I could use `numpy.random.seed(some_value)` to force the results to be the same each time. This would simplify my explanations in some cases, but would ruin the interactive nature of this chapter. To get a real feel for how normal distributions and Kalman filters work you will probably want to run cells several times, observing what changes, and what stays roughly the same.

The dotted line (labeled ‘track’) shows the actual position of the dog, and the circles (labeled ‘sensor’) depict the noisy signal produced by the simulated RFID sensor. Please note that the track was manually plotted - we do not yet have a filter that recovers that information from the sensor measurements!

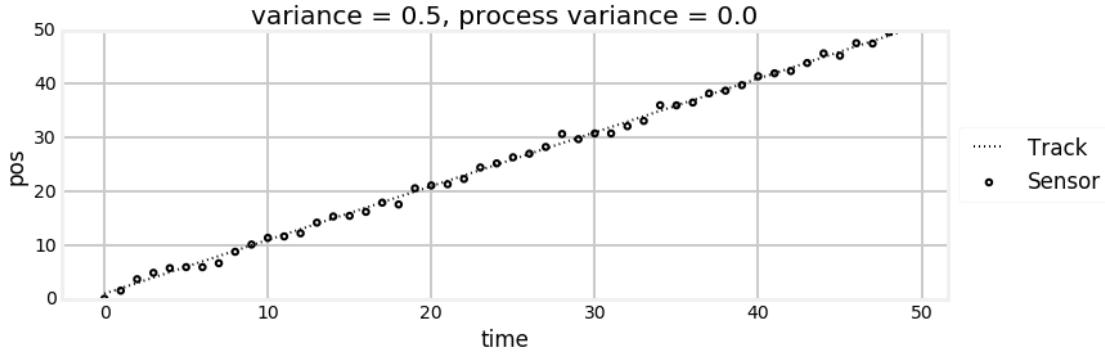
If you are running this in an interactive Jupyter Notebook, I strongly urge you to run the script several times in a row. You can do this by putting the cursor in the cell containing the Python code and pressing CTRL+Enter. Each time it runs you should see a different sensor output.

I also urge you to adjust the noise setting to see the result of various values. However, since you may be reading this in a static format I will show several examples. The first plot shows the noise set to 100.0, and the second shows noise set to 0.5.

In [5]: `test_sensor(measurement_var=100.0)`



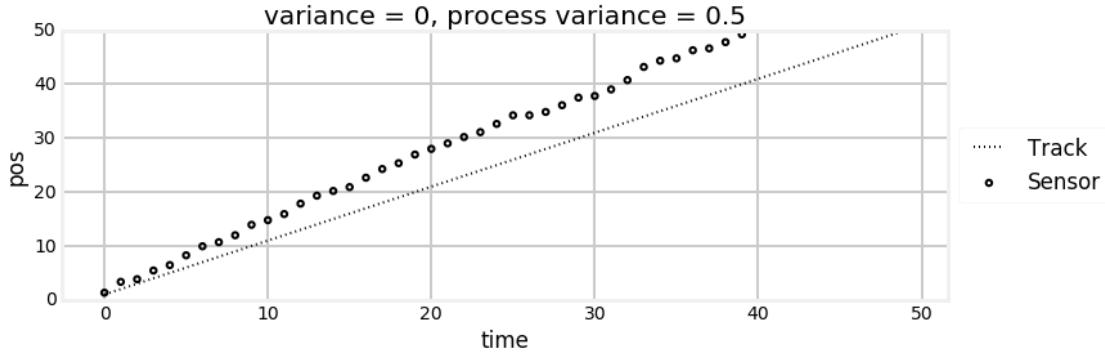
In [6]: `test_sensor(measurement_var=0.5)`



Now let's see the effect of the process noise. My simulation is not meant to exactly model any given physical process. On each call to `sense_position()` I modify the current velocity by randomly generated normal process noise. However, I strive to keep the velocity close to the initial value, as I am assuming that there is a control mechanism in place trying to maintain the same speed. In this case, the control mechanism would be the dog's brain! For an automobile it could be a cruise control or the human driver.

So let's first look at the plot with some process noise but no measurement noise to obscure the results.

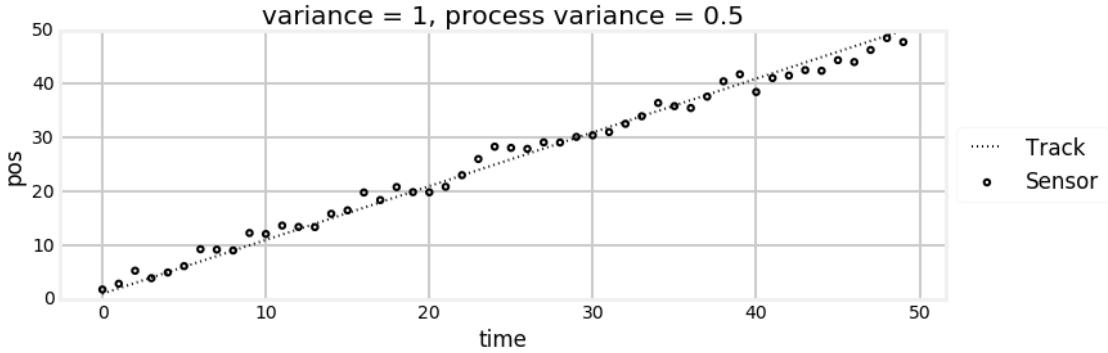
```
In [7]: np.random.seed(1234) # this will make the dog wander away
test_sensor(measurement_var=0, process_var=0.5)
```



We can see that the position wanders from the ideal track. You may have thought that the track would have wandered back and forth like the measurement noise did, but recall that we are modifying velocity on each, not the position. So once the track has deviated, it will remain deviated until the random changes in velocity happen to result in the track going back to the original track.

Finally, let's look at the combination of measurement noise and process noise. Depending on the random number generation the sensor reading may lie on top of the track, or wander away from it. Try running it several times to observe the various possibilities.

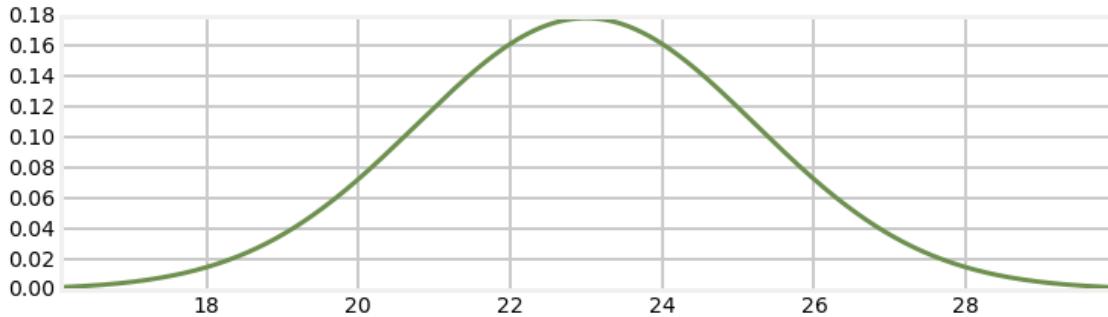
```
In [8]: test_sensor(measurement_var=1, process_var=0.5)
```



## 4.3 Math with Gaussians

We can express our belief in the dog's position with a Gaussian. Say we believe that our dog is at 23 meters, and the variance in that belief is  $5 \text{ m}^2$ , or  $\mathcal{N}(23, 5)$ . We can represent that in a plot:

```
In [9]: from book_format import set_figsize, figsize
import filterpy.stats as stats
with figsize(y=3):
    stats.plot_gaussian_pdf(mean=23, variance=5)
```



This chart depicts our belief, in other words our uncertainty about the dog's position. Notice how this relates to the bar charts from the Discrete Bayes chapter. In that chapter we drew bars of various heights to depict the probability that the dog was at any given position. In many cases those bars took on a shape very similar to this chart. The differences are that the bars depict the probability of a discrete position, whereas this chart depicts the probability distribution of a continuous range of positions.

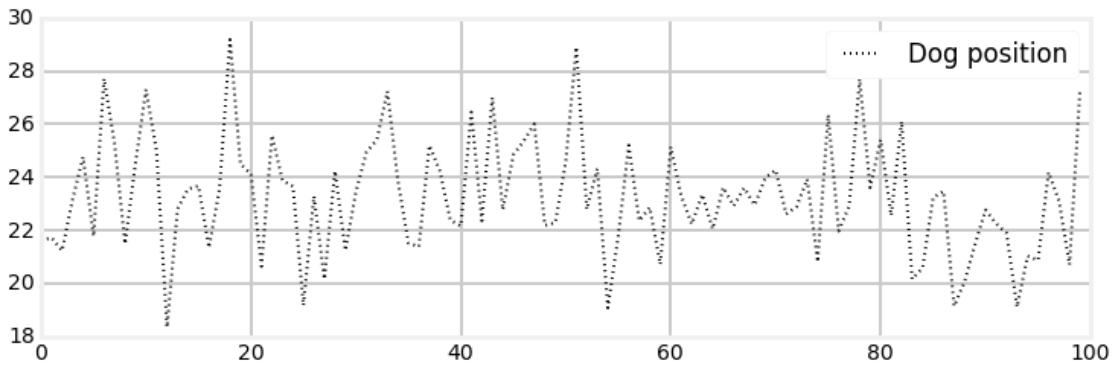
This graph corresponds to a fairly inexact belief. While we believe that the dog is at 23 m, note that roughly speaking positions 21 m to 25 m are quite likely as well. Let's assume for the moment our dog is standing still, and we query the sensor again. This time it returns 23.2 m as the position. Can we use this additional information to improve our estimate of the dog's position?

Intuition suggests we can. Consider: if we read the sensor 100 times and each time it returned a value between 21 and 25, all centered around 23, we should be very confident that the dog is near 23. Of course,

a different physical interpretation is possible. Perhaps our dog was randomly wandering back and forth in a way that exactly emulated a normal distribution. But that seems extremely unlikely - I certainly have never seen a dog do that. So the only reasonable assumption is that the dog was mostly standing still at 23.0.

Let's look at 100 sensor readings in a plot, where we assume the dog stands still at 23 m.

```
In [10]: dog = DogSimulation(23, 0, measurement_var=5, process_var=0.0)
xs = range(100)
ys = []
for _ in xs:
    dog.move()
    ys.append(dog.sense_position())
bp.plot_track(xs, ys, label='Dog position')
plt.legend(loc='best');
```



Eyeballing this confirms our intuition - no dog moves like this. However, noisy sensor data certainly looks like this. So let's proceed and try to solve this mathematically. But how?

Recall the code for adding a measurement to a preexisting belief from the **Discrete Bayes** chapter:

```
def update(pos, measure, p_hit, p_miss):
    q = array(pos, dtype=float)
    for i in range(len(hallway)):
        if hallway[i] == measure:
            q[i] = pos[i] * p_hit
        else:
            q[i] = pos[i] * p_miss
    normalize(q)
    return q
```

The algorithm is computing:

```
new_belief = prior_belief * measurement * sensor_error
```

The `measurement * sensor_error` term might not be obvious, but recall that measurement in this case was always 1 or 0, and so it was left out for convenience. Here `prior` carries the both the colloquial sense of `previous`, but also the Bayesian meaning. In Bayesian terms the `prior` is the probability after the prediction and before the measurements have been incorporated.

If we are implementing this with Gaussians, we might expect it to be implemented as:

```
new_gaussian = measurement * old_gaussian
```

where measurement is a Gaussian returned from the sensor. But does that make sense? Can we multiply Gaussians? If we multiply a Gaussian with a Gaussian is the result another Gaussian, or something else?

The algebra for this is not extremely hard, but it is messy. It uses Bayes theorem to compute the posterior (new Gaussian) given the prior (old Gaussian) and the likelihood (the probability for the measurement). I've placed the math at the bottom of the chapter if you are interested in the details. Here I will present the results. The subscript  $z$  stands for the measurement.

$$\begin{aligned}\mathcal{N}(\mu_{\text{posterior}}, \sigma_{\text{posterior}}^2) &= \mathcal{N}(\mu_{\text{prior}}, \sigma_{\text{prior}}^2) \times \mathcal{N}(\mu_z, \sigma_z^2) \\ &= \mathcal{N}\left(\frac{\sigma_{\text{prior}}^2 \mu_z + \sigma_z^2 \mu_{\text{prior}}}{\sigma_{\text{prior}}^2 + \sigma_z^2}, \frac{1}{\frac{1}{\sigma_{\text{prior}}^2} + \frac{1}{\sigma_z^2}}\right)\end{aligned}$$

In other words the result of multiplying two Gaussians is a Gaussian is

$$\begin{aligned}\mu &= \frac{\sigma_{\text{prior}}^2 \mu_z + \sigma_z^2 \mu_{\text{prior}}}{\sigma_{\text{prior}}^2 + \sigma_z^2}, \\ \sigma^2 &= \frac{1}{\frac{1}{\sigma_{\text{prior}}^2} + \frac{1}{\sigma_z^2}} = \frac{\sigma_{\text{prior}}^2 \sigma_z^2}{\sigma_{\text{prior}}^2 + \sigma_z^2}\end{aligned}$$

Without doing a deep analysis we can immediately infer some things. First and most importantly the result of multiplying two Gaussians is another Gaussian. The mean is a scaled sum of the prior and the measurement. The variance is a combination of the variances of the variances of the input. We conclude from this that the variances are completely unaffected by the values of the mean!

Let's examine the result of multiplying  $\mathcal{N}(23, 5)$  with itself. This corresponds to getting 23.0 as the sensor value twice in a row. But before you look at the result, what do you think the result will look like? What should the new mean be? Will the variance be wider, narrower, or the same?

```
In [11]: import filterpy.stats as stats

def multiply(mu1, var1, mu2, var2):
    # avoid division by zero
    if var1 == 0.0: var1 = 1.e-80
    if var2 == 0.0: var2 = 1.e-80

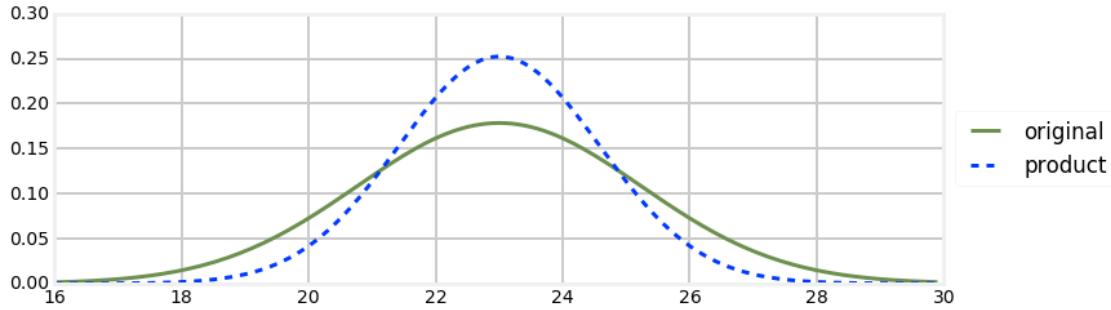
    mean = (var1*mu2 + var2*mu1) / (var1 + var2)
    variance = 1 / (1/var1 + 1/var2)
    return (mean, variance)

xs = np.arange(16, 30, 0.1)

mean1, var1 = 23, 5
mean, var = multiply(mean1, var1, mean1, var1)

ys = [stats.gaussian(x, mean1, var1) for x in xs]
plt.plot(xs, ys, label='original')

ys = [stats.gaussian(x, mean, var) for x in xs]
plt.plot(xs, ys, label='product', ls='--')
bp.show_legend();
```



The result of the multiplication is taller and narrower than the original Gaussian but the mean is the same. Does this match your intuition of what the result should have been?

If we think of the Gaussians as two measurements, this makes sense. If I measure twice and get 23 meters each time, I should conclude that the length is close to 23 meters. Thus the mean should be 23. I am more confident with two measurements than with one, so the variance of the result should be smaller.

“Measure twice, cut once” is a useful saying and practice due to this fact! The Gaussian is a mathematical model of this physical fact, so we should expect the math to follow our physical process.

Now let’s multiply two Gaussians (or equivalently, two measurements) that are partially separated. In other words, their means will be different, but their variances will be the same. What do you think the result will be? Think about it, and then look at the graph.

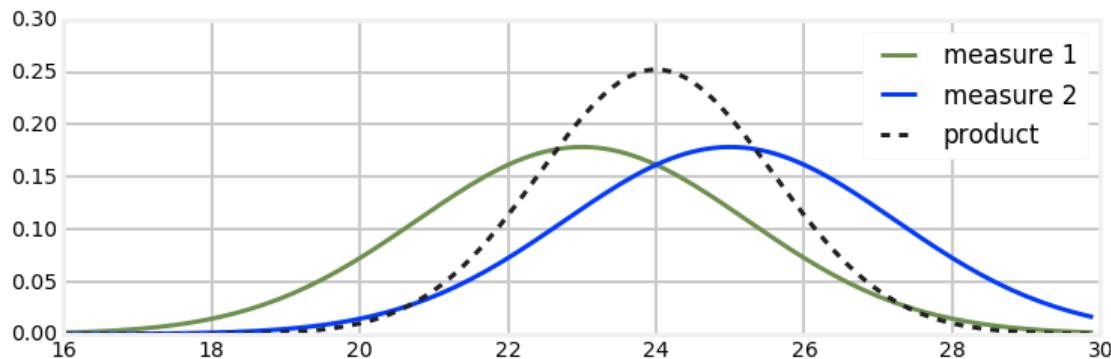
```
In [12]: xs = np.arange(16, 30, 0.1)

mean1, var1 = 23, 5
mean2, var2 = 25, 5
mean, var = multiply(mean1, var1, mean2, var2)

ys = [stats.gaussian(x, mean1, var1) for x in xs]
plt.plot(xs, ys, label='measure 1')

ys = [stats.gaussian(x, mean2, var2) for x in xs]
plt.plot(xs, ys, label='measure 2')

ys = [stats.gaussian(x, mean, var) for x in xs]
plt.plot(xs, ys, label='product', ls='--')
plt.legend();
```



Another beautiful result! If I handed you a measuring tape and asked you to measure the distance from table to a wall, and you got 23 meters, and then a friend make the same measurement and got 25 meters, your best guess must be 24 meters if you trust the skills of your friends equally.

That is fairly counter-intuitive in more complicated situations, so let's consider it further. Perhaps a more reasonable assumption would be that either you or your coworker made a mistake, and the true distance is either 23 or 25, but certainly not 24. Surely that is possible. However, suppose the two measurements you reported as 24.01 and 23.99. In that case you would agree that in this case the best guess for the correct value is 24? Which interpretation we choose depends on the properties of the sensors we are using.

This topic is fairly deep, and I will explore it once we have completed our Kalman filter. For now I will merely say that the Kalman filter requires the interpretation that measurements are accurate, with Gaussian noise, and that a large error caused by misreading a measuring tape is not Gaussian noise.

Dealing with incorrect measurements, such as measuring the distance to the wrong object, is a topic we will need to treat separately. The filter assumes that the data you feed it is correct, just noisy. Engineers that design radar systems probably spend much of their effort on the problem of incorrect measurements; it is a deep and difficult topic which we will ignore for now.

One final test of your intuition. What if the two measurements are widely separated?

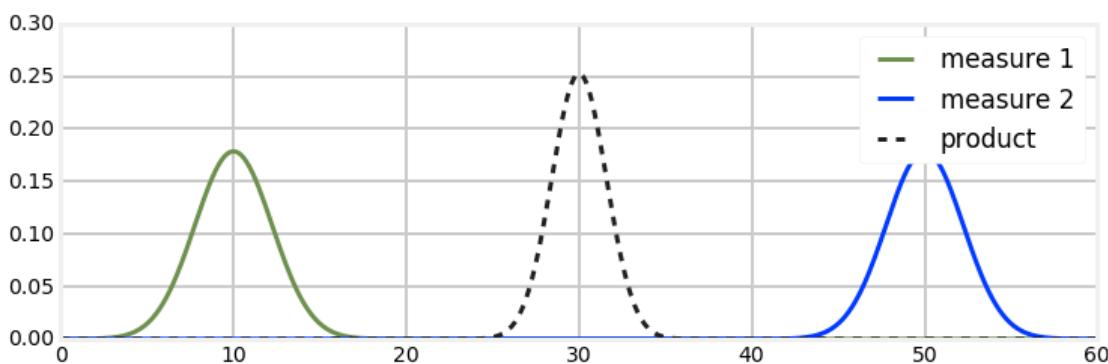
```
In [13]: xs = np.arange(0, 60, 0.1)

mean1, var1 = 10, 5
mean2, var2 = 50, 5
mean, var = multiply(mean1, var1, mean2, var2)

ys = [stats.gaussian(x, mean1, var1) for x in xs]
plt.plot(xs, ys, label='measure 1')

ys = [stats.gaussian(x, mean2, var2) for x in xs]
plt.plot(xs, ys, label='measure 2')

ys = [stats.gaussian(x, mean, var) for x in xs]
plt.plot(xs, ys, label='product', ls='--')
plt.legend();
```



This result bothered me quite a bit when I first learned it. If my first measurement was 10, and the next one was 50, why would I choose 30 as a result? And why would I be more confident? Doesn't it make sense

that either one of the measurements is wrong, or that I am measuring a moving object? Shouldn't the result be nearer 50? And, shouldn't the variance be larger, not smaller?

Well, no. Recall the g-h filter chapter. We agreed that if I weighed myself on two scales, and the first read 160 lbs while the second read 170 lbs, and both were equally accurate, the best estimate was 165 lbs. Furthermore I should be a bit more confident about 165 lbs vs 160 lbs or 170 lbs because I know have two readings, both near this estimate, increasing my confidence that neither is wildly wrong.

Of course, this example is quite exaggerated. The width of the Gaussians is fairly narrow, so this combination of measurements is extraordinarily unlikely. It is hard to eyeball this, but the measurements are well over  $3\sigma$  apart, so the probability of this happening is much less than 1%. Still, it can happen, and the math is correct.

Let's look at the math again to convince ourselves that the physical interpretation of the Gaussian equations makes sense. I'm going to switch back to using priors and measurements. The math and reasoning is the same whether you are using a prior and incorporating a measurement, or just trying to compute the mean and variance of two measurements. For Kalman filters we will be doing a lot more of the former than the latter, so let's get used to it.

$$\mu = \frac{\sigma_{\text{prior}}^2 \mu_z + \sigma_z^2 \mu_{\text{prior}}}{\sigma_{\text{prior}}^2 + \sigma_z^2}$$

If both have the same accuracy, then  $\sigma_{\text{prior}}^2 = \sigma_z^2$ , and the resulting equation is

$$\mu = \frac{\sigma_z^2 (\mu_{\text{prior}} + \mu_z)}{2\sigma_z^2} = \frac{1}{2}(\mu_{\text{prior}} + \mu_z)$$

which is the average of the two means. If we look at the extreme cases, assume the first scale is very much more accurate than than the second one. At the limit, we can set  $\sigma_{\text{prior}}^2 = 0$ , yielding

$$\mu = \frac{0 * \mu_z + \sigma_z^2 \mu_{\text{prior}}}{\sigma_z^2},$$

or

$$\mu = \mu_{\text{prior}}$$

Finally, if we set  $\sigma_{\text{prior}}^2 = 9\sigma_z^2$ , then the resulting equation is

$$\mu = \frac{9\sigma_z^2 \mu_z + \sigma_z^2 \mu_{\text{prior}}}{9\sigma_z^2 + \sigma_z^2}$$

or

$$\mu = \frac{1}{10}\mu_{\text{prior}} + \frac{9}{10}\mu_z$$

This again fits our physical intuition of favoring the second, accurate scale over the first, inaccurate scale.

### 4.3.1 Interactive Example

Rather than going through a dozen examples, lets use interactive code. Here you can use sliders to alter the mean and variance of two Gaussians. As you change the values a plot of the two Gaussians along with the product of the two is displayed.

```
In [14]: from IPython.html.widgets import interact, interactive, fixed
        from IPython.html.widgets import FloatSlider
```

```

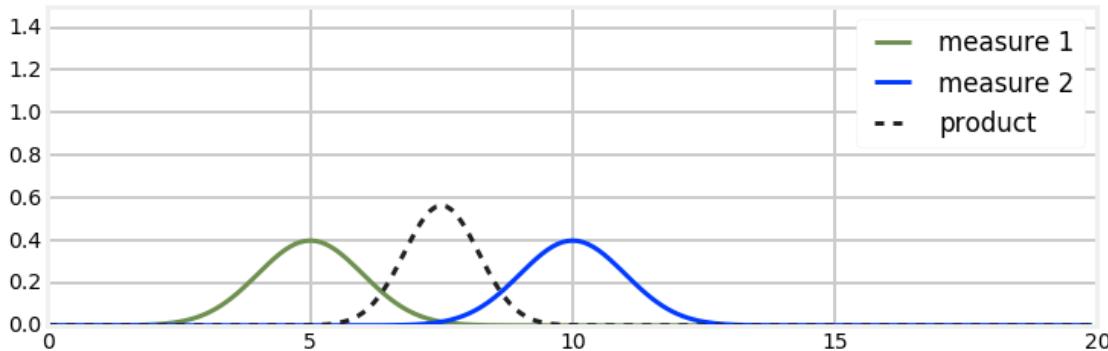
def plot_products(mu_1, mu_2, var_1, var_2):
    xs = np.arange(0, 20, min(var_1, var_2)/10)
    ys = [stats.gaussian(x, mu_1, var_1) for x in xs]
    plt.plot(xs, ys, label='measure 1')

    ys = [stats.gaussian(x, mu_2, var_2) for x in xs]
    plt.plot(xs, ys, label='measure 2')

    mean, var = multiply(mu_1, var_1, mu_2, var_2)
    ys = [stats.gaussian(x, mean, var) for x in xs]
    plt.plot(xs, ys, label='product', ls='--')
    plt.ylim(0, 1.5)
    plt.legend();

interact(plot_products,
         mu_1=FloatSlider(value=5, min=0., max=15),
         mu_2=FloatSlider(value=10, min=0., max=15),
         var_1=FloatSlider(value=1, min=.1, max=5.),
         var_2=FloatSlider(value=1, min=.1, max=5.));

```



## 4.4 Implementing the Update Step

The discrete Bayes filter uses a NumPy array to encode our belief about the position of our dog at any time. That array stored our belief of our dog's position in the hallway using 10 discrete positions. This was very crude, because with a 100 meter hallway that corresponded to positions 10 meters apart. It would have been trivial to expand the number of positions to say 1,000, and that is what we would do if using it for a real problem. But the problem remains that the distribution is discrete and multimodal - it can express strong belief that the dog is in two positions at the same time.

In this chapter we use a Gaussian to reflect our belief of the dog's position. In other words, we will use  $\text{pos}_{\text{dog}} = \mathcal{N}(\mu, \sigma^2)$ . Gaussians extend to infinity on both sides of the mean, so the single Gaussian will cover the entire hallway. They are unimodal, and seem to reflect the behavior of real-world sensors - most errors are small and clustered around the mean. Here is the entire implementation of the update function for a Kalman filter:

```
In [15]: def update(mean, var, measurement, measurement_var):
    return multiply(mean, var, measurement, measurement_var)
```

Kalman filters are supposed to be hard! But this is very short and straightforward. All we are doing is multiplying the Gaussian that reflects our belief of where the dog is with the new measurement.

Perhaps this would be clearer if we used more specific names:

```
def update_dog(dog_pos, dog_var, measurement, measurement_var):
    return multiply(dog_pos, dog_var,
                    measurement, measurement_var)
```

That is less abstract, which perhaps helps with comprehension, but it is poor coding practice. We are writing a Kalman filter that works for any problem, not just tracking dogs in a hallway, so we won't use variable names with 'dog' in them. Still, the `update_dog()` function should make what we are doing very clear.

Let's look at an example. We will suppose that our current belief for the dog's position is  $\mathcal{N}(2, 5)$ . Don't worry about where that number came from as I'll answer that shortly. We will create a `DogSimulation` object initialized to be at position 0.0 with no velocity, and modest noise. This corresponds to the dog standing still at the far left side of the hallway. Note that the belief  $\mathcal{N}(2, 5)$  means that we mistakenly believe the dog is at position 2.0, not 0.0.

```
In [16]: dog = DogSimulation(velocity=0., measurement_var=5, process_var=0.0)
        pos, var = 2, 5
        for i in range(20):
            dog.move()
            pos, var = update(pos, var, dog.sense_position(), 5)
            print('time:', "%2d" % i,
                  '\tposition:', "%.3f" % pos,
                  '\tvariance:', "%.3f" % var)

time: 0      position: -0.033      variance: 2.500
time: 1      position: -0.969      variance: 1.667
time: 2      position: -0.985      variance: 1.250
time: 3      position: -0.167      variance: 1.000
time: 4      position: 0.024      variance: 0.833
time: 5      position: 0.419      variance: 0.714
time: 6      position: 0.190      variance: 0.625
time: 7      position: 0.109      variance: 0.556
time: 8      position: 0.177      variance: 0.500
time: 9      position: -0.024     variance: 0.455
time: 10     position: -0.042     variance: 0.417
time: 11     position: 0.002      variance: 0.385
time: 12     position: -0.048     variance: 0.357
time: 13     position: -0.054     variance: 0.333
time: 14     position: -0.002     variance: 0.312
time: 15     position: -0.042     variance: 0.294
time: 16     position: -0.195     variance: 0.278
time: 17     position: -0.001     variance: 0.263
time: 18     position: -0.020     variance: 0.250
time: 19     position: -0.077     variance: 0.238
```

Because of the random numbers I do not know the exact values that you see, but the position should have converged to near 0 m despite the initial error of believing that the position was 2.0 m. Furthermore, the variance should have converged from the initial value of 5.0 to 0.238.

By now the fact that we converged to a position of 0.0 should not be terribly surprising. All we are doing is computing `new_pos = old_pos * measurement` and the measurement is a normal distribution around 0, so

we should get very close to 0 after 20 iterations. But the truly amazing part of this code is how the variance became 0.238 despite every measurement having a variance of 5.0.

If we think about the physical interpretation of this is should be clear that this is what should happen. If you sent 20 people into the hall with a tape measure to physically measure the position of the dog you would be very confident in the result after 20 measurements - more confident than after 1 or 2 measurements. So it makes sense that as we make more measurements the variance gets smaller.

Mathematically it makes sense as well. Recall the computation for the variance after the multiplication:  $\sigma^2 = 1/(\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2})$ . We take the reciprocals of the sigma from the measurement and prior belief, add them, and take the reciprocal of the result. Think about that for a moment, and you will see that this will always result in smaller numbers as we proceed.

## 4.5 Implementing Predictions

That is a beautiful result, but it is not yet a filter. We assumed that the dog was sitting still, an extremely dubious assumption. The discrete Bayes filter used a loop of predict and update functions, and we must do the same to accommodate movement.

How do we perform the predict function with Gaussians? Recall the discrete Bayes method:

```
def predict(pos, move, p_correct, p_over, p_under):
    n = len(pos)
    result = array(pos, dtype=float)
    for i in range(n):
        result[i] = \
            pos[(i-move) % n] * p_correct + \
            pos[(i-move-1) % n] * p_over + \
            pos[(i-move+1) % n] * p_under
    return result
```

In a nutshell, we shift the probability array by the amount we believe the dog moved, and adjust the probability. How do we do that with Gaussians?

It turns out that we add them. Think of the case without Gaussians. I think my dog is at 7.3 meters, and he moves 2.6 meters to right, where is he now? Obviously,  $7.3 + 2.6 = 9.9$ . He is at 9.9 meters. The algorithm is `new_pos = old_pos + dist_moved`. It does not matter if we use floating point numbers or Gaussians for these values, the algorithm must be the same.

How is addition for Gaussians performed? It turns out to be very simple:

$$\mathcal{N}(\mu_1, \sigma_1^2) + \mathcal{N}(\mu_2, \sigma_2^2) = \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

In other words

$$\mu = \mu_1 + \mu_2 \sigma^2 = \sigma_1^2 + \sigma_2^2$$

We can make this more meaningful by changing the subscripts.

$$\begin{aligned}\mu_{\text{after}} &= \mu_{\text{before}} + \mu_{\text{movement}} \\ \sigma_{\text{after}}^2 &= \sigma_{\text{before}}^2 + \sigma_{\text{movement}}^2\end{aligned}$$

All we do is add the means and the variance separately! Does that make sense? Think of the physical representation of this abstract equation.  $\mu_{\text{before}}$  is the old position, and  $\mu_{\text{movement}}$  is how far we moved. Surely it makes sense that our new position is  $\mu_{\text{before}} + \mu_{\text{movement}}$ .

What about the variance? It is perhaps harder to form an intuition about this. However, recall that with the `predict()` function for the discrete Bayes filter we always lost information - our confidence after the update was lower than our confidence before the update. We don't really know where the dog is moving, so the confidence should get smaller (variance gets larger).  $\sigma_{\text{movement}}^2$  is the amount of uncertainty added to the system due to the imperfect prediction about the movement, and so we would add that to the existing uncertainty.

I assure you that the equation for Gaussian addition is correct, and derived by basic algebra. Therefore it is reasonable to expect that if we are using Gaussians to model physical events, the results must correctly describe those events.

Now is a good time to either work through the algebra to convince yourself of the mathematical correctness of the algorithm, or to work through some examples and see that it behaves reasonably. This book will do the latter.

So, here is our implementation of the `predict` function:

```
In [17]: def predict(pos, variance, movement, movement_variance):
    return (pos + movement, variance + movement_variance)
```

What is left? Just calling these functions. Discrete Bayes did nothing more than loop over the `update()` and `predict()` functions, so let's do the same.

```
In [18]: import kf_internal
```

```
np.random.seed(13)
pos = (0., 400.)    # Gaussian N(0, 400)
velocity = 1.

# variance in process model and the RFID sensor
process_var = 1.
sensor_var = 2.

dog = DogSimulation(pos[0], velocity, sensor_var, process_var)

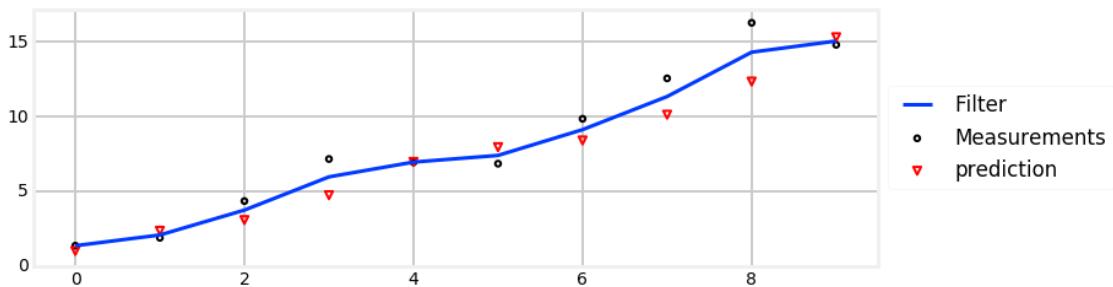
# simulate dog and get measurements
N = 10
zs = [dog.move_and_sense() for _ in range(N)]

positions, predictions = [], []
print('PREDICT:      x      var\t UPDATE:     x      var      z')
for i, z in enumerate(zs):
    pred_pos = predict(pos[0], pos[1], velocity, process_var)
    predictions.append(pred_pos[0])

    pos = update(pred_pos[0], pred_pos[1], z, sensor_var)
    positions.append(pos[0])

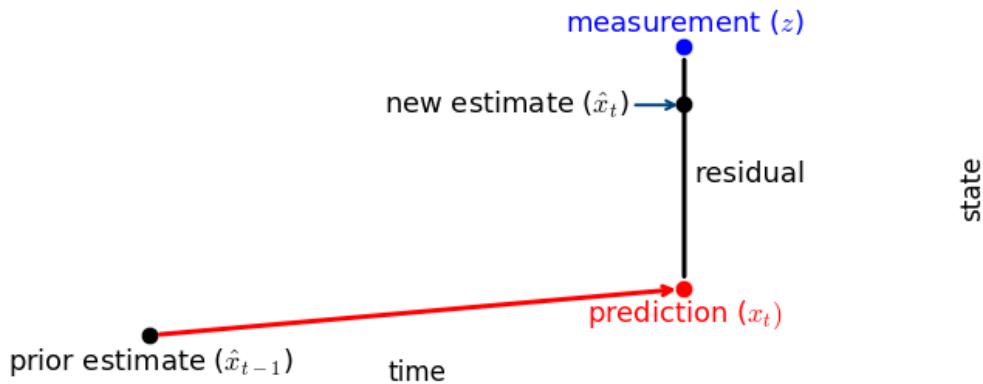
    kf_internal.print_gh(pred_pos, pos, z)
print()
print('final estimate:      {:.10.3f}'.format(pos[0]))
print('actual final position: {:.10.3f}'.format(dog.x))
bp.plot_filter(positions)
bp.plot_measurements(zs)
bp.plot_predictions(predictions)
bp.show_legend();
```

PREDICT:	x	var	UPDATE:	x	var	z
	1.000	401.000		1.352	1.990	1.354
	2.352	2.990		2.070	1.198	1.882
	3.070	2.198		3.736	1.047	4.341
	4.736	2.047		5.960	1.012	7.156
	6.960	2.012		6.949	1.003	6.939
	7.949	2.003		7.396	1.001	6.844
	8.396	2.001		9.122	1.000	9.847
	10.122	2.000		11.338	1.000	12.553
	12.338	2.000		14.305	1.000	16.273
	15.305	2.000		15.053	1.000	14.800
final estimate:						15.053
actual final position:						14.838



Before we discuss the code and results, let's try to understand this better with this chart showing one predict update cycle.

```
In [19]: from mkf_internal import show_residual_chart
show_residual_chart()
```



This chart is very similar to one we already saw in the g-h Filter chapter. At each predict step we take our existing estimate and update it based on our process model. We then make a measurement. The difference between the measurement and the prediction is called the residual. We know neither the measurement or

prediction is perfect, so we choose a value in between the two. This is scaled so that the new estimate is closer to the one we trust more. That scaling depends on the variance of the measurement vs the variance of the process model. We then step time forward one step, and what was the new estimate is considered to be the prior estimate, and the process repeats.

Now let's walk through the code and output.

```
pos = (0., 400.) # Gaussian N(0, 400)
```

Here we set the initial position at 0 m. We set the variance to 400 m<sup>2</sup>, which is a standard deviation of 20 meters. You can think of this as saying "I believe with 99% accuracy the position is 0 plus or minus 60 meters". This is because with Gaussians ~99% of values fall within  $3\sigma$  of the mean.

```
velocity = 1.
```

So how do I know the velocity? Magic? Consider it a prediction, or perhaps we have a secondary velocity sensor. If this is a robot then this would be a control input to the robot. In subsequent chapters we will learn how to handle situations where you don't have a velocity sensor or input, so please accept this simplification for now.

```
process_var = 1.
sensor_var = 2.
```

Here the variance for the predict step is `process_var` and the variance for the sensor is `sensor_var`. The meaning of sensor variance should be clear - it is how much variance there is in each sensor reading. The process variance is how much error there is in the prediction model. We are predicting that at each time step the dog moves forward one meter. Dogs rarely do what we expect, and of course things like hills or the whiff of a squirrel will change his progress. If this was a robot responding to digital commands the performance would be better, but not perfect. Perhaps a hand made robot would have a variance of  $\sigma^2 = .1$ , and an industrial robot might have  $\sigma^2 = 0.02$ . These are not 'magic' numbers; the square root of the express the distance error in meters. It is easy to get a Kalman filter working by just plugging in numbers, but if the numbers do not reflect reality the performance of the filter will be poor.

Next we initialize the simulator and run it with

```
dog = DogSimulation(pos[0], velocity, sensor_var, process_var)
zs = [dog.move_and_sense() for _ in range(N)]
```

It may seem very 'convenient' to set the simulator to the same position as our guess, and it is. Do not fret. In the next example we will see the effect of a wildly inaccurate guess for the dog's initial position. We move the dog for 10 steps forward and store the measurements in `zs` via a list comprehension.

It is traditional to label the measurement  $z$ , and so I follow that convention here. When I was learning this topic I found the standard variable names very obscure. However, if you wish to read the literature you will have to become used to them, so I will not use a more readable variable name such as `m` or `measure`. It is something you have to memorize if you want to work with Kalman filters. As we continue I will introduce more of the standard naming into the code. The `s` at the end of the name is a plural - I'm saying there are many `zs` in the list. This is a common convention with mathematical Python code.

The next line allocates an array to store the output of the filtered positions.

```
positions = np.zeros(N)
```

Now we enter our `predict()` ... `update()` loop.

```

for i in range(N):
    pred_pos = predict(pos[0], pos[1], vel, process_var)
    pos = update(pred_pos[0], pred_pos[1], z, sensor_var)
    positions.append(pos[0])

```

We call the `predict()` function with the Gaussian for the dog's position, and another Gaussian for its movement. `pos[0]`, `pos[1]` is the mean and variance for the position, and 'vel, process\_var' is the mean and variance for the movement.

The first time through the loop I get `pred_pos = (1.0, 401.0)`, as can be seen in the printed table.

What is this saying? After the prediction, we believe that we are at 1.0, and the variance is now 401. Recall we started at 400. The variance got worse, which is always what happens during the prediction step because it involves a loss of information.

Then we call the update function of our filter, using `pred_pos` as the current position and save the result in our `positions` array.

For this I get this as the result: `pos = (1.352, 1.990), z = 1.354`.

What is happening? The dog is actually at 1.0 but the measured position is 1.354 due to the sensor noise. That is pretty far off from the predicted value of 1. The filter incorporates that measurement into its prediction. This is the first time through the loop and so the variance is  $401 \text{ m}^2$ . A large variance implies that the current prediction is very poor, so the filter estimates the position to be 1.352 - very close to the measurement. Intuition tells us that the results should get better as we make more measurements, so let's hope that this is true for our filter as well.

Now look at the variance:  $1.99 \text{ m}^2$ . It has dropped tremendously from  $401 \text{ m}^2$ . Why? Well, the RFID has a reasonably small variance of  $2.0 \text{ m}^2$ , so we trust it far more than our previous belief. However, the previous belief does contain a bit of useful information, so our variance is now slightly smaller than 2.0, which is the variance of the measurement alone.

Now the software loops, calling `predict()` and `update()` in turn. By the end the final estimated position is 15.053 vs the actual position of 14.838. The variance has converged to  $1.0 \text{ m}^2$ .

Now look at the plot. The noisy measurements are plotted in with a dotted red line, and the filter results are in the solid blue line. Both are quite noisy, but notice how much noisier the measurements (red line) are. This is your first Kalman filter and it seems to work!

Before we continue I want to point out how few lines of code actually perform the filtering. Most of the code is either initialization, storing of data, simulating the dog movement, and printing results. The code that performs the filtering is the very succinct.

```

for i in range(N):
    pos = predict(pos[0], pos[1], vel, process_var)
    pos = update(pos[0], pos[1], z, sensor_var)

```

If we didn't use pre-written functions the the code would be:

```

for i in range(N):
    # predict
    pos = pos + vel
    var = var + sensor_var

    # update
    pos = (var*vel + process_var*pos) / (var + process_var)
    var = 1 / (1/var + 1/process_var)

```

Just 4 lines of very simple math implements the entire filter!

In this example I only plotted 10 data points so the output from the print statements would not overwhelm us. Now let's look at the filter's performance with more data. This time we will plot both the output of the filter and the variance. The variance is plotted as a lightly shaded yellow area between dotted lines. I've increased the size of the process and sensor variance so they are easier to see on the chart - for a real Kalman filter of course you will not be randomly changing these values.

```
In [20]: process_var = 2
        sensor_var = 4.5
        pos = (0, 400)    # gaussian N(0, 400)
        N = 25

        dog = DogSimulation(pos[0], velocity, sensor_var, process_var)
        zs = [dog.move_and_sense() for _ in range(N)]

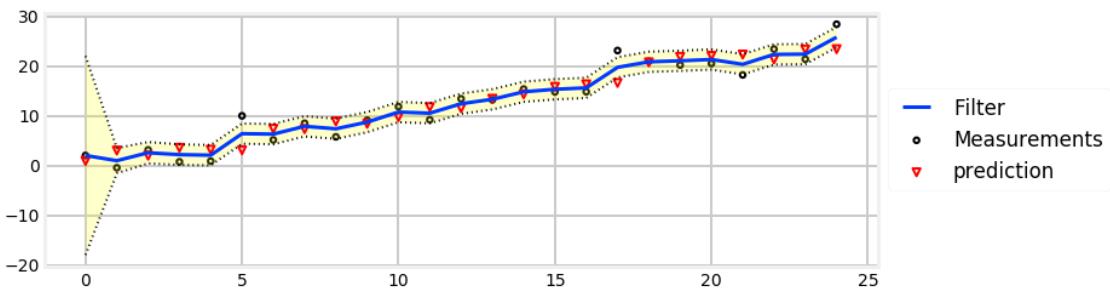
        positions, predictions = np.zeros((N, 2)), np.zeros((N, 2))
        for i, z in enumerate(zs):
            pos = predict(pos[0], pos[1], velocity, process_var)
            predictions[i] = pos

            pos = update(pos[0], pos[1], z, sensor_var)
            positions[i] = pos

        bp.plot_measurements(zs)
        bp.plot_filter(positions[:, 0], var=predictions[:, 1])
        bp.plot_predictions(predictions[:, 0])
        bp.show_legend()
        kf_internal.print_variance(positions)
```

Variance:

```
4.4502 2.6507 2.2871 2.1955 2.1712
2.1647 2.1629 2.1625 2.1623 2.1623
2.1623 2.1623 2.1623 2.1623 2.1623
2.1623 2.1623 2.1623 2.1623 2.1623
2.1623 2.1623 2.1623 2.1623 2.1623
```



Here we can see that the variance converges very quickly to roughly 2.1623 in 10 steps. We interpret this as meaning that we become very confident in our position estimate very quickly. The first few measurements are unsure due to our uncertainty in our guess at the initial position, but the filter is able to quickly determine an accurate estimate.

Before I go on, I want to emphasize that this code fully implements a one dimensional Kalman filter. “One dimensional” just refers to the fact we are tracking one variable. If you have tried to

read the literature you are perhaps surprised, because this looks nothing like the endless pages of math in those books. To be fair, the math gets more complicated in multiple dimensions, but not by much. So long as we worry about using the equations rather than deriving them we can create Kalman filters without a lot of effort. Moreover, I hope you'll agree that you have a decent intuitive grasp of what is happening. We represent our beliefs with Gaussians, and our beliefs get better over time because more measurements means we have more data to work with.

### 4.5.1 Exercise: Modify Variance Values

Modify the values of `process_var` and `sensor_var` and note the effect on the filter and on the variance. Which has a larger effect on the variance convergence?. For example, which results in a smaller variance:

```
process_var = 40
sensor_var = 2
```

or:

```
process_var = 2
sensor_var = 40
```

### 4.5.2 Tracking Animation

If you are reading this in a browser you will be able to see an animation of the filter tracking the dog directly below this sentence.

The top plot shows the output of the filter in green, and the measurements with a dashed red line. The bottom plot shows the Gaussian at each step.

When the track first starts you can see that the measurements varies quite a bit from the initial prediction. At this point the Gaussian probability is small (the curve is low and wide) so the filter does not trust its prediction. As a result, the filter adjusts its estimate a large amount. As the filter innovates you can see that as the Gaussian becomes taller, indicating greater certainty in the estimate, the filter's output becomes very close to a straight line. At  $x = 15$  and greater you can see that there is a large amount of noise in the measurement, but the filter does not react much to it compared to how much it changed for the first noisy measurement.

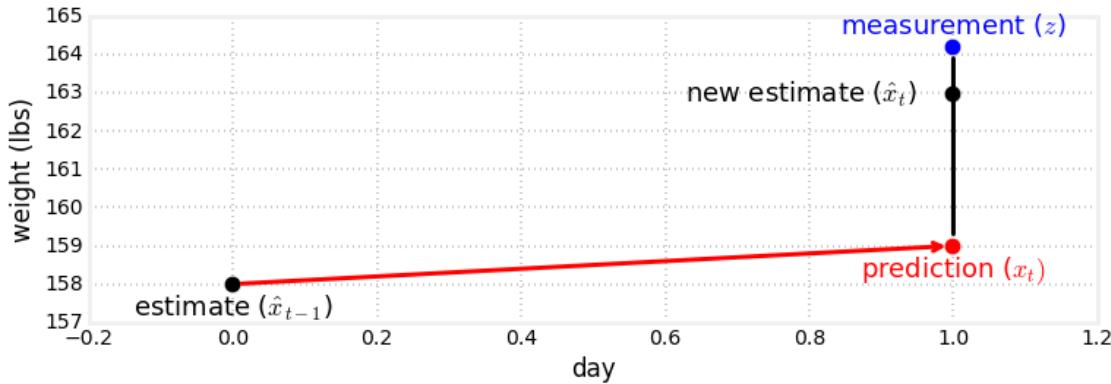
## 4.6 The Kalman Filter is a g-h Filter

Earlier we set  $\sigma_{\text{prior}}^2 = 9\sigma_z^2$ , to get

$$\mu = \frac{1}{10}\mu_{\text{prior}} + \frac{9}{10}\mu_z$$

We can see that we are scaling the prior and the measurement by some factor. Think back to the **g-h Filter** chapter and this plot:

```
In [21]: import gh_internal
gh_internal.plot_estimate_chart_3()
```



The Kalman filter equation is performing this computation! We have a measurement and a prediction, and the new estimate is chosen as some value between the two. So, as promised, the Kalman filter is a g-h filter. I have not made it explicit in this chapter, but the scaling factor is called the **Kalman gain**. In this equation

$$\mu = \frac{1}{10}\mu_{\text{prior}} + \frac{9}{10}\mu_z$$

if set the Kalman gain  $K = 0.9$  we can rewrite the equation as

$$\mu = (1 - K)\mu_{\text{prior}} + K\mu_z$$

That is a general equation for the Kalman filter in the one dimensional case.

How does this relate to the g-h filter? In the g-h filter we chose the constants  $g$  and  $h$  and didn't vary them. Here the scaling factors are chosen dynamically based on the current variance of the filter and the variance of the measurement. in other words the Kalman gain is determined by a ratio of the variances of the prior and measurement. As the Kalman filter innovates the scaling factors change as we become more (or less) certain about our answer.

Some algebra will reveal the relationship between the Kalman filter and the g-h filter. If you are curious, I present it in **Kalman Filter Math** chapter.

At each time step we can show that  $g$  equals

$$g = \frac{\sigma_{\text{prior}}^2}{\sigma_z^2}$$

This make intuitive sense. When we perform an update we choose a value somewhere between the measurement and prediction. If the measurement is more accurate we will choose a value nearer it, and vice versa. This computation for  $g$  is a ratio of the errors in the prediction and the measurement.

The equation for  $h$  is similar.

$$h = \frac{COV(x, \dot{x})}{\sigma_z^2}$$

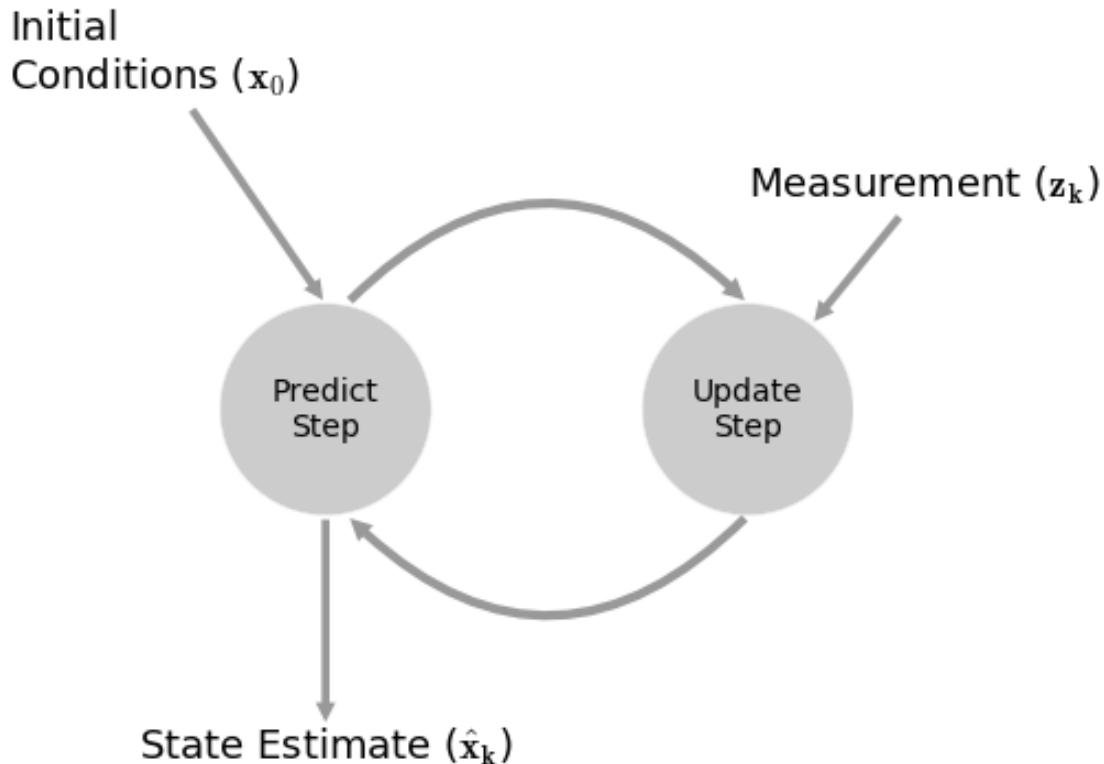
We haven't discussed covariance (denoted by  $COV$ ), so you may not fully understand this equation. Think of the numerator as the variance between the position and its velocity. Recall that  $g$  is used to scale the state, and  $h$  is used to scale the velocity of the state, and this equation should seem reasonable. In the next chapter it will become very clear.

The takeaway point is that  $g$  and  $h$  are specified fully by the variance and covariances of the measurement and predictions at time  $n$ . That is all the Kalman filter is. The Kalman filter is a g-h filter where  $g$  and  $h$  vary over time according to the variance of the measurements and our current beliefs.

## 4.7 Full Description of the Algorithm

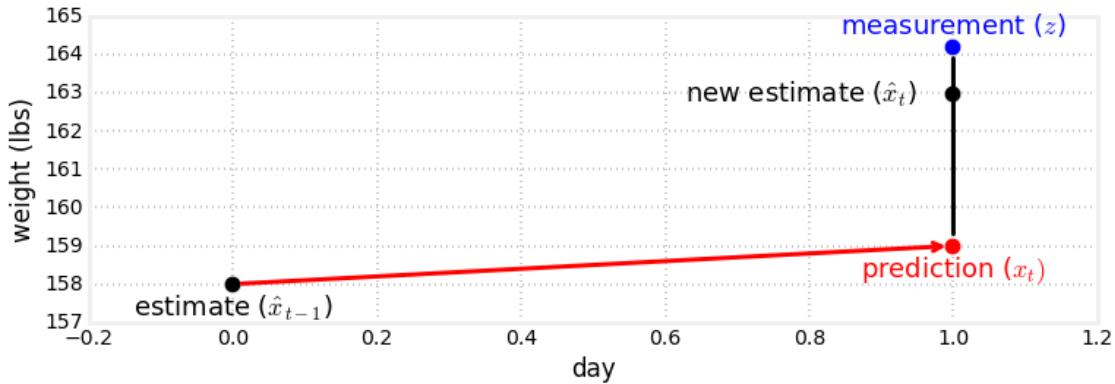
Recall the diagram we used for the g-h filter:

```
In [22]: import gh_internal
gh_internal.create_predict_update_chart()
```



We've been doing the same thing in this chapter. And, as already mentioned, each filter makes a prediction, takes a measurement, and then forms a new estimate somewhere between the two.

```
In [23]: gh_internal.plot_estimate_chart_3()
```



**This is extremely important to understand:** Every filter in this book implements the **same algorithm**, just with different mathematical details. The math can become challenging in later chapters, but the idea is easy to understand.

It is important to see past the details of the equations of a specific filter and understand what the equations are calculating and why. This is important not only so you can understand the filters in this book, but so that you can understand filters that you will see in other situations. There are a tremendous number of filters that we do not describe - the SVD filter, the unscented particle filter, the  $H_\infty$  filter, the ensemble filter, and so many more. They all use different math to implement the same algorithm. The choice of math affects the quality of results and what problems can be represented, but not the underlying ideas.

Here is the generic algorithm:

#### Initialization

1. Initialize the state of the filter
2. Initialize our belief in the state

#### Predict

1. Use system behavior to predict state at the next time step
2. Adjust belief to account for the uncertainty in prediction

#### Update

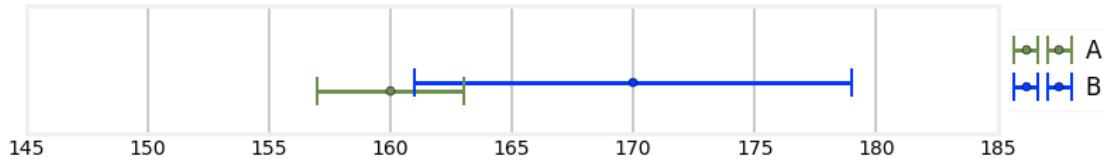
1. Get a measurement and associated belief about its accuracy
2. Compute residual between estimated state and measurement
3. Compute scaling factor based on whether the measurement or prediction is more accurate
4. set state between the prediction and measurement based on scaling factor
5. update belief in the state based on how certain we are in the measurement

You will be hard pressed to find a Bayesian filter algorithm that does not fit into this form. Some filters will not include some aspect, such as error in the prediction, and others will have very complicated methods of computation, but this is what they all do. If you go back to the **g-h Filter** chapter and reread it you'll recognize that the very simple thought experiments we did there result in this algorithm.

## 4.8 Comparison with g-h Filter and discrete Bayes Filter

Now is a good time to be understand the differences between these three filters in terms of how we model errors. For the g-h filter we modeled our measurements as shown in this graph

```
In [24]: import gh_internal as gh
gh.plot_errorbar2()
```



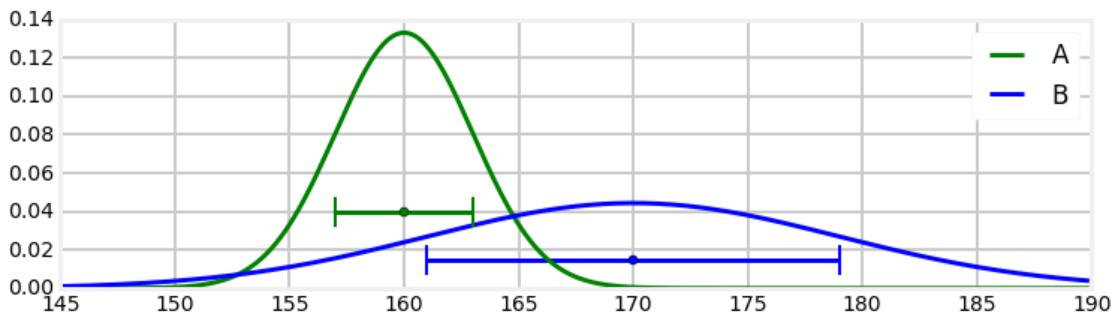
In this graph sensor A returned a measurement of 160, and sensor B returned 170. The bars are error bars - they illustrate the possible range of error for the measurement. Hence, the actual value that A is measuring can be between 157 to 163, and B is measuring a value between 161 to 179.

I did not define it at the time, but this is a uniform distribution. A uniform distribution assigns equal probability to any event in the range. According to this model it is equally likely for sensor A to read 157, 160, or 163. Likewise, 161, 170, and 179 are equally likely. Any value outside these ranges have 0 probability.

We can model this situation with Gaussians. I'll use  $\mathcal{N}(160, 3^2)$  for sensor A, and  $\mathcal{N}(170, 9^2)$  for sensor B. I've plotted these below with the uniform distribution error bars for comparison.

```
In [25]: import book_format
with book_format figsize(y=3):
    xs = np.arange(145, 190, 0.1)
    ys = [stats.gaussian(x, 160, 3**2) for x in xs]
    plt.plot(xs, ys, label='A', color='g')

    ys = [stats.gaussian(x, 170, 9**2) for x in xs]
    plt.plot(xs, ys, label='B', color='b')
    plt.legend();
    plt.errorbar(160, [0.04], xerr=[3], fmt='o', color='g', capthick=2, capsiz
    e=10)
    plt.errorbar(170, [0.015], xerr=[9], fmt='o', color='b', capthick=2, capsiz
    e=10)
```

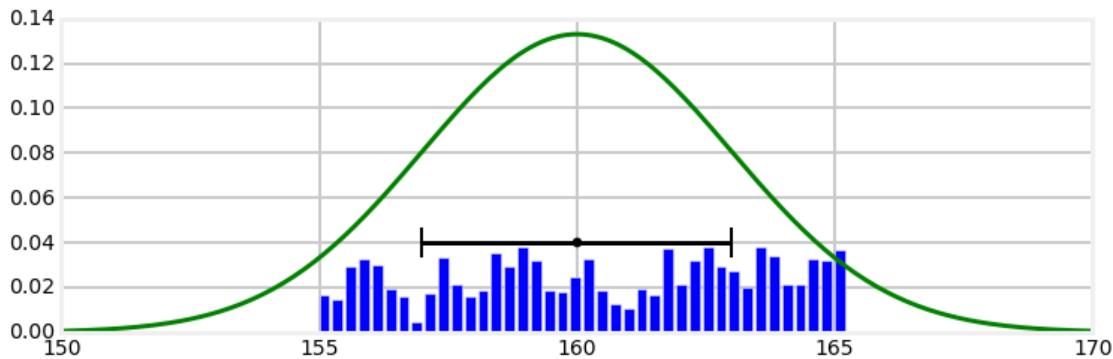


Using a uniform or Gaussian distribution is a modeling choice. Neither exactly describes reality. In most cases the Gaussian distribution is more realistic. Most sensors are more likely to return readings near the value being measured, and unlikely to return a reading far from that value. The Gaussian models this tendency. In contrast the uniform distribution assumes that any measurement within a range is equally likely.

Now let's see the discrete distribution used in the discrete Bayes filter. This model divides the range of possible values into discrete ranges and assigns a probability to each bucket. This assignment can be entirely arbitrary so long as the probabilities sum to one.

Let's plot the data for one sensor using a uniform distribution, a Gaussian distribution, and a discrete distribution.

```
In [26]: from random import random
xs = np.arange(145, 190, 0.1)
ys = [stats.gaussian(x, 160, 3**2) for x in xs]
belief = np.array([random() for _ in range(40)])
belief = belief / sum(belief)
x = np.linspace(155, 165, len(belief))
plt.gca().bar(x, belief, width=0.2)
plt.plot(xs, ys, label='A', color='g')
plt.errorbar(160, [0.04], xerr=[3], fmt='o', color='k', capthick=2, capsiz=10)
plt.xlim(150, 170);
```



I used random numbers to form the discrete distribution to illustrate that it can model any arbitrary probability distribution. This provides it with enormous power. With enough discrete buckets we can model the error characteristics of any sensor no matter how complicated. But with this power comes mathematical intractability. Multiplying or adding Gaussians takes two lines of math. Multiplying or adding a discrete distribution requires looping over the data, and we have no easy way to characterize the result. In contrast, multiplying two Gaussians results in another Gaussian, as does adding them. This regularity allows us to perform powerful analysis on the performance and behavior of our filters. Analyzing the performance characteristics of a filter based on a discrete distribution is extremely difficult to impossible.

There is no ‘correct’ choice here. Later in the book we will introduce the particle filter which uses a discrete distribution. It is an extremely powerful technique because it can handle arbitrarily complex situations. This comes at the cost of very slow performance, and resistance to analytical analysis.

For now we will ignore these matters and return to using Gaussians for the next several chapters. As we progress you will learn the strengths and limitations of using Gaussians in our mathematical models.

## 4.9 Implementation in a Class

For many purposes the code above suffices. However, if you write enough of these filters the functions will become a bit annoying. For example, having to write

```
pos = predict(pos[0], pos[1], movement, movement_variance)
```

is a bit cumbersome and error prone. Let's investigate how we might implement this in a form that makes our lives easier.

First, values for the process error and the measurement errors are often constant for a given problem, so we only want to specify them once. We can store them in instance variables in the class. Second, it is annoying to have to pass in the state (pos in the code snippet above) and then remember to assign the output of the function back to that state, so the state should also be an instance variable. Our first attempt might look like:

```
class KalmanFilter1D:
    def __init__(self, state, state_var, measurement_var, movement_var):
        self.state = state
        self.state_var = state_var
        self.measurement_var = measurement_var
        self.movement_var = movement_var
```

That works, but I am going to use different naming. The Kalman filter literature uses math, not code. It uses single letter variables and I will start exposing you to them now. At first it seems impossibly terse, but as you become familiar with the notation you'll see that the math formulas in the textbooks will have an exact one-to-one correspondence with the code. Unfortunately there is not a lot of meaning behind the notation unless you have experience in the history of control theory; you will have to memorize them. If you do not make this effort you will never be able to read the literature or code written by others.

We will use  $x$  for the state (estimated value of the filter) and  $P$  for the variance of the state.  $R$  is the measurement error, and  $Q$  is the process (prediction) error. Almost every publication uses this notation, so learn it. This gives us:

```
class KalmanFilter1D:
    def __init__(self, x0, P, R, Q):
        self.x = x0
        self.P = P
        self.R = R
        self.Q = Q
```

Now we can implement the `update()` and `predict()` methods. In the literature the measurement is usually named either  $z$  or  $y$ ; I find  $y$  is too easy to confuse with the  $y$  axis of a plot, so I like  $z$ . I like to think I can hear a  $z$  in measurement, which helps me remember what  $z$  stands for. So for the update method we might write:

```
def update(z):
    self.x = (self.P*z + self.x*self.R) / (self.P + self.R)
    self.P = 1 / (1/self.P + 1/self.R)
```

Finally, the movement is usually called  $u$ , and so we will use that. So for the predict method we might write:

```
def predict(self, u):
    self.x += u
    self.P += self.Q
```

That give us the following code. Production code would require significant comments and error checking. However, in subsequent chapters we will develop Kalman filter code that works for any dimension. This class will never be more than a stepping stone for us, so I'll keep things simple.

```
In [27]: class KalmanFilter1D:
    def __init__(self, x0, P, R, Q):
        self.x = x0
        self.P = P
        self.R = R
        self.Q = Q

    def update(self, z):
        self.x = (self.P*z + self.x*self.R) / (self.P + self.R)
        self.P = 1. / (1./self.P + 1./self.R)

    def predict(self, u=0.0):
        self.x += u
        self.P += self.Q
```

## 4.10 Introduction to Designing a Filter

So far we have developed filters for a position sensor. We are used to this problem by now, and may feel ill-equipped to implement a Kalman filter for a different problem. To be honest, there is still quite a bit of information missing from this presentation. Following chapters will fill in the gaps. Still, lets get a feel for it by designing and implementing a Kalman filter for a thermometer. The sensor for the thermometer outputs a voltage that corresponds to the temperature that is being measured. We have read the manufacturer's specifications for the sensor, and it tells us that the sensor exhibits white noise with a standard deviation of 0.13 volts.

We can simulate the temperature sensor measurement with this function:

```
In [28]: def volt(voltage, std):
    return voltage + (randn() * std)
```

Now we need to write the Kalman filter processing loop. As with our previous problem, we need to perform a cycle of predicting and updating. The sensing step probably seems clear - call `volt()` to get the measurement, pass the result into `update()` method, but what about the predict step? We do not have a sensor to detect 'movement' in the voltage, and for any small duration we expect the voltage to remain constant. How shall we handle this?

As always, we will trust in the math. We have no known movement, so we will set that to zero. However, that means that we are predicting that the temperature will never change. If that is true, then over time we should become extremely confident in our results. Once the filter has enough measurements it will become very confident that it can predict the subsequent temperatures, and this will lead it to ignoring measurements that result due to an actual temperature change. This is called a smug filter, and is something you want to avoid. So we will add a bit of error to our prediction step to tell the filter not to discount changes in voltage over time. In the code below I set `process_var = .05**2`. This is the expected variance in the change of voltage over each time step. I chose this value merely to be able to show how the variance changes through the update and predict steps. For an real sensor you would set this value for the actual amount of change you expect. For example, this would be an extremely small number if it is a thermometer for ambient air temperature in a house, and a high number if this is a thermocouple in a chemical reaction chamber. We will say more about selecting the actual value in the later chapters.

Let's see what happens.

```
In [29]: temp_change = 0
voltage_std = .13
process_var = .05**2
actual_voltage = 16.3

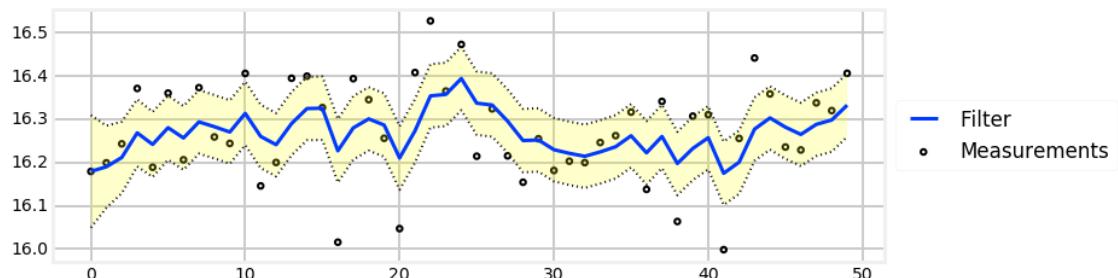
N = 50
zs = [volt(actual_voltage, voltage_std) for i in range(N)]
ps = []
estimates = []

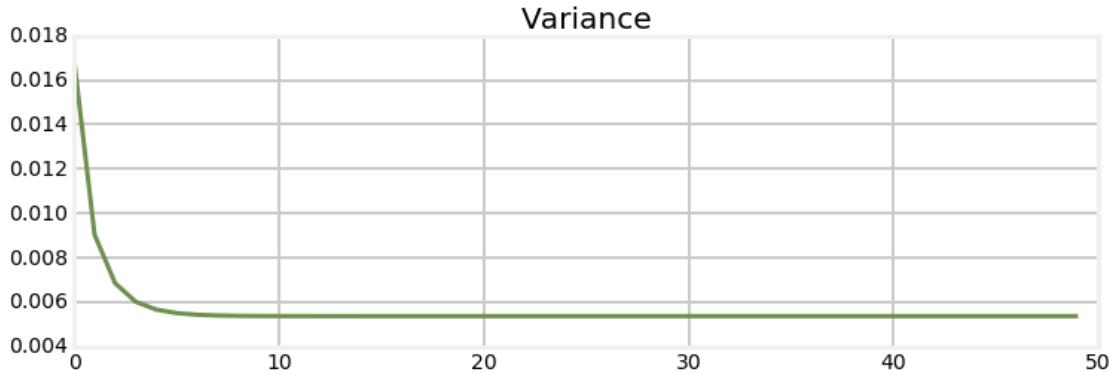
kf = KalmanFilter1D(x0=25,          # initial state
                     P=1000,        # initial variance
                     R=voltage_std**2, # sensor noise
                     Q=process_var) # error in prediction

for i in range(N):
    kf.predict(temp_change)
    kf.update(zs[i])

    # save for latter plotting
    estimates.append(kf.x)
    ps.append(kf.P)

# plot the filter output and the variance
bp.plot_measurements(zs)
bp.plot_filter(estimates, var=np.array(ps))
bp.show_legend()
plt.show()
plt.plot(ps)
plt.title('Variance')
plt.show()
print('Variance converges to {:.3f}'.format(ps[-1]))
```





Variance converges to 0.005

The first plot shows the individual sensor measurements vs the filter output. Despite a lot of noise in the sensor we quickly discover the approximate voltage of the sensor. In the run I just completed at the time of authorship, the last voltage output from the filter is 16.213, which is quite close to the 16.4 used by the `volt()` function. On other runs I have gotten larger and smaller results.

Spec sheets are what they sound like - specifications. Any individual sensor will exhibit different performance based on normal manufacturing variations. Values are often maximums - the spec is a guarantee that the performance will be at least that good. If you buy an expensive piece of equipment it often comes with a sheet of paper displaying the test results of your specific item; this is usually very trustworthy. On the other hand, if this is a cheap sensor it is likely it received little to no testing prior to being sold. Manufacturers typically test a small subset of their output to verify that a sample falls within the desired performance range. If you have a critical application you will need to read the specification sheet carefully to figure out exactly what they mean by their ranges. Do they guarantee their number is a maximum, or is it, say, the  $3\sigma$  error rate? Is every item tested? Is the variance normal, or some other distribution? Finally, manufacturing is not perfect. Your part might be defective and not match the performance on the sheet.

For example, I am looking at a data sheet for an airflow sensor. There is a field Repeatability, with the value  $\pm 0.50\%$ . Is this a Gaussian? Is there a bias? For example, perhaps the repeatability is nearly 0.0% at low temperatures, and always nearly +0.50 at high temperatures. Data sheets for electrical components often contain a section of “Typical Performance Characteristics”. These are used to capture information that cannot be easily conveyed in a table. For example, I am looking at a chart showing output voltage vs current for a LM555 timer. There are three curves showing the performance at different temperatures. The response is ideally linear, but all three lines are curved. This clarifies that errors in voltage outputs are probably not Gaussian - in this chip’s case higher temperatures leads to lower voltage output, and the voltage output is quite nonlinear if the input current is very high.

As you might guess, modeling the performance of your sensors is one of the harder parts of creating a Kalman filter that performs well.

#### 4.10.1 Animation

For those reading this in a browser here is an animation showing the filter working. If you are not using a browser you can see this plot at [https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python/blob/master/animations/05\\_volt\\_animate.gif](https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python/blob/master/animations/05_volt_animate.gif).

The top plot in the animation draws a green line for the predicted next voltage, then a red '+' for the actual measurement, draws a light red line to show the residual, and then draws a blue line to the filter’s output.

You can see that when the filter starts the corrections made are quite large, but after only a few updates the filter only adjusts its output by a small amount even when the measurement is far from it.

The lower plot shows the Gaussian belief as the filter innovates. When the filter starts the Gaussian curve is centered over 25, our initial guess for the voltage, and is very wide and short due to our initial uncertainty. But as the filter innovates, the Gaussian quickly moves to about 16.0 and becomes taller, reflecting the growing confidence that the filter has in its estimate for the voltage. You will also note that the Gaussian's height bounces up and down a little bit. If you watch closely you will see that the Gaussian becomes a bit shorter and more spread out during the prediction step, and becomes taller and narrower as the filter incorporates another measurement.

Think of this animation in terms of the g-h filter. At each step the g-h filter makes a prediction, takes a measurement, computes the residual (the difference between the prediction and the measurement), and then selects a point on the residual line based on the scaling factor  $g$ . The Kalman filter is doing exactly the same thing, except that the scaling factor  $g$  varies with time. As the filter becomes more confident in its state the scaling factor favors the filter's prediction over the measurement.

## 4.11 Exercise: Plot Histogram of Measurements

Write a function that runs the Kalman filter many times and record what value the voltage converges to each time. Plot this as a histogram. Use `plt.hist(data, bins=100)` to plot the histogram. After 10,000 runs do the results look normally distributed? Does this match your intuition of what should happen?

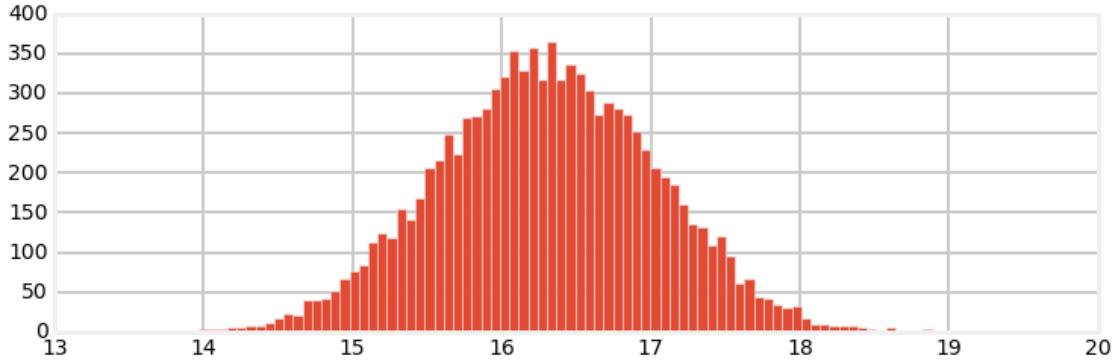
In [30]: #Your code here

### 4.11.1 Solution

```
In [31]: variance = 2.13**2
process_var = 0.1**2
actual_voltage = 16.3

def VKF():
    N = 50
    voltage = (14, 1000)
    for i in range(N):
        Z = volt(actual_voltage, variance)
        voltage = predict(voltage[0], voltage[1], 0, process_var)
        voltage = update(voltage[0], voltage[1], Z, variance)
    return voltage[0]

vs = []
for i in range(10000):
    vs.append(VKF())
plt.hist(vs, bins=100, color='#e24a33');
```



#### 4.11.2 Discussion

The results do in fact look like a normal distribution. Each voltage is Gaussian, and the Central Limit Theorem guarantees that a large number of Gaussians is normally distributed.

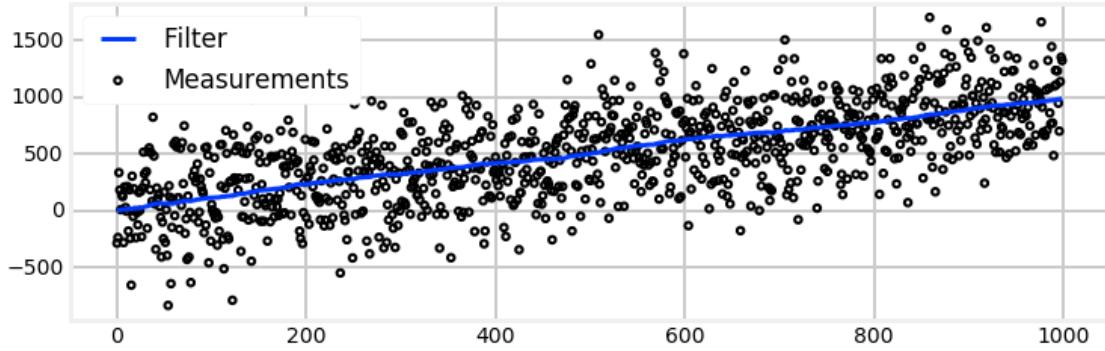
### 4.12 Example: Extreme Amounts of Noise

With the dog filter I didn't put a lot of noise in the signal, and I also 'correctly guessed' that the dog was at position 0. How does the filter perform in real world conditions? I will start by injecting more noise in the RFID sensor while leaving the process variance at  $2 \text{ m}^2$ . I will inject an extreme amount of noise - noise that apparently swamps the actual measurement. What does your intuition say about the filter's performance if the has a standard deviation of 300 meters?. In other words, an actual position of 1.0 m might be reported as 287.9 m, or -589.6 m, or any other number in roughly that range. Think about it before you scroll down.

```
In [32]: velocity = 1
sensor_var = 300**2
process_var = 2
pos = (0,500)
N = 1000
dog = DogSimulation(pos[0], velocity, sensor_var, process_var)
zs = [dog.move_and_sense() for _ in range(N)]
ps = []

for i in range(N):
    pos = predict(pos[0], pos[1], velocity, process_var)
    pos = update(pos[0], pos[1], zs[i], sensor_var)
    ps.append(pos[0])

bp.plot_measurements(zs, lw=1)
bp.plot_filter(ps)
plt.legend(loc='best');
```



In this example the noise is extreme yet the filter still outputs a nearly straight line! This is an astonishing result! What do you think might be the cause of this performance?

We get a nearly straight line because our process error is small. A small process error tells the filter that the prediction is very trustworthy, and the prediction is a straight line, so the filter outputs a straight line.

## 4.13 Example: Incorrect Process Variance

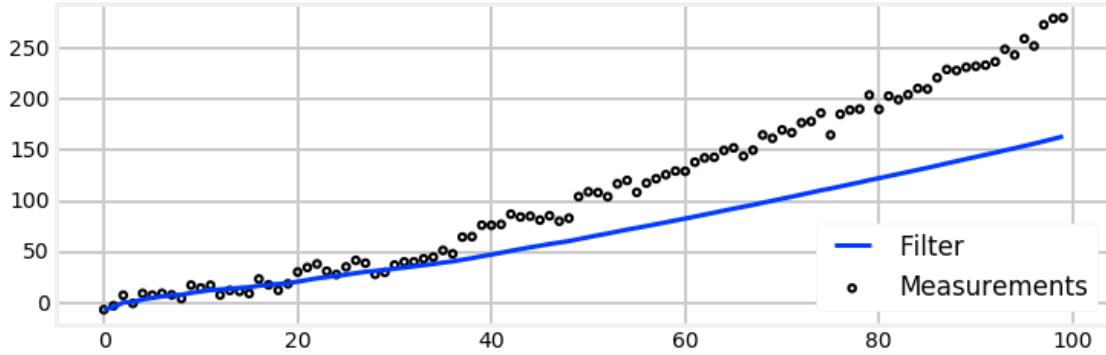
That last filter looks fantastic! Why wouldn't we set the process variance very low, as it guarantees the result will be straight and smooth?

The process variance tells the filter how much the system is changing over time. If you lie to the filter by setting this number artificially low the filter will not be able to react to changes that are happening. Let's have the dog increase his velocity by a small amount at each time step and see how the filter performs with a process variance of  $0.001 \text{ m}^2$ .

```
In [33]: velocity = 1
sensor_var = 20
process_var = .001
pos = (0,500)
N = 100
dog = DogSimulation(pos[0], velocity, sensor_var, process_var*10000)
zs = []
for _ in range(N):
    dog.velocity += 0.04
    zs.append(dog.move_and_sense())
ps = []

for i in range(N):
    pos = predict(pos[0], pos[1], velocity, process_var)
    pos = update(pos[0], pos[1], zs[i], sensor_var)
    ps.append(pos[0])

bp.plot_measurements(zs, lw=1)
bp.plot_filter(ps)
plt.legend(loc=4);
```



It is easy to see that the filter is not correctly responding to the measurements. The measurements clearly indicate that the dog is changing speed but the filter has been told that it's predictions are nearly perfect so it almost entirely ignores the measurements. I encourage you to adjust the amount of movement in the dog vs process variance. We will also be studying this topic much more in the later chapters. The key point is the filter will do what you tell it to. There is nothing in the math that allows the filter to recognize that you have lied to it about the variance in the measurements and process model.

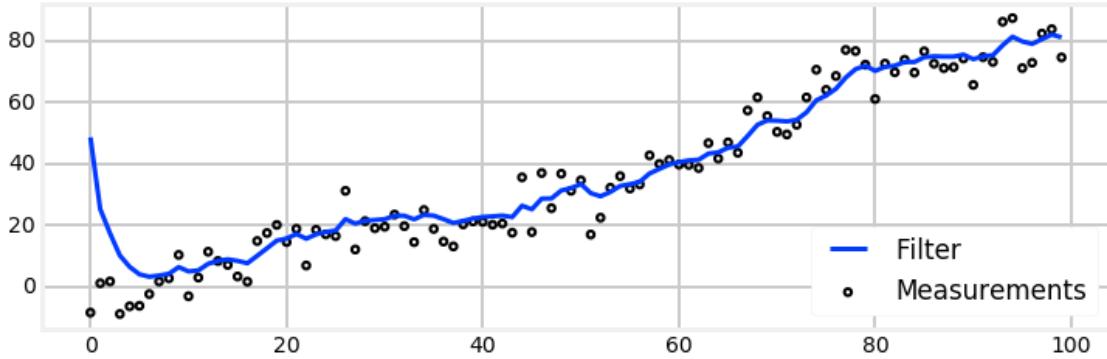
## 4.14 Example: Bad Initial Estimate

Now let's lets look at the results when we make a bad initial estimate of position. To avoid obscuring the results I'll reduce the sensor variance to 30, but set the initial position to 1000 meters. Can the filter recover from a 1000 meter error?

```
In [34]: velocity = 1
sensor_var = 30
process_var = 2
pos = (1000, 500)
N = 100
dog = DogSimulation(0, velocity, sensor_var, process_var)
zs = [dog.move_and_sense() for _ in range(N)]
ps = []

for i in range(N):
    pos = predict(pos[0], pos[1], velocity, process_var)
    pos = update(pos[0], pos[1], zs[i], sensor_var)
    ps.append(pos[0])

bp.plot_measurements(zs, lw=1)
bp.plot_filter(ps)
plt.legend(loc=4);
```



Again the answer is yes! Because we are relatively sure about our belief in the sensor ( $\sigma^2 = 30$ ) after only the first step we have changed our position estimate from 1000 m to roughly 60 m. After another 5-10 measurements we have converged to the correct value. This is how we get around the chicken and egg problem of initial guesses. In practice we would likely assign the first measurement from the sensor as the initial value, but you can see it doesn't matter much if we wildly guess at the initial conditions - the Kalman filter still converges so long as the filter variances are chosen to match the actual process and measurement variances.

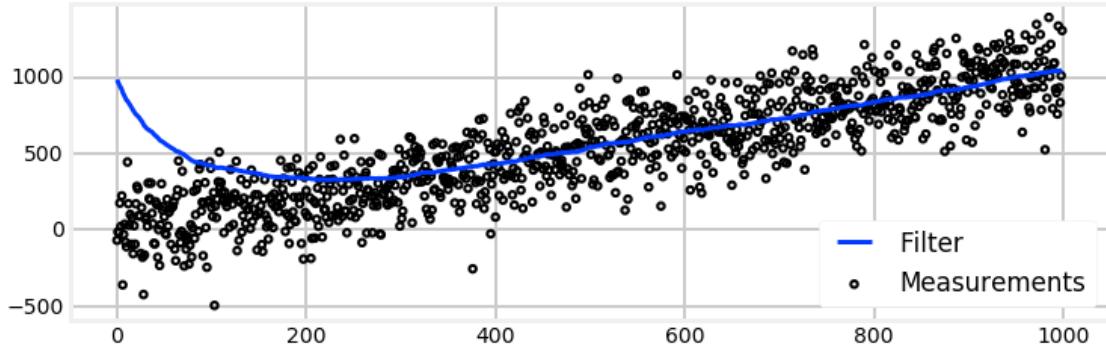
## 4.15 Example: Large Noise and Bad Initial Estimate

What about the worst of both worlds, large noise and a bad initial estimate?

```
In [35]: sensor_var = 30000
process_var = 2
pos = (1000, 500)
N = 1000
dog = DogSimulation(0, velocity, sensor_var, process_var)
zs = [dog.move_and_sense() for _ in range(N)]
ps = []

for i in range(N):
    pos = predict(pos[0], pos[1], velocity, process_var)
    pos = update(pos[0], pos[1], zs[i], sensor_var)
    ps.append(pos[0])

bp.plot_measurements(zs, lw=1)
bp.plot_filter(ps)
plt.legend(loc=4);
```



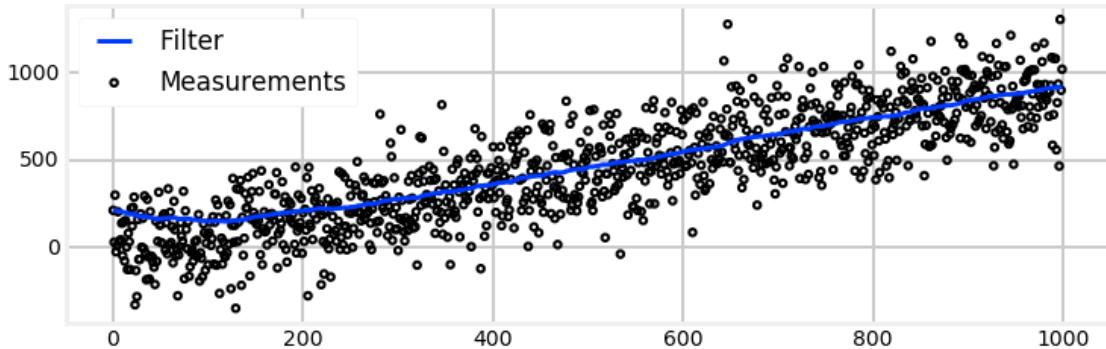
This time the filter struggles. Notice that the previous example only computed 100 updates, whereas this example uses 1000. By my eye it takes the filter 400 or so iterations to become reasonable accurate, but maybe over 600 before the results are good. Kalman filters are good, but we cannot expect miracles. If we have extremely noisy data and extremely bad initial conditions, this is as good as it gets.

Finally, let's implement the suggestion of using the first measurement as the initial position.

```
In [36]: sensor_var = 30000
process_var = 2
N = 1000
dog = DogSimulation(0, velocity, sensor_var, process_var)
zs = [dog.move_and_sense() for _ in range(N)]

pos = (zs[0], 500)
ps = []
for i in range(N):
    pos = predict(pos[0], pos[1], velocity, process_var)
    pos = update(pos[0], pos[1], zs[i], sensor_var)
    ps.append(pos[0])

bp.plot_measurements(zs, lw=1)
bp.plot_filter(ps)
plt.legend(loc='best');
```



This simple change significantly improves the results. On some runs it takes 200 iterations or so to settle to a good solution, but other runs it converges very rapidly. This all depends on the amount of noise in the first measurement. A large amount of noise causes the initial estimate to be far from the dog's position.

200 iterations may seem like a lot, but the amount of noise we are injecting is truly huge. In the real world we use sensors like thermometers, laser range finders, GPS satellites, computer vision, and so on. None have the enormous errors in these examples. A reasonable variance for a cheap thermometer might be  $0.2 \text{ C}^{\circ 2}$ , and our code is using  $30,000 \text{ C}^{\circ 2}$ .

## 4.16 Exercise: Interactive Plots

Implement the Kalman filter using Jupyter Notebook's animation features to allow you to modify the various constants in real time using sliders. Refer to the section **Interactive Gaussians** in the **Gaussians** chapter to see how to do this. You will use the `interact()` function to call a calculation and plotting function. Each parameter passed into `interact()` automatically gets a slider created for it. I have written the boilerplate for this; you fill in the required code.

```
In [37]: from IPython.html.widgets import interact, interactive, fixed
         from IPython.html.widgets import FloatSlider

         def plot_kalman_filter(start_pos,
                                sensor_noise,
                                velocity,
                                process_noise):
            # your code goes here
            pass

            interact(plot_kalman_filter,
                     start_pos=(-10, 10),
                     sensor_noise=FloatSlider(value=5, min=0, max=100),
                     velocity=FloatSlider(value=1, min=-2., max=2.),
                     process_noise=FloatSlider(value=5, min=0, max=100.));
```

### 4.16.1 Solution

One possible solution follows. We have sliders for the start position, the amount of noise in the sensor, the amount we move in each time step, and how much movement error there is. Process noise is perhaps the least clear - it models how much the dog wanders off course at each time step, so we add that into the dog's position at each step. I set the random number generator seed so that each redraw uses the same random numbers, allowing us to compare the graphs as we move the sliders.

```
In [38]: from numpy.random import seed
         def plot_kalman_filter(start_pos,
                                sensor_noise,
                                velocity,
                                process_noise):
            N = 20
            zs, ps = [], []
            seed(303)
            dog = DogSimulation(start_pos, velocity, sensor_noise, process_noise)
            zs = [dog.move_and_sense() for _ in range(N)]
            pos = (0., 1000.) # mean and variance
```

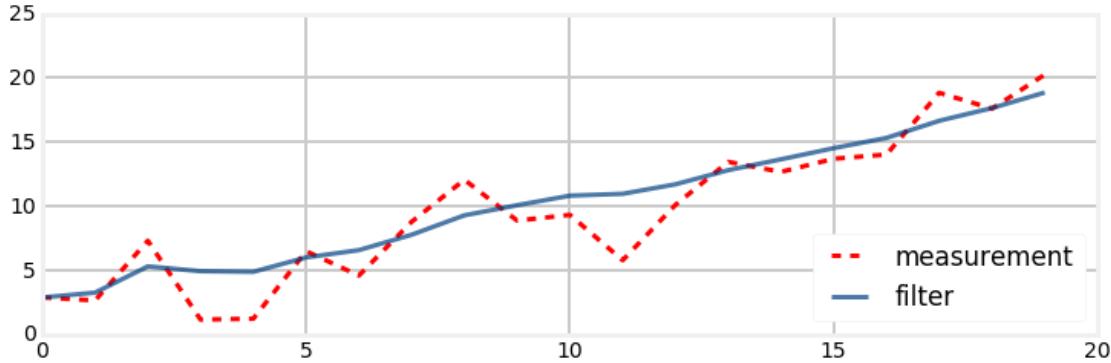
```

for i in range(N):
    pos = predict(pos[0], pos[1], velocity, process_noise)
    pos = update(pos[0], pos[1], zs[i], sensor_noise)
    ps.append(pos[0])

plt.plot(zs, c='r', linestyle='dashed', label='measurement')
plt.plot(ps, c='#004080', alpha=0.7, label='filter')
plt.legend(loc=4);

interact(plot_kalman_filter,
         start_pos=(-10, 10),
         sensor_noise=FloatSlider(value=5, min=0., max=100),
         velocity=FloatSlider(value=1, min=-2., max=2.),
         process_noise=FloatSlider(value=.1, min=0, max=40));

```



## 4.17 Exercise - Nonlinear Systems

Our equations for the Kalman filter are linear:

$$\begin{aligned}\mathcal{N}(\mu, \sigma^2) &= \mathcal{N}(\mu, \sigma^2) + \mathcal{N}(\mu_{\text{move}}, \sigma_{\text{move}}^2) \\ \mathcal{N}(\mu, \sigma^2) &= \mathcal{N}(\mu, \sigma^2) \times \mathcal{N}(\mu_z, \sigma_z^2)\end{aligned}$$

Do you suppose that this filter works well or poorly with nonlinear systems?

Implement a Kalman filter that uses the following equation to generate the measurement value

```

for i in range(100):
    Z = math.sin(i/3.) * 2

```

Adjust the variance and initial positions to see the effect. What is, for example, the result of a very bad initial guess?

In [39]: #enter your code here.

### 4.17.1 Solution

```
In [40]: sensor_var = 30
process_var = 2
pos = (100,500)

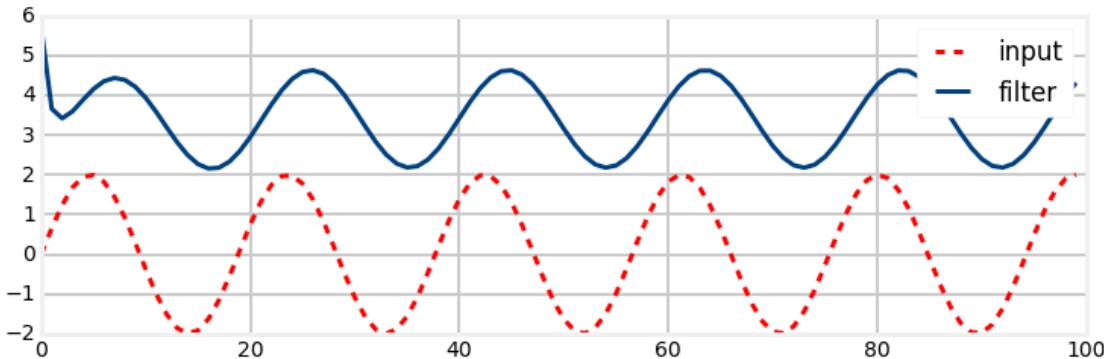
zs, ps = [], []

for i in range(100):
    pos = predict(pos[0], pos[1], velocity, process_var)

    Z = math.sin(i/3.)*2
    zs.append(Z)

    pos = update(pos[0], pos[1], Z, sensor_var)
    ps.append(pos[0])

plt.plot(zs, c='r', linestyle='dashed', label='input')
plt.plot(ps, c='#004080', label='filter')
plt.legend(loc='best');
```



### 4.17.2 Discussion

Here we use a bad initial guess of 100. We can see that the filter never ‘acquires’ the signal. Note now the peak of the filter output always lags the peak of the signal by a small amount, and how the filtered signal does not come very close to capturing the high and low peaks of the input signal.

If we recall the **g-h Filter** chapter we can understand what is happening here. The structure of the g-h filter requires that the filter output chooses a value part way between the prediction and measurement. A varying signal like this one is always accelerating, whereas our process model assumes constant velocity, so the filter is mathematically guaranteed to always lag the input signal.

Maybe we didn’t adjust things ‘quite right’. After all, the output looks like an offset sin wave. Let’s test this assumption.

## 4.18 Exercise - Noisy Nonlinear Systems

Implement the same system, but add noise to the measurement.

In [41]: #enter your code here

### 4.18.1 Solution

```
In [42]: sensor_var = 30
process_var = 2
pos = (100,500)

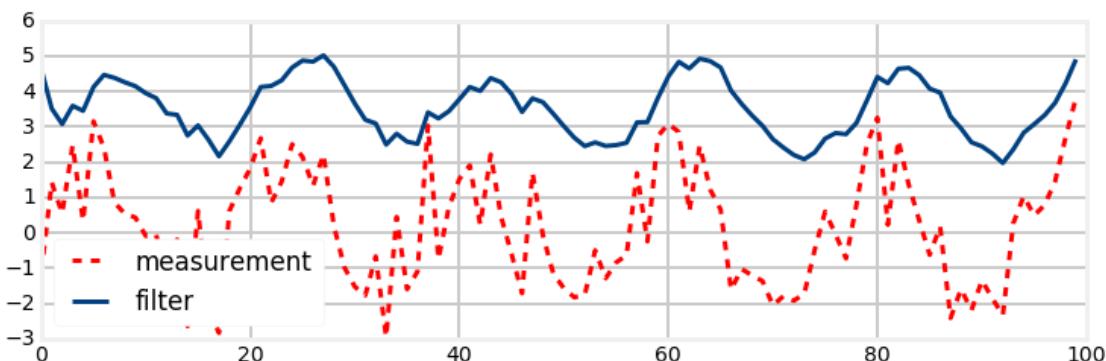
zs, ps = [], []

for i in range(100):
    pos = predict(pos[0], pos[1], velocity, process_var)

    Z = math.sin(i/3.)*2 + randn()*1.2
    zs.append(Z)

    pos = update(pos[0], pos[1], Z, sensor_var)
    ps.append(pos[0])

p1, = plt.plot(zs, c='r', linestyle='dashed', label='measurement')
p2, = plt.plot(ps, c='#004080', label='filter')
plt.legend(loc='best');
```



### 4.18.2 Discussion

This is terrible! The output is not at all like a sin wave, except in the grossest way. With linear systems we could add extreme amounts of noise to our signal and still extract a very accurate result, but here even modest noise creates a very bad result.

Very shortly after practitioners began implementing Kalman filters they recognized the poor performance of them for nonlinear systems and began devising ways of dealing with it. Later chapters are devoted to this problem.

## 4.19 Fixed Gain Filters

Embedded computers usually have extremely limited processors; many do not have floating point circuitry. These simple equations can impose a heavy burden on the chip. This is less true as technology advances, but do not underestimate the value of spending 10 cents less on a processor when millions will be used.

In the example above the variance of the filter converged to a fixed value. This will always happen if the variance of the measurement and process is a constant. You can take advantage of this fact by running simulations to determine what the variance converges to. Then you can hard code this value into your filter. So long as you initialize the filter to a good starting guess (I recommend using the first measurement as your initial value) the filter will perform very well.

So for embedded applications it is common to perform many simulations to determine a fixed Kalman gain value, and then hard code this into the application. Besides reducing chip cost this also makes verifying correctness easier. If you are making a medical device, a device that is being sent into space, or a device that will have one million copies made then an error can cost lives, end a scientific mission, or and/or cost millions of dollars. A professional engineer does everything she can to reduce and eliminate risks.

## 4.20 Derivation From Bayes Theorem (Optional)

I gave you the equations for the product of two Gaussians but did not derive them for you. They are a direct result of applying Bayes rules to Gaussians, however. You don't need to understand this section to read the rest of the book or to use Kalman filters. You may need to know this if you read the literature, much of which is written in terms of Bayes theorem. It's a bit of neat algebra, so why not read on?

We can state the problem as: let the prior be  $N(\mu_p, \sigma_p^2)$ , and measurement be  $z \propto N(z, \sigma_z^2)$ . What is the posterior  $x$  given the measurement  $z$ ?

Write the posterior as  $P(x|z)$ . Now we can use Bayes Theorem to state

$$P(x|z) = \frac{P(z|x)P(x)}{P(z)}$$

$P(z)$  is a normalizing constant, so we can create a proportionality

$$P(x|z) \propto P(z|x)P(x)$$

Now we substitute in the equations for the Gaussians, which are

$$\begin{aligned} P(z|x) &= \frac{1}{\sqrt{2\pi\sigma_z^2}} \exp\left[-\frac{(z-x)^2}{2\sigma_z^2}\right] \\ P(x) &= \frac{1}{\sqrt{2\pi\sigma_p^2}} \exp\left[-\frac{(x-\mu_p)^2}{2\sigma_p^2}\right] \end{aligned}$$

We can drop the leading terms, as they are constants, giving us

$$\begin{aligned} P(x|z) &\propto \exp\left[-\frac{(z-x)^2}{2\sigma_z^2}\right] \exp\left[-\frac{(x-\mu_p)^2}{2\sigma_p^2}\right] \\ &\propto \exp\left[-\frac{(z-x)^2}{2\sigma_z^2} - \frac{(x-\mu_p)^2}{2\sigma_p^2}\right] \\ &\propto \exp\left[-\frac{1}{2\sigma_z^2\sigma_p^2}[\sigma_p^2(z-x)^2 - \sigma_z^2(x-\mu_p)^2]\right] \end{aligned}$$

Now we multiply out the squared terms and group in terms of the posterior  $x$ .

$$\begin{aligned} P(x|z) &\propto \exp \left[ -\frac{1}{2\sigma_z^2\sigma_p^2} [\sigma_p^2(z^2 - 2xz + x^2) + \sigma_z^2(x^2 - 2x\mu_p + \mu_p^2)] \right] \\ &\propto \exp \left[ -\frac{1}{2\sigma_z^2\sigma_p^2} [x^2(\sigma_p^2 + \sigma_z^2) - 2x(\sigma_z^2\mu_p + \sigma_p^2z) + (\sigma_p^2z^2 + \sigma_z^2\mu_p^2)] \right] \end{aligned}$$

The last parentheses do not contain the posterior  $x$ , so it can be treated as a constant and discarded.

$$P(x|z) \propto \exp \left[ -\frac{1}{2} \frac{x^2(\sigma_p^2 + \sigma_z^2) - 2x(\sigma_z^2\mu_p + \sigma_p^2z)}{\sigma_z^2\sigma_p^2} \right]$$

Divide numerator and denominator by  $\sigma_p^2 + \sigma_z^2$  to get

$$P(x|z) \propto \exp \left[ -\frac{1}{2} \frac{x^2 - 2x \left( \frac{\sigma_z^2\mu_p + \sigma_p^2z}{\sigma_p^2 + \sigma_z^2} \right)}{\frac{\sigma_z^2\sigma_p^2}{\sigma_p^2 + \sigma_z^2}} \right]$$

Proportionality lets us create or delete constants at will, so we can factor this into

$$P(x|z) \propto \exp \left[ -\frac{1}{2} \frac{(x - \frac{\sigma_z^2\mu_p + \sigma_p^2z}{\sigma_p^2 + \sigma_z^2})^2}{\frac{\sigma_z^2\sigma_p^2}{\sigma_p^2 + \sigma_z^2}} \right]$$

A Gaussian is

$$N(m, s^2) \propto \exp \left[ -\frac{1}{2} \frac{(x - m)^2}{s^2} \right]$$

So we can see that  $P(x|z)$  has a mean of

$$\mu_{\text{posterior}} = \frac{\sigma_z^2\mu_p + \sigma_p^2z}{\sigma_p^2 + \sigma_z^2}$$

and a variance of

$$\sigma_{\text{posterior}} = \frac{\sigma_z^2\sigma_p^2}{\sigma_p^2 + \sigma_z^2} = \frac{1}{\frac{1}{\sigma_p^2} + \frac{1}{\sigma_z^2}}$$

which are the equations that we have been using.

## 4.21 Summary

The one dimensional Kalman filter that we describe in this chapter is a special, restricted case of the more general filter. Most texts do not discuss the one dimensional form. However, I think it is a vital stepping stone. We started the book with the g-h filter, then implemented the discrete Bayes filter, and now implemented the one dimensional Kalman filter. I have tried to show you that each of these filters use the same algorithm and reasoning. The mathematics of the Kalman filter that we will learn shortly is fairly sophisticated, and it can be difficult to understand the underlying simplicity of the filter. That sophistication comes with significant benefits: the generalized filter will markedly outperform the filters in this chapter.

This chapter takes time to assimilate. To truly understand it you will probably have to work through this chapter several times. I encourage you to change the various constants in the code and observe the results.

Convince yourself that Gaussians are a good representation of a unimodal belief of the position of a dog in a hallway, the position of an aircraft in the sky, or the temperature of a chemical reaction chamber. Then convince yourself that multiplying Gaussians truly does compute a new belief from your prior belief and the new measurement. Finally, convince yourself that if you are measuring movement, that adding the Gaussians together updates your belief.

If you do not fully understand this, reread this chapter. Try implementing the filter from scratch, just by looking at the equations and reading the text. Change the constants. Maybe try to implement a different tracking problem, like tracking stock prices. Experimentation will build your intuition and understanding of how these marvelous filters work.

Most of all, spend enough time with the **Full Description of the Algorithm** section to ensure you understand the algorithm and how it relates to the g-h filter and discrete Bayes filter. There is just one ‘trick’ here - selecting a value somewhere between a prediction and a measurement. Each algorithm performs that trick with different math, but all use the same logic.



# Chapter 5

## Multivariate Gaussians - Modeling Uncertainty in Multiple Dimensions

### 5.1 Introduction

The techniques in the last chapter are very powerful, but they only work with one variable or dimension. Gaussians represent a mean and variance that are scalars - real numbers. They provide no way to represent multidimensional data, such as the position of a dog in a field. You may retort that you could use two Kalman filters from the last chapter. One would track the x coordinate and the other the y coordinate. That does work, but suppose we want to track position, velocity, acceleration, and attitude. These values are related to each other, and as we learned in the g-h chapter we should never throw away information. Through one key insight we will achieve markedly better filter performance than was possible with the equations from the last chapter.

In this chapter I will introduce you to multivariate Gaussians - Gaussians for more than one variable, and the key insight I mention above. Then, in the next chapter we will use the math from this chapter to write a complete filter in just a few lines of code.

### 5.2 Multivariate Normal Distributions

In the last two chapters we used Gaussians for a scalar (one dimensional) variable, expressed as  $\mathcal{N}(\mu, \sigma^2)$ . A more formal term for this is univariate normal, where univariate means ‘one variable’. The probability distribution of the Gaussian is known as the univariate normal distribution.

What might a multivariate normal distribution be? Multivariate means multiple variables. Our goal is to be able to represent a normal distribution across multiple dimensions. I don’t necessarily mean spatial dimensions - it could be position, velocity, and acceleration. Consider a two dimensional case. Let’s say we believe that  $x = 2$  and  $y = 17$ . This might be the x and y coordinates for the position of our dog, it might be the position and velocity of our dog on the x-axis, or the temperature and wind speed at our weather station. It doesn’t really matter. We can see that for  $N$  dimensions, we need  $N$  means, which we will arrange in a column matrix (vector) like so:

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

Therefore for this example we would have

$$\mu = \begin{bmatrix} 2 \\ 17 \end{bmatrix}$$

The next step is representing our variances. At first blush we might think we would also need N variances for N dimensions. We might want to say the variance for x is 10 and the variance for y is 4, like so.

$$\sigma^2 = \begin{bmatrix} 10 \\ 4 \end{bmatrix}$$

This is incomplete because it does not consider the more general case. In the **Gaussians** chapter we computed the variance in the heights of students. We can do the same for their weights. But there is also a relationship between height and weight. In general, a taller person weighs more than a shorter person. We say they are correlated. We want a way to express not only what we think the variance is in the height and the weight, but also the degree to which they are correlated.

Before we can understand multivariate normal distributions we need to understand the mathematics behind correlations.

### 5.3 Correlation and Covariance

The covariance describes how much two variables vary together. Covariance is short for correlated variances. In other words, variance is a measure for how a population vary amongst themselves, and covariance is a measure for how much two variables change in relation to each other. Generally speaking as height increases weight also increases. These variables are correlated. It is also called positively correlated because as one variable gets larger so does the other. As the outdoor temperature decreases home heating bills increase. These are inversely correlated because as one variable gets larger the other variable lowers. This is also called negatively correlated. The price of tea and the number of tail wags my dog makes have no relation to each other, and we say they are uncorrelated - each can change independent of the other.

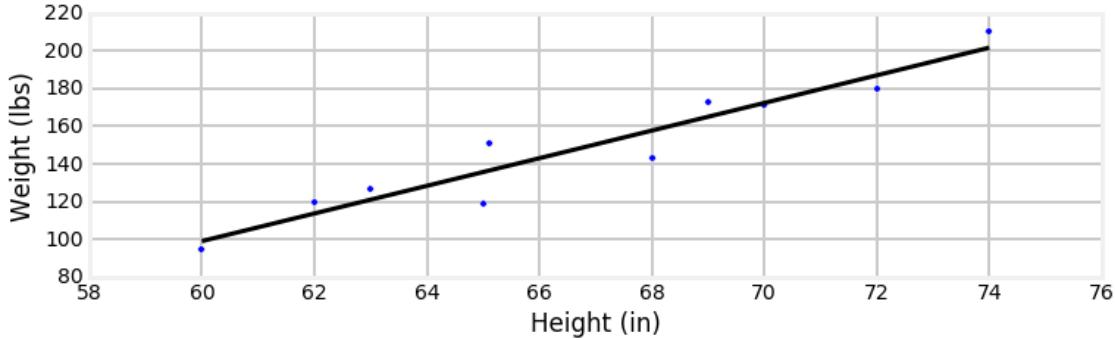
Correlation implies prediction. If you are significantly taller than me I can predict that you also weigh more. As winter comes I predict that I will be spending more to heat my house. If my dog wags his tail more I don't conclude that tea prices will be increasing.

For example, here is a plot of height and weight of students on the school's track team. If a student is 67 inches tall I can predict they weigh roughly 140 pounds. Since the correlation is not perfect neither is my prediction.

```
In [2]: from gaussian_internal import plot_correlated_data

height = [60, 62, 63, 65, 65.1, 68, 69, 70, 72, 74]
weight = [95, 120, 127, 119, 151, 143, 173, 171, 180, 210]

plot_correlated_data(height, weight, 'Height (in)', 'Weight (lbs)', False)
```



In this book we will always be using linear correlation. We assume that the relationship between variables is linear. That is, a straight line is a good fit for the data. I've fit a straight line through the data in the above chart. The concept of nonlinear correlation exists, but we will not be using it.

The equation for the covariance between  $X$  and  $Y$  is

$$COV(X, Y) = \sigma_{xy} = \mathbb{E}[(X - \mu_x)(Y - \mu_y)]$$

Where  $\mathbb{E}[X]$  is the expected value of  $X$  defined as

$$\mathbb{E}[X] = \begin{cases} \sum_{i=1}^n p_i x_i & \text{discrete} \\ \int_{-\infty}^{\infty} x f(x) & \text{continuous} \end{cases}$$

Compare the covariance equation to the equation for the variance:

$$VAR(X) = \sigma_x^2 = \mathbb{E}[(X - \mu)^2]$$

As you can see they are very similar.

We use a covariance matrix to denote covariances of a multivariate normal distribution, and it looks like this:

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & \sigma_n^2 \end{bmatrix}$$

If you haven't seen this before it is probably a bit confusing. Instead of continuing with the mathematical definitions I will build your intuition with thought experiments. At this point, note that the diagonal contains the variance for each variable, and that all off-diagonal elements (covariances) represent how much the  $i$ th (horizontal row) and  $j$ th (vertical column) variables are linearly correlated to each other. So  $\sigma_3^2$  is the variance of the third variable, and  $\sigma_{13}$  is the covariance between the first and third variables.

A covariance of 0 indicates no correlation. If the variance for  $x$  is 10, the variance for  $y$  is 4, and there is no linear correlation between  $x$  and  $y$ , then we would write

$$\Sigma = \begin{bmatrix} 10 & 0 \\ 0 & 4 \end{bmatrix}$$

If there was a small amount of positive correlation between  $x$  and  $y$  we might have

$$\Sigma = \begin{bmatrix} 10 & 1.2 \\ 1.2 & 4 \end{bmatrix}$$

where 1.2 is the covariance between  $x$  and  $y$ . I say the correlation is “small” because the covariance of 1.2 is small relative to the variances of 10.

If there was a large amount of negative correlation between  $x$  and  $y$  we might have

$$\Sigma = \begin{bmatrix} 10 & -9.7 \\ -9.7 & 4 \end{bmatrix}$$

Note that the matrix always symmetric. Logically the covariance between  $x$  and  $y$  is always equal to the covariance between  $y$  and  $x$ . That is,  $\sigma_{xy} = \sigma_{yx}$  for any  $x$  and  $y$ .

Let’s work through a few examples. `numpy.cov` computes the covariance matrix. Lets create data with perfectly correlated data. By this I mean that the data perfectly fits on a line - there is no variance from the line.

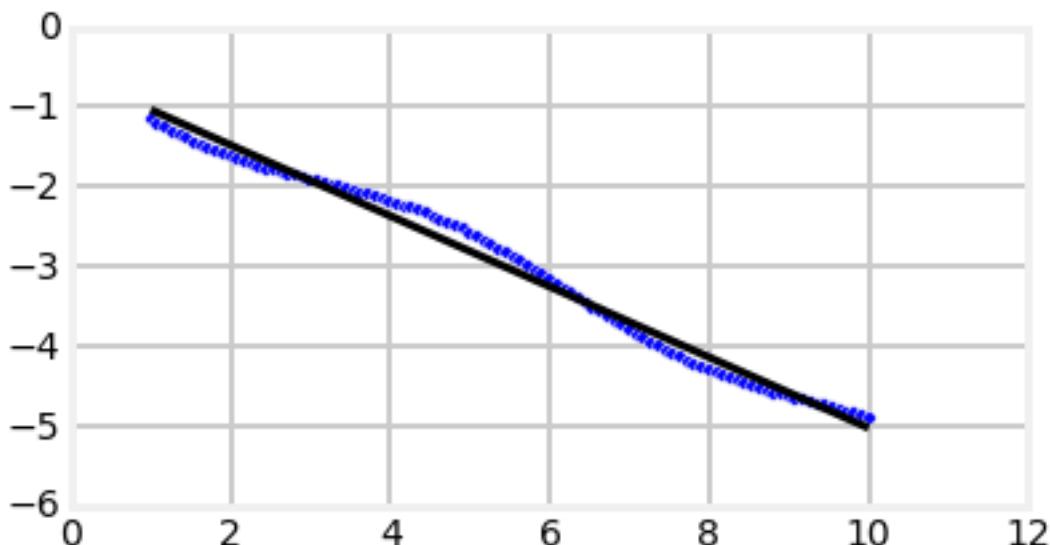
```
In [3]: import numpy as np
X = np.linspace(1, 10, 100)
Y = np.linspace(1, 10, 100)
print(np.cov(X, Y))

[[ 6.956  6.956]
 [ 6.956  6.956]]
```

We can see from the covariance matrix that the covariance is equal to the variance in  $x$  and in  $y$ .

Now let’s add some noise to one of the variables so that they are no longer perfectly correlated. I will make  $Y$  negative to create a negative correlation.

```
In [4]: X = np.linspace(1, 10, 100)
Y = -(np.linspace(1, 5, 100) + np.sin(X)*.2)
plot_correlated_data(X, Y)
print(np.cov(X, Y))
```

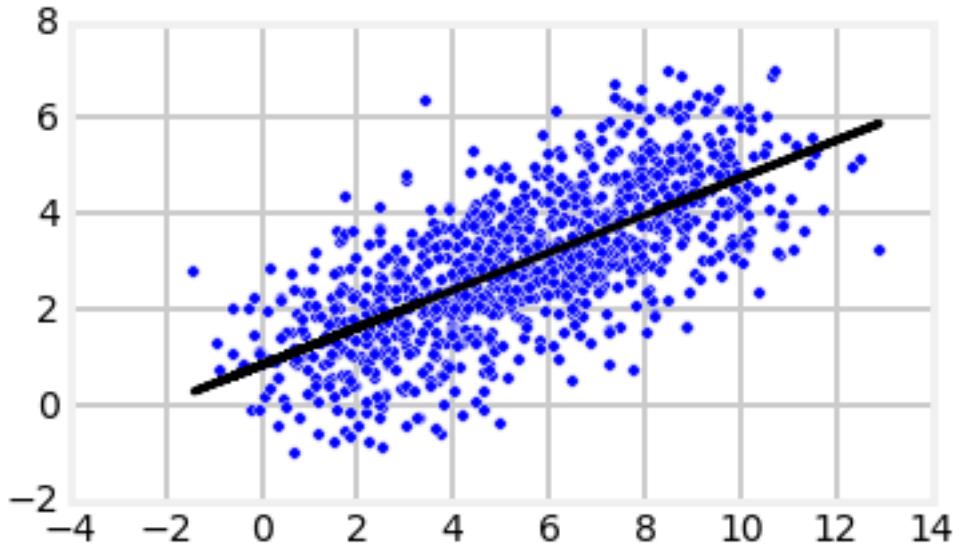


```
[[ 6.956 -3.084]
 [-3.084  1.387]]
```

The data no longer forms a straight line. The covariance is  $\sigma_{xy} = -3.08$ . The absolute value is smaller than the variance of  $X$  or  $Y$ . It is not close to zero, and so we know there is still a high degree of correlation. We can verify this by looking at the chart. The data forms nearly a straight line.

Now I will add random noise to a straight line.

```
In [5]: from numpy.random import randn
X = np.linspace(1, 10, 1000) + randn(1000)
Y = np.linspace(1, 5, 1000) + randn(1000)
plot_correlated_data(X, Y)
print(np.cov(X, Y))
```

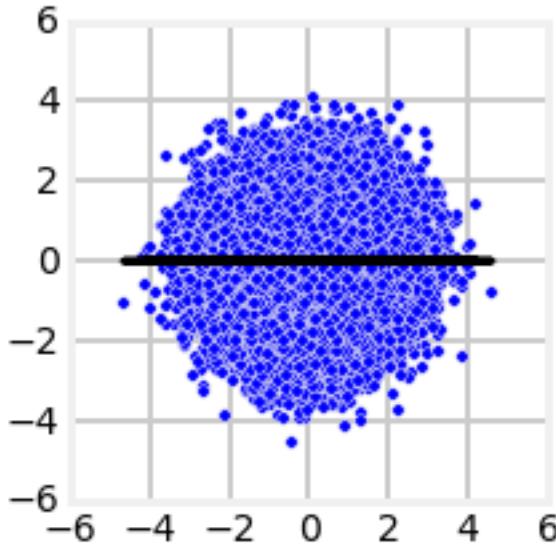


```
[[ 7.78   3.043]
 [ 3.043  2.439]]
```

We see that the variance is smaller in relation to the variances, reflecting the lower correlation between  $X$  and  $Y$ . We can still fit a straight line through this data, but there is much greater variation in the data.

Finally, here is the covariance between completely random data.

```
In [6]: X = randn(100000)
Y = randn(100000)
plot_correlated_data(X, Y)
print(np.cov(X, Y))
```



```
[[ 1.002  0.    ]
 [ 0.      1.011]]
```

Here the covariances are very near zero. As you can see with the plot, there is no clear way to draw a line to fit the data. A vertical line would be as convincing as the horizontal line shown.

## 5.4 Multivariate Normal Distribution Equation

Here is the multivariate normal distribution in  $n$  dimensions.

$$f(\mathbf{x}, \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right]$$

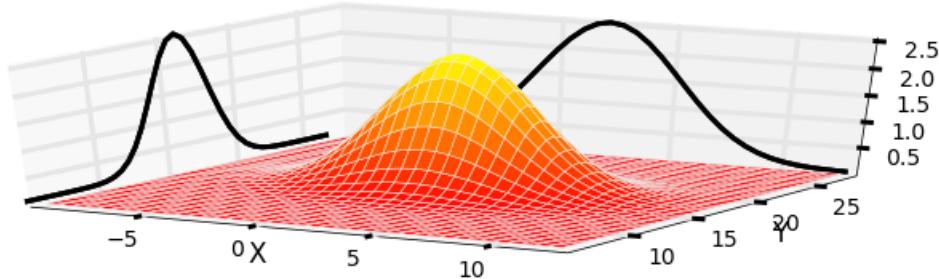
I urge you to not try to remember this equation. We will program it in a Python function and then call it if we need to compute a specific value. However, note that it has the same form as the univariate normal distribution:

$$f(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{1}{2} (x - \mu)^2 / \sigma^2 \right]$$

The multivariate version merely replaces the scalars of the univariate equations with matrices. If you are reasonably well-versed in linear algebra this equation should look quite manageable; if not, don't worry, we have code to compute it for us! Let's plot it and see what it looks like.

```
In [7]: import mkf_internal
mean = (2, 17)
cov = [[10., 0],
       [0, 4.]]
```

```
mkf_internal.plot_3d_covariance(mean, cov)
```



This is a plot of multivariate Gaussian with a mean of  $\mu = [2, 7]$  and a covariance of  $\Sigma = [10, 0; 0, 4]$ . The three dimensional shape shows the probability density of for any value of  $(X, Y)$  in the z-axis. I have projected the variance for  $x$  and  $y$  onto the walls of the chart - you can see that they take on the normal Gaussian bell curve shape. The curve for  $X$  is wider than the curve for  $Y$ , which is explained by  $\sigma_x^2 = 10$  and  $\sigma_y^2 = 4$ . The highest point of the curve is at the the means for  $X$  and  $Y$ .

All multivariate Gaussians have this shape. If we think of this as a the Gaussian for the position of a dog, the z-value at each point of  $(X, Y)$  is the probability density of it being at that position. Strictly speaking this is the joint probability density function, which I will define soon. So, the dog has the highest probability of being near  $(2, 17)$ , a modest probability of being near  $(5, 14)$ , and a very low probability of being near  $(10, 10)$ . As with the univariate case this is a probability density, not a probability. Continuous distributions have an infinite range, and so the probability of being exactly at  $(2, 17)$ , or any point, is 0%. We can compute the probability of being within a given range by computing the area under the curve.

FilterPy [2] implements the equation with the function `filterpy.stats.multivariate_gaussian`. SciPy's `stats` module implements the multivariate normal equation with `multivariate_normal()`. It implements a 'frozen' form where you set the mean and covariance once, and then calculate the probability for any number of values for  $x$  over any arbitrary number of calls. This is much more efficient then recomputing everything in each call. I named my function `multivariate_gaussian()` to ensure it is never confused with the SciPy version. I will say that for a single call, where the frozen variables do not matter, mine consistently runs faster as measured by the `timeit` function.

The tutorial[1] for the `scipy.stats` module explains 'freezing' distributions and other very useful features.

```
In [8]: from filterpy.stats import gaussian, multivariate_gaussian
```

I'll demonstrate using it, and then move on to more interesting things.

First, let's find the probability density for our dog being at  $(2.5, 7.3)$  if we believe he is at  $(2, 7)$  with a variance of 8 for  $x$  and a variance of 3 for  $y$ .

Start by setting  $x$  to  $(2.5, 7.3)$ . You can use a tuple, list, or NumPy array.

```
In [9]: x = [2.5, 7.3]
```

Next, we set the mean of our belief:

```
In [10]: mu = [2.0, 7.0]
```

Finally, we have to define our covariance matrix. In the problem statement we did not mention any correlation between  $x$  and  $y$ , and we will assume there is none. This makes sense; a dog can choose to independently

wander in either the  $x$  direction or  $y$  direction without affecting the other. If there is no correlation between the values place the variances in the diagonal, and set off-diagonal elements to zero. I will use name  $P$ . Kalman filters use the name  $\mathbf{P}$  for the covariance matrix, and we need to become familiar with the conventions.

```
In [11]: P = [[8., 0.],
           [0., 3.]]
```

Now call the function

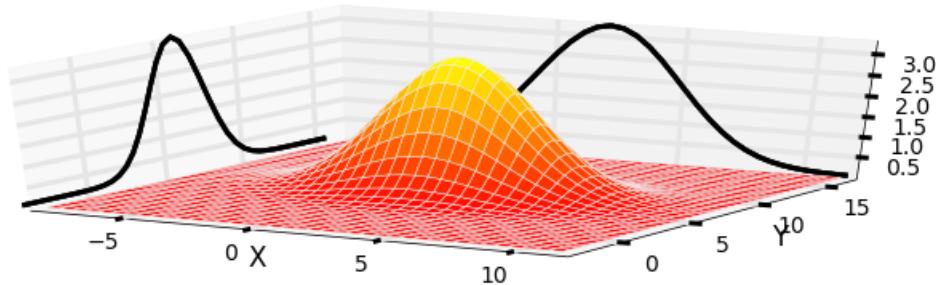
```
In [12]: print('{:.4}'.format(multivariate_gaussian(x, mu, P)))
```

0.03151

These numbers are not easy to interpret. Let's view a plot of it.

```
In [13]: import mkf_internal
import matplotlib.pyplot as plt

ax = mkf_internal.plot_3d_covariance(mu, P)
```



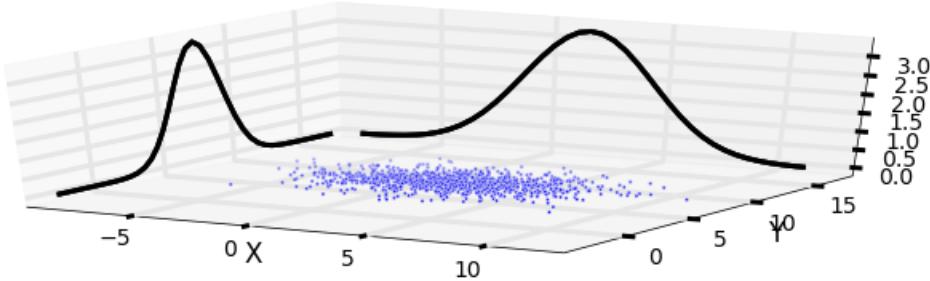
The result is clearly a 3D bell shaped curve. We can see that the Gaussian is centered around  $(2, 7)$ , and that the probability density quickly drops away in all directions. This plot shows what is known as the joint probability. On the sides of the plot I have drawn the marginal probability for  $X$  and  $Y$  (defined below).

It's time to define some terms. The joint probability, denoted  $P(x, y)$ , is the probability of both  $x$  and  $y$  happening. For example, if you have two die  $P(2, 5)$  is the probability of the first die rolling a 2 and the second die rolling a 5. Assuming the die are six sided and fair, the probability  $P(2, 5) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$ . The 3D chart above shows the joint probability density function.

The marginal probability is the probability of an event happening without regard of any other event. In the chart above the Gaussian curve drawn to the left is the marginal for  $Y$ . This is the probability for the dog being at any position in  $Y$  disregarding the value for  $X$ . In back I've plotted the marginal for  $X$ , which is the probability of the dog being at any position in  $X$  disregarding  $Y$ .

Let's look at this in a slightly different way. Instead of plotting a surface showing the probability distribution I will generate 1,000 points with the distribution  $\begin{bmatrix} 8 & 0 \\ 0 & 3 \end{bmatrix}$ .

```
In [14]: mkf_internal.plot_3d_sampled_covariance(mu, P)
```

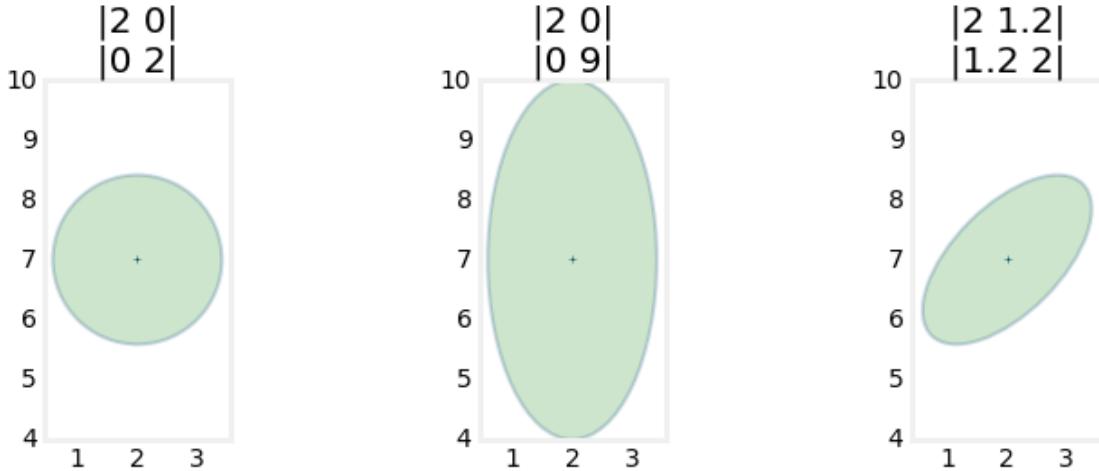


We can think of the sampled points as being possible locations for our dog given those particular mean and covariances. The contours on the side show the marginal probability for  $X$  and  $Y$ . We can see that he is far more likely to be at  $(2, 7)$  where there are many points, than at  $(-5, 5)$  where there are few.

As beautiful as these plots are, it is hard to get useful information from them. For example, it is not easy to tell if  $X$  and  $Y$  both have the same variance. In most of the book I'll display Gaussians using contour plots using helper functions from FilterPy.

```
In [15]: from book_format import set_figsize, figsize
```

```
with figsize(y=4):
    mkf_internal.plot_3_covariances()
```

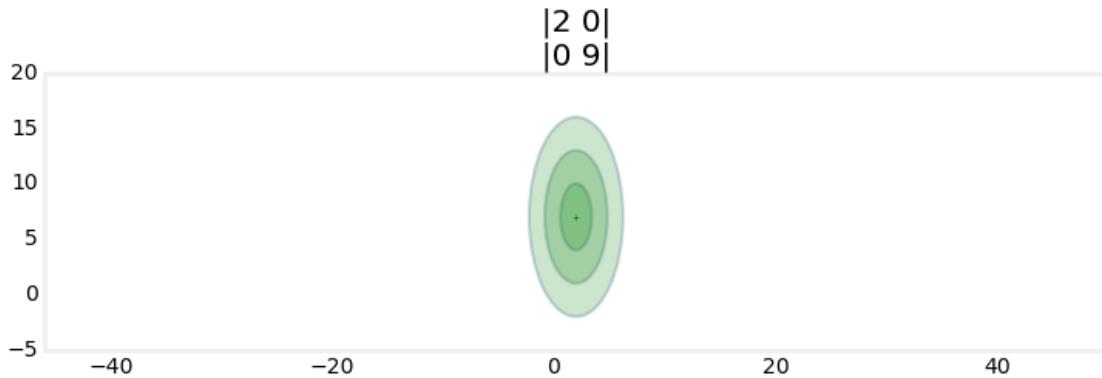


For those of you viewing this online or in Jupyter Notebook on your computer, here is an animation of varying the covariance while holding the variance constant.

(source: <http://git.io/vqxLS>)

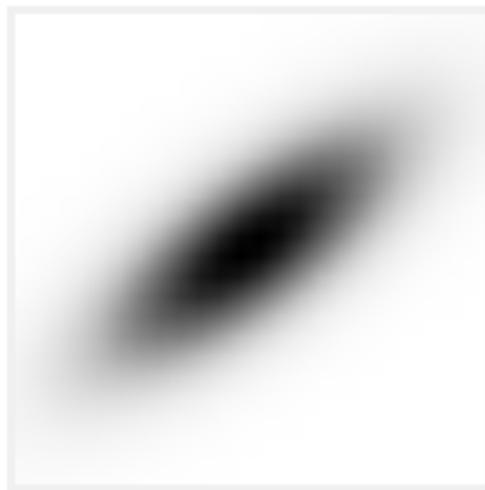
The contour plots display the range of values that the multivariate Gaussian takes for a specific standard deviation. This is like taking a horizontal slice out of the 3D plot. By default it displays one standard deviation, but you can use either the `variance` or `std` parameter to control what is displayed. For example, `variance=3**2` or `std=3` would display the 3rd standard deviation, and `variance=[1, 4, 9]` or `std=[1, 2, 3]` would display the 1st, 2nd, and 3rd standard deviations.

```
In [16]: from filterpy.stats import plot_covariance_ellipse
P = [[2, 0], [0, 9]]
plot_covariance_ellipse((2, 7), P, fc='g', alpha=0.2,
                        std=[1, 2, 3],
                        title='|2 0|\n|0 9|')
plt.gca().grid(b=False)
```



The solid colors may suggest to you that the probability distribution is constant between the standard deviations. This is not true, as you can tell from the 3D plot of the Gaussian. Here is a 2D shaded representation of the probability distribution for the covariance  $\begin{pmatrix} 2 & 1.2 \\ 1.2 & 1.3 \end{pmatrix}$ . Darker gray corresponds to higher probability density.

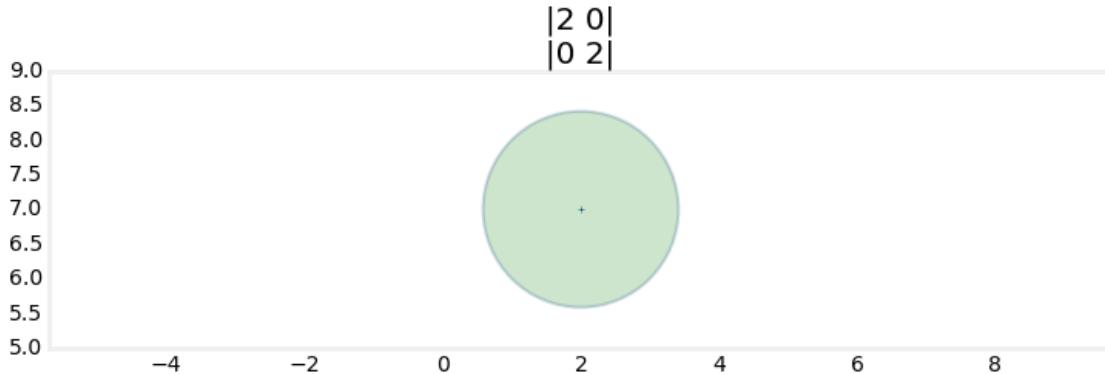
```
In [17]: from nonlinear_plots import plot_cov_ellipse_colormap
plot_cov_ellipse_colormap(cov=[[2, 1.2], [1.2, 1.3]])
```



Thinking about the physical interpretation of these plots clarifies their meaning. The mean and covariance of the first plot is

$$\mu = \begin{bmatrix} 2 \\ 7 \end{bmatrix}, \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

```
In [18]: x = [2, 7]
P = [[2, 0], [0, 2]]
plot_covariance_ellipse(x, P, fc='g', alpha=0.2,
                        title='|2 0|\n|0 2|')
plt.gca().grid(b=False)
```

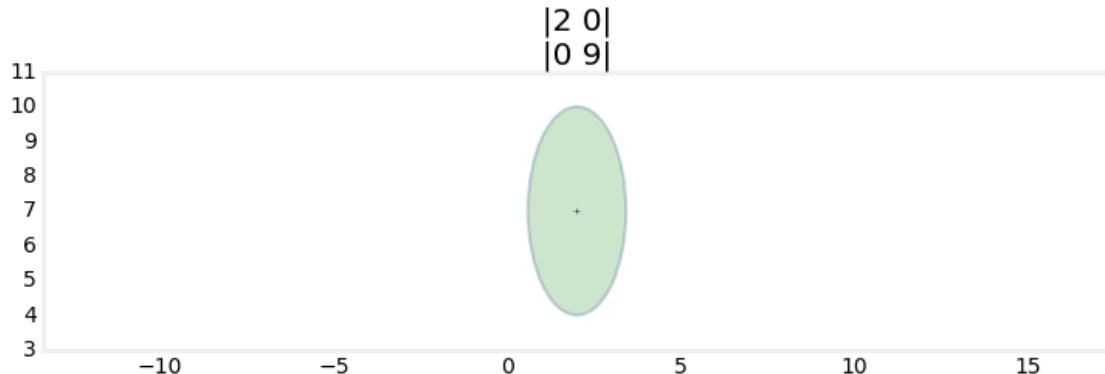


A Bayesian way of thinking about this is that the ellipse shows us the amount of error in our belief. A tiny circle would indicate that we have a very small error, and a very large circle indicates a lot of error in our belief. The shape of the ellipse shows us the geometric relationship of the errors in  $X$  and  $Y$ . Here we have a circle so errors in  $X$  and  $Y$  are equally likely.

The second plot is for the mean and covariance

$$\mu = \begin{bmatrix} 2 \\ 7 \end{bmatrix}, \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 9 \end{bmatrix}$$

```
In [19]: x = [2, 7]
P = [[2, 0], [0, 9]]
plot_covariance_ellipse(x, P, fc='g', alpha=0.2,
                        title='|2 0|\n|0 9|')
plt.gca().grid(b=False)
```

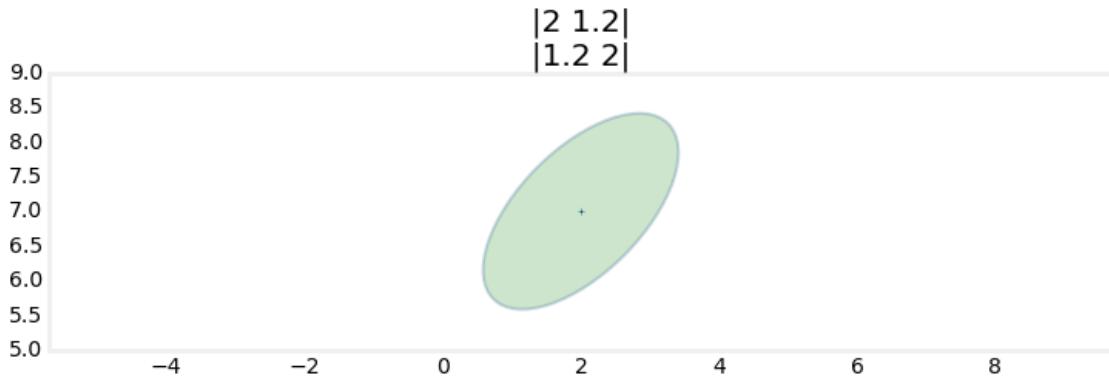


This time we use a different variance for  $X$  ( $\sigma_x^2 = 2$ ) vs  $Y$  ( $\sigma_y^2 = 9$ ). The result is a tall and narrow ellipse. We can see that a lot more uncertainty in  $Y$  vs  $X$ . In both cases we believe the dog is at  $(2, 7)$ , but the uncertainties are different.

The third plot shows the mean and covariance

$$\mu = \begin{bmatrix} 2 \\ 7 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} 2 & 1.2 \\ 1.2 & 2 \end{bmatrix}$$

```
In [20]: x = [2, 7]
P = [[2, 1.2], [1.2, 2]]
plot_covariance_ellipse(x, P, fc='g', alpha=0.2,
                        title='|2 1.2|\n|1.2 2|')
plt.gca().grid(b=False)
```

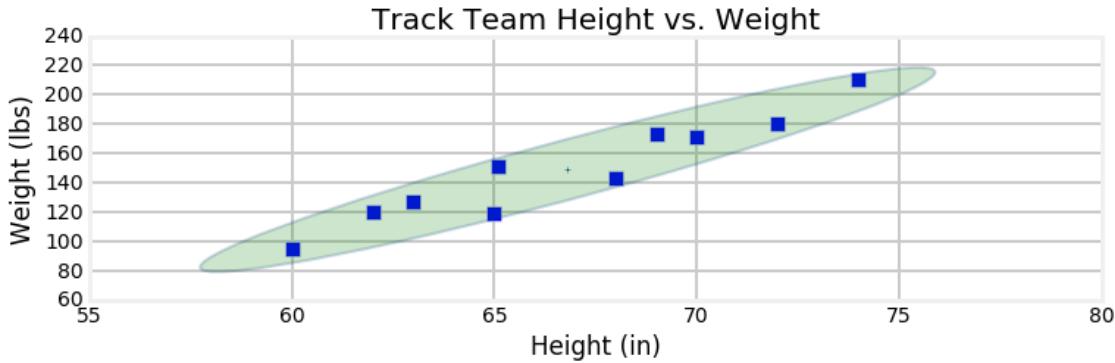


This is the first contour that has values in the off-diagonal elements of the covariance, and this is the first contour plot with a slanted ellipse. This is not a coincidence. The two facts are telling us the same thing. A slanted ellipse tells us that the  $x$  and  $y$  values are somehow correlated. The off-diagonal elements in the covariance matrix are non-zero, indicating that a correlation exists.

Recall the plot for height versus weight. It formed a slanted grouping of points. We can use NumPy's `cov()` function to compute the covariance of two or more variables by placing them into a 2D array. Let's do that, then plot the  $2\sigma$  covariance ellipse on top of the data.

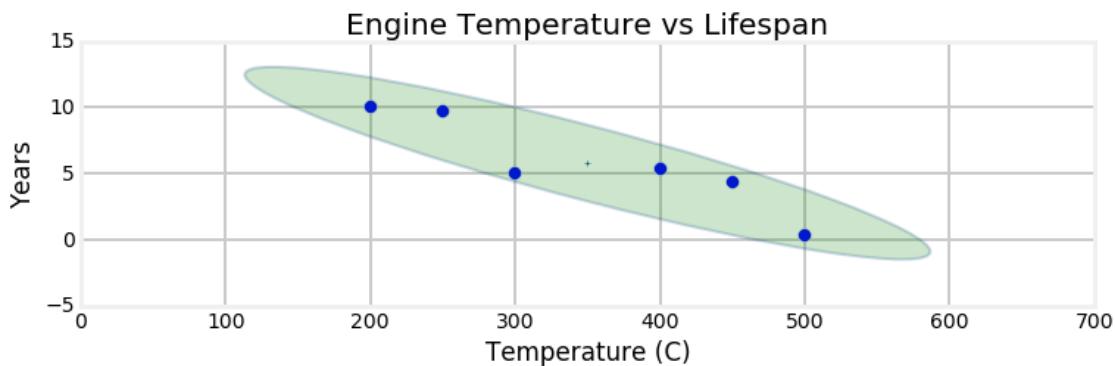
```
In [21]: import numpy as np  
cov_hw = np.cov(np.vstack((height, weight)))  
cov_hw
```

```
Out[21]: array([[ 20.583,  150.779],  
                 [ 150.779, 1213.656]])
```



This should help you form a strong intuition on the meaning and use of covariances. The covariance ellipse shows you how the data is ‘scattered’ in relation to each other. A narrow ellipse like this tells you that the data is very correlated. There is only a narrow range of weights for any given height. The ellipse leans towards the right, telling us there is a positive correlation - as x increases y also increases. If the ellipse leaned towards the left then the correlation would be negative - as x increases y decreases. We can see this in the following plot:

```
In [23]: max_temp = [200, 250, 300, 400, 450, 500]
lifespan = [10, 9.7, 5, 5.4, 4.3, 0.3]
plt.scatter(max_temp, lifespan, s=80)
cov = np.cov(np.vstack((max_temp, lifespan)))
plot_covariance_ellipse((np.mean(max_temp), np.mean(lifespan)), cov, fc='g',
alpha=0.2, axis_equal=False, std=2)
plt.title('Engine Temperature vs Lifespan')
plt.xlabel('Temperature (C)'); plt.ylabel('Years');
```



The relationships between variances and covariances can be hard to puzzle out by inspection, so here is an interactive plot. (If you are reading this in a static form instructions to run this online are here: <http://git.io/vmv6e>)

```
In [24]: from IPython.html.widgets import interact, interactive, fixed
from IPython.html.widgets import FloatSlider
```

```

def plot_covariance(var_x, var_y, cov_xy):
    P1 = [[var_x, cov_xy], [cov_xy, var_y]]

    plot_covariance_ellipse((10, 10), P1, axis_equal=False,
                           show_semiaxis=True)

    plt.xlim(4, 16)
    plt.gca().set_aspect('equal')
    plt.ylim(4, 16)
    plt.show()

interact (plot_covariance,
          var_x=FloatSlider(value=5., min=0, max=20.),
          var_y=FloatSlider(value=5., min=0., max=20.),
          cov_xy=FloatSlider(value=1.5, min=0.0, max=50, step=.2));

```



### 5.4.1 Pearson's Correlation Coefficient

The correlation between two variables can be given a numerical value with Pearson's Correlation Coefficient. It is defined as

$$\rho_{xy} = \frac{COV(X, Y)}{\sigma_x \sigma_y}$$

This value can range in value from -1 to 1. If the covariance is 0 than  $\rho = 0$ . A value greater than 0 indicates that the relationship is a positive correlation, and a negative value indicates that there is a negative correlation. Values near -1 or 1 indicate a very strong correlation, and values near 0 indicate a very weak correlation.

Correlation and covariance are very closely related. Covariance has units associated with it, and correlation is a unitless ratio. For example, for our dog  $\sigma_{xy}$  has units of meters squared.

We can use `scipy.stats.pearsonr` function to compute the Pearson coefficient. It returns a tuple of the Pearson coefficient and of the 2 tailed p-value. The latter is not used in this book. Here we compute  $\rho$  for height vs weight of student athletes:

```
In [25]: from scipy.stats import pearsonr
pearsonr(height, weight)[0]
```

```
Out[25]: 0.95397310960801929
```

Here we compute the correlation between engine temperature and lifespan.

```
In [26]: pearsonr(max_temp, lifespan)[0]
```

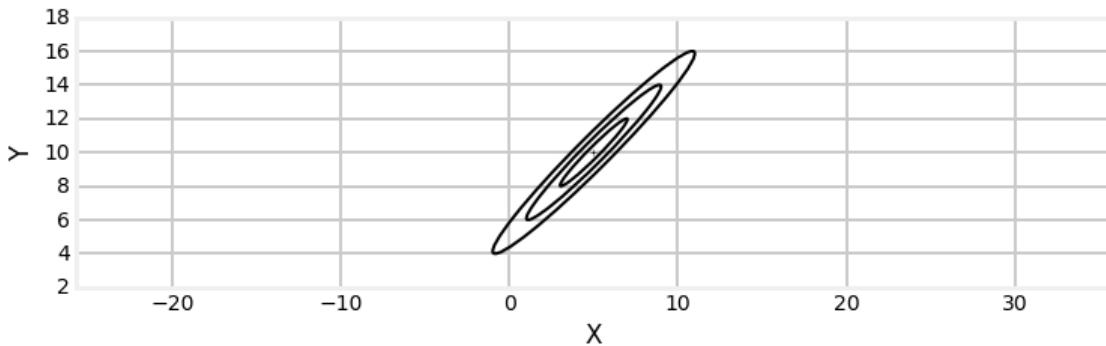
```
Out[26]: -0.91782234535272555
```

We will not be using this coefficient much in this book, but you may see it elsewhere.

## 5.5 Using Correlations to Improve Estimates

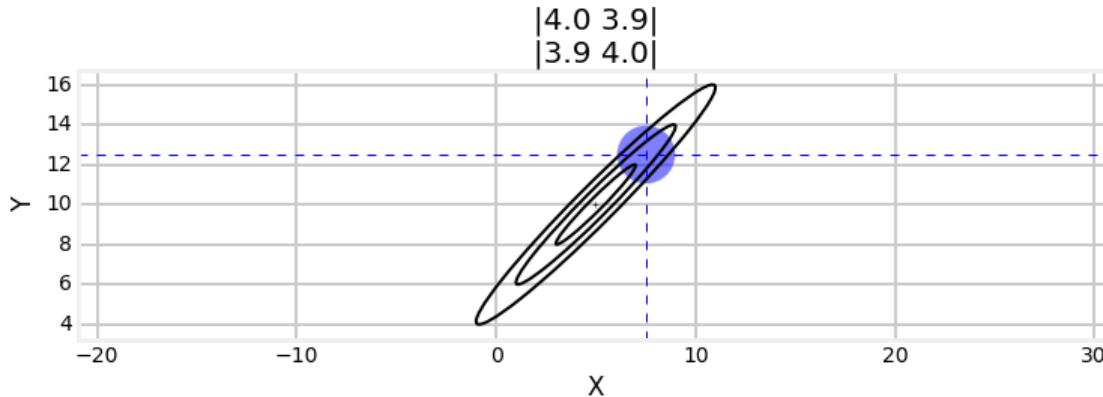
Suppose we believe our dog is at position  $(5, 10)$  with some given covariance. If the standard deviation in  $x$  and  $y$  is each 2 meters, but they are strongly correlated, the covariance contour would look something like this.

```
In [27]: P = [[4, 3.9], [3.9, 4]]
plot_covariance_ellipse((5, 10), P, ec='k', std=[1, 2, 3])
plt.xlabel('X')
plt.ylabel('Y');
```



Now suppose I were to tell you that we know that  $x = 7.5$ . What can we infer about the value for  $y$ ? The position is extremely likely to lie within the  $3\sigma$  covariance ellipse. We can infer the position in  $y$  based on the covariance matrix because there is a correlation between  $x$  and  $y$ . I've illustrated the likely range of values for  $y$  as a blue filled circle.

```
In [28]: mkf_internal.plot_correlation_covariance()
```



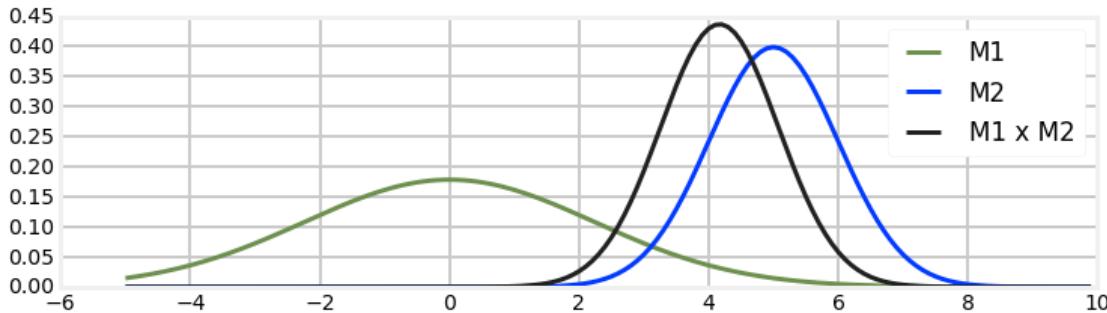
A circle not mathematically correct, but it gets the idea across. We will tackle the mathematics in the next section. For now recognize that we can predict that  $y$  is likely near 12. A value of  $y = -10$  is extremely improbable.

A word about correlation and independence. If variables are independent they can vary separately. If you walk in an open field, you can move in the  $x$  direction (east-west), the  $y$  direction(north-south), or any combination thereof. Independent variables are always also uncorrelated. Except in special cases, the reverse does not hold true. Variables can be uncorrelated, but dependent. For example, consider  $y = x^2$ . Correlation is a linear measurement, so  $x$  and  $y$  are uncorrelated. However, they are obviously dependent on each other.

## 5.6 Multiplying Multidimensional Gaussians

In the previous chapter we incorporated an uncertain measurement with an uncertain estimate by multiplying their Gaussians together. The result was another Gaussian with a smaller variance. If two pieces of uncertain information corroborate each other we should be more certain in our conclusion. The graphs look like this:

In [29]: `mkf_internal.plot_gaussian_multiply()`



The combination of measurements 1 and 2 yields more certainty, so the new Gaussian is taller and narrower - the variance became smaller. The same happens in multiple dimensions with multivariate Gaussians.

Here are the equations for multiplying multivariate Gaussians. They are generated by plugging the Gaussians for the prior and the estimate into Bayes Theorem. I gave you the algebra for the univariate case in the last section of the **Univariate Kalman Filter** chapter. You will not need to remember these equations, as they are computed by Kalman filter equations that will be presented shortly. This computation is also available in FilterPy using the `multivariate_multiply()` method, which you can import from `filterpy.stats`.

$$\begin{aligned}\mu &= \Sigma_2(\Sigma_1 + \Sigma_2)^{-1}\mu_1 + \Sigma_1(\Sigma_1 + \Sigma_2)^{-1}\mu_2 \\ \Sigma &= \Sigma_1(\Sigma_1 + \Sigma_2)^{-1}\Sigma_2\end{aligned}$$

To give you some intuition about this, recall the equations for multiplying univariate Gaussians:

$$\begin{aligned}\mu &= \frac{\sigma_1^2\mu_2 + \sigma_2^2\mu_1}{\sigma_1^2 + \sigma_2^2}, \\ \sigma^2 &= \frac{1}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}} = \frac{\sigma_1^2\sigma_2^2}{\sigma_1^2 + \sigma_2^2}\end{aligned}$$

This looks similar to the equations for the multivariate equations. This will be more obvious if you recognize that matrix inversion, denoted by the  $-1$  power, is like division since  $AA^{-1} = I$ . I will rewrite the inversions as divisions - this is not a mathematically correct thing to do but it does help us see what is going on.

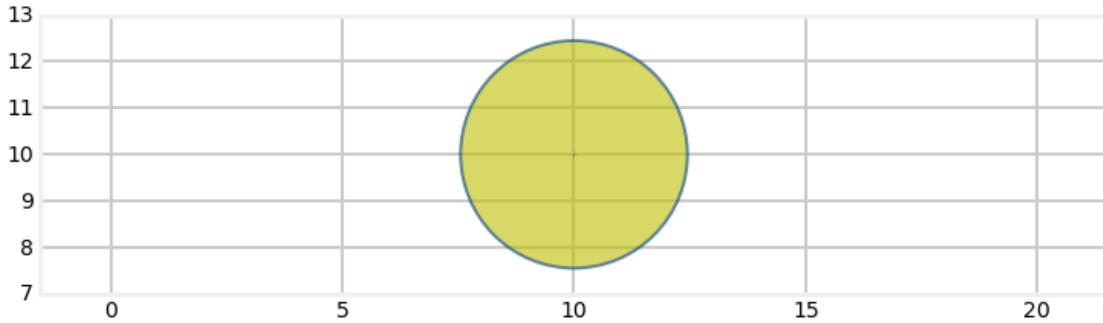
$$\mu \approx \frac{\Sigma_2\mu_1 + \Sigma_1\mu_2}{\Sigma_1 + \Sigma_2}$$

$$\Sigma \approx \frac{\Sigma_1\Sigma_2}{(\Sigma_1 + \Sigma_2)}$$

In this form the relationship between the univariate and multivariate equations is clear.

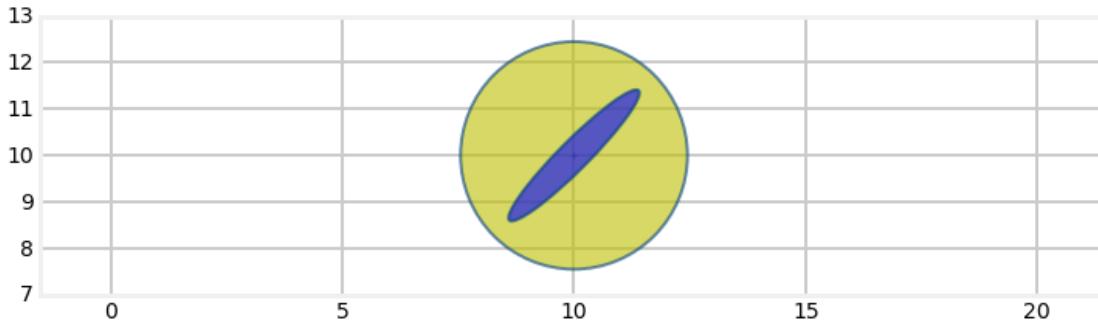
Now let's explore multivariate Gaussians in terms of a concrete example. Suppose that we are tracking an aircraft with two radar systems. I will ignore altitude so I can use two dimensional plots. Radars give us the range and bearing to a target. We start out being uncertain about the position of the aircraft, so the covariance, which is our uncertainty about the position, might look like this. In the language of Bayesian statistics this is our prior.

```
In [30]: P0 = [[6, 0], [0, 6]]
plot_covariance_ellipse((10, 10), P0, fc='y', alpha=0.6)
```



Now suppose that there is a radar to the lower left of the aircraft. Further suppose that the radar's bearing measurement is accurate, but the range measurement is inaccurate. The covariance for the error in the measurement might look like this (plotted in blue):

```
In [31]: P1 = [[2, 1.9], [1.9, 2]]
plot_covariance_ellipse((10, 10), P0, fc='y', alpha=0.6)
plot_covariance_ellipse((10, 10), P1, fc='b', alpha=0.6)
```



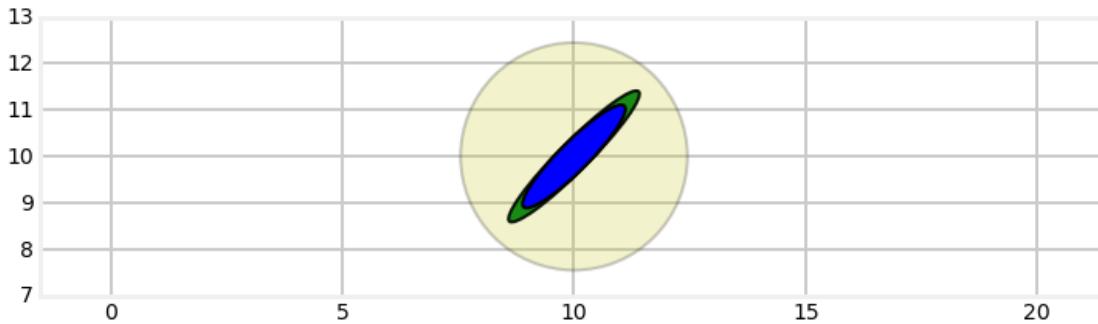
Recall that Bayesian statistics calls this the evidence. The ellipse points towards the radar. It is very long because the range measurement is inaccurate, and the aircraft could be within a considerable distance of the measured range. It is very narrow because the bearing estimate is very accurate and thus the aircraft must be very close to the bearing estimate.

We want to find the posterior - the mean and covariance that results from incorporating the evidence into the prior. As in every other chapter we combine evidence by multiplying them together.

```
In [32]: from filterpy.stats import multivariate_multiply
```

```
P2 = multivariate_multiply((10, 10), P0, (10, 10), P1)[1]

plot_covariance_ellipse((10, 10), P0, fc='y', alpha=0.2)
plot_covariance_ellipse((10, 10), P1, fc='g', alpha=0.9)
plot_covariance_ellipse((10, 10), P2, fc='b')
```



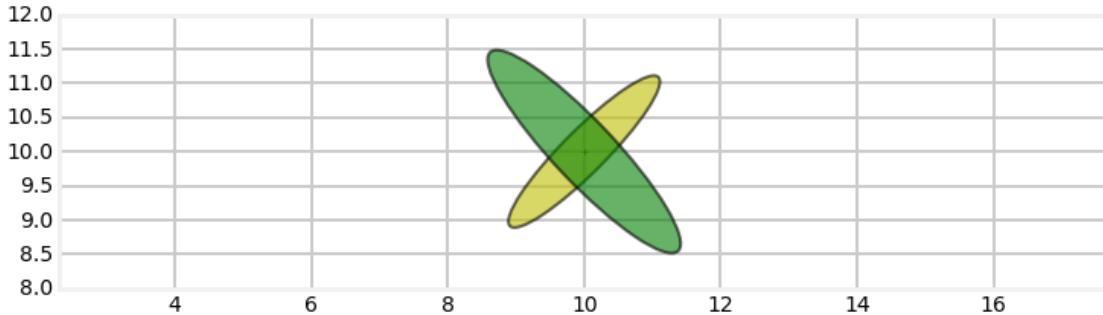
I have plotted the original estimate (prior) in a very transparent yellow, the radar reading in green (evidence), and the finale estimate (posterior) in blue.

The posterior retained the same shape and position as the radar measurement, but is smaller. We've seen this with one dimensional Gaussians. Multiplying two Gaussians makes the variance smaller because we are incorporating more information, hence we are less uncertain. Another point to recognize is that the

covariance shape reflects the physical layout of the aircraft and the radar system. The importance of this will become clear in the next step.

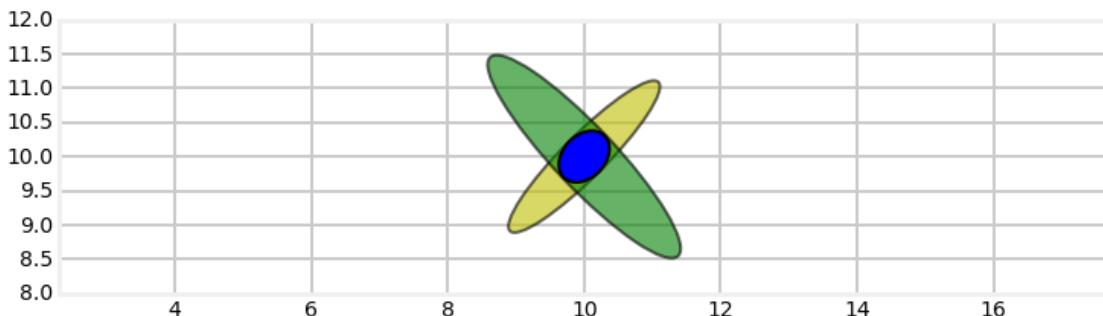
Now lets say we get a measurement from a second radar, this one to the lower right, which I will plot in green against the yellow covariance of our current belief.

```
In [33]: P3 = [[2, -1.9], [-1.9, 2.2]]
plot_covariance_ellipse((10, 10), P2, ec='k', fc='y', alpha=0.6)
plot_covariance_ellipse((10, 10), P3, ec='k', fc='g', alpha=0.6)
```



To incorporate this new information we will multiply the Gaussians together.

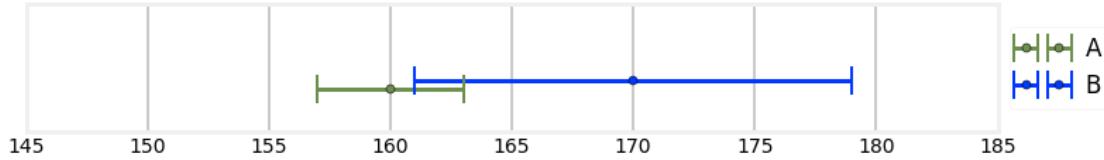
```
In [34]: P4 = multivariate_multiply((10, 10), P2, (10, 10), P3)[1]
plot_covariance_ellipse((10, 10), P2, ec='k', fc='y', alpha=0.6)
plot_covariance_ellipse((10, 10), P3, ec='k', fc='g', alpha=0.6)
plot_covariance_ellipse((10, 10), P4, ec='k', fc='b')
```



The only likely place for the aircraft is where the two ellipses intersect. The shapes reflects the geometry of the problem. This allows us to triangulate on the aircraft, resulting in a very accurate estimate. We didn't explicitly write any code to perform triangulation; it was a natural outcome of multiplying the Gaussians of each measurement together.

Think back to the **g-h Filter** chapter where we displayed the error bars of two weighings on a scale. The estimate must fall somewhere within the region where the error bars overlap. Here the estimate must fall between 161 to 163 pounds.

```
In [35]: import gh_internal as gh
gh.plot_errorbar2()
```

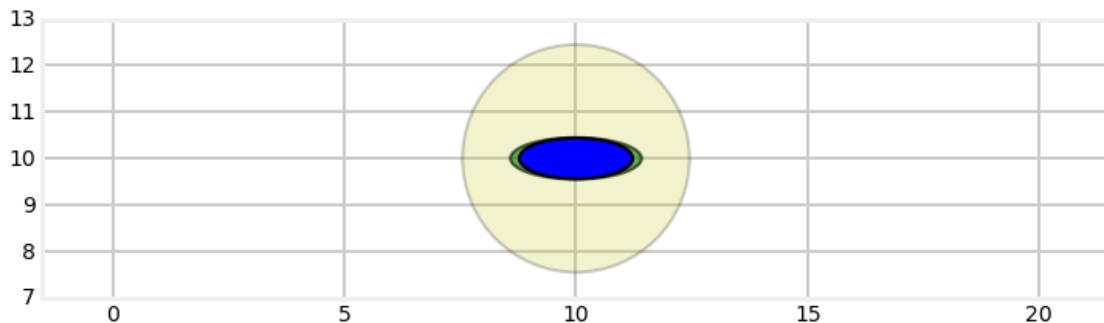


Now suppose the first radar is to the left of the aircraft. I can model the measurement error with

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 0.2 \end{bmatrix}$$

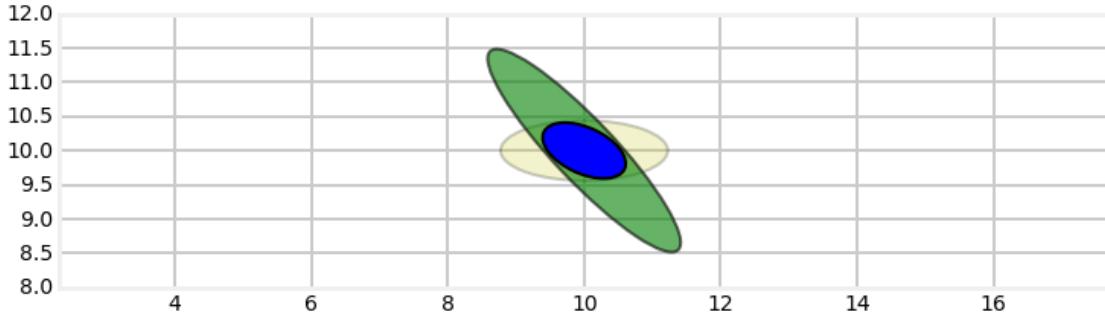
Here we see the result of multiplying the prior with the measurement.

```
In [36]: P1 = [[2, 0], [0, .2]]
P2 = multivariate_multiply((10, 10), P0, (10, 10), P1)[1]
plot_covariance_ellipse((10, 10), P0, ec='k', fc='y', alpha=0.2)
plot_covariance_ellipse((10, 10), P1, ec='k', fc='g', alpha=0.6)
plot_covariance_ellipse((10, 10), P2, ec='k', fc='b')
```



Now we can incorporate the measurement from the second radar system, which we will leave in the same position as before.

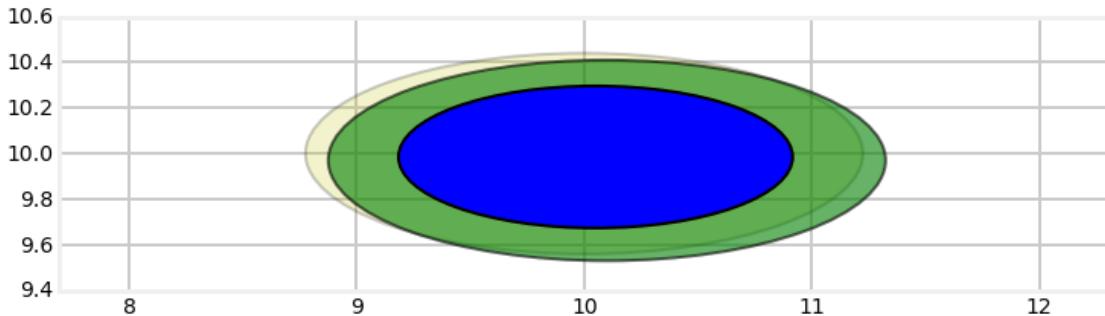
```
In [37]: P3 = [[2, -1.9], [-1.9, 2.2]]
P4 = multivariate_multiply((10, 10), P2, (10, 10), P3)[1]
plot_covariance_ellipse((10, 10), P2, ec='k', fc='y', alpha=0.2)
plot_covariance_ellipse((10, 10), P3, ec='k', fc='g', alpha=0.6)
plot_covariance_ellipse((10, 10), P4, ec='k', fc='b')
```



Our estimate is not as accurate as the previous example. The two radar stations are no longer orthogonal to each other relative to the aircraft's position so the triangulation is not optimal.

For a final example, imagine taking two measurements from the same radar a short time apart. The covariance ellipses will nearly overlap, leaving a very large error in our new estimate:

```
In [38]: P5 = multivariate_multiply((10,10), P2, (10.1, 9.97), P2)
        plot_covariance_ellipse((10, 10), P2, ec='k', fc='y', alpha=0.2)
        plot_covariance_ellipse((10.1, 9.97), P2, ec='k', fc='g', alpha=0.6)
        plot_covariance_ellipse(P5[0], P5[1], ec='k', fc='b')
```

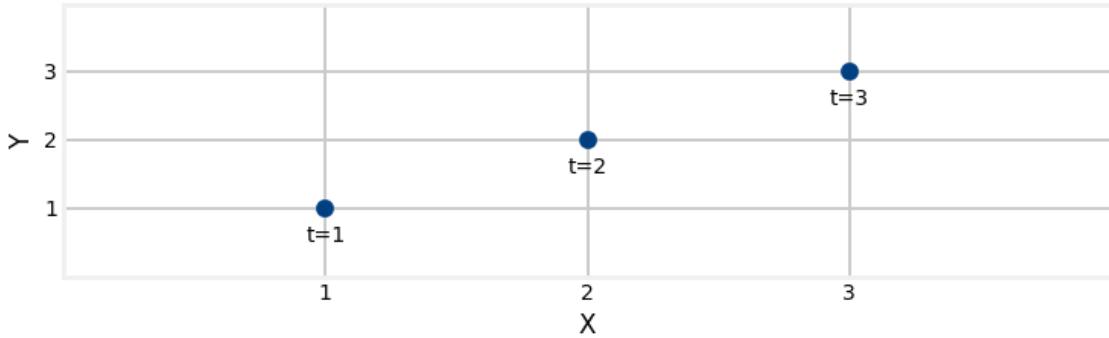


## 5.7 Hidden Variables

You can already see why a multivariate Kalman filter can perform better than a univariate one. The last section demonstrated how we can use correlations between variables to significantly improve our estimates. We can take this much further. This section contains the key insight to this chapter, so read carefully.

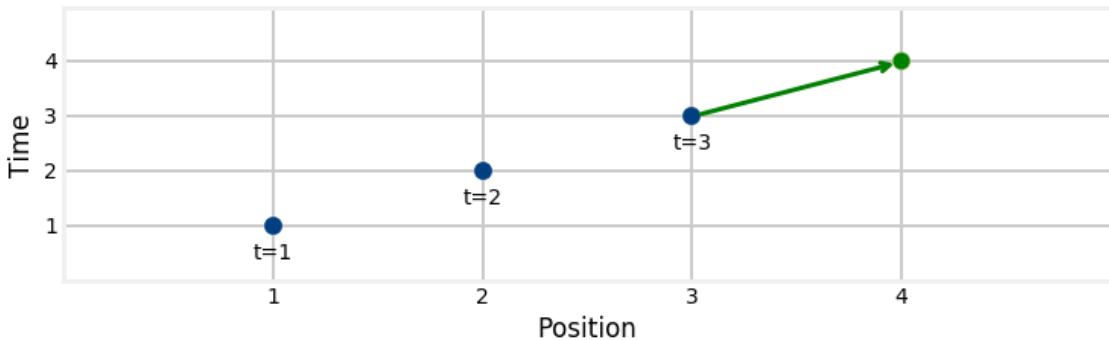
Let's say we are tracking an aircraft and we get the following data for the  $x$  and  $y$  coordinates at time  $t=1, 2$ , and  $3$  seconds. What does your intuition tell you the value of  $x$  will be at time  $t=4$  seconds?

```
In [39]: import mkf_internal
        mkf_internal.show_position_chart()
```



It appears that the aircraft is flying in a straight line and we know that aircraft cannot turn on a dime. The most reasonable guess is that at  $t=4$  the aircraft is at  $(4,4)$ . I will depict that with a green arrow.

In [40]: `mkf_internal.show_position_prediction_chart()`



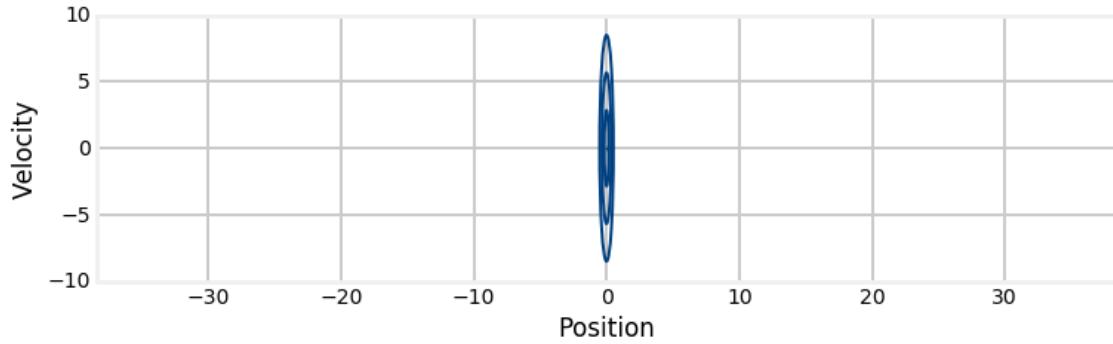
You made this inference because you inferred a constant velocity for the airplane. The reasonable assumption is that the aircraft is moving one unit each in x and y per time step.

Think back to the **g-h Filter** chapter when we were trying to improve the weight predictions of a noisy scale. We incorporated weight gain into the equations because it allowed us to make a better prediction of the weight the next day. The g-h filter uses the  $g$  parameter to scale the amount of significance given to the current weight measurement, and the  $h$  parameter scaled the amount of significance given to the weight gain.

We are going to do the same thing with our Kalman filter. After all, the Kalman filter is a form of a g-h filter. In this case we are tracking an airplane, so instead of weight and weight gain we need to track position and velocity. Weight gain is the derivative of weight, and of course velocity is the derivative of position. It's impossible to plot and understand the 4D chart that would be needed to plot x and y and their respective velocities so let's do it for x, knowing that the math generalizes to more dimensions.

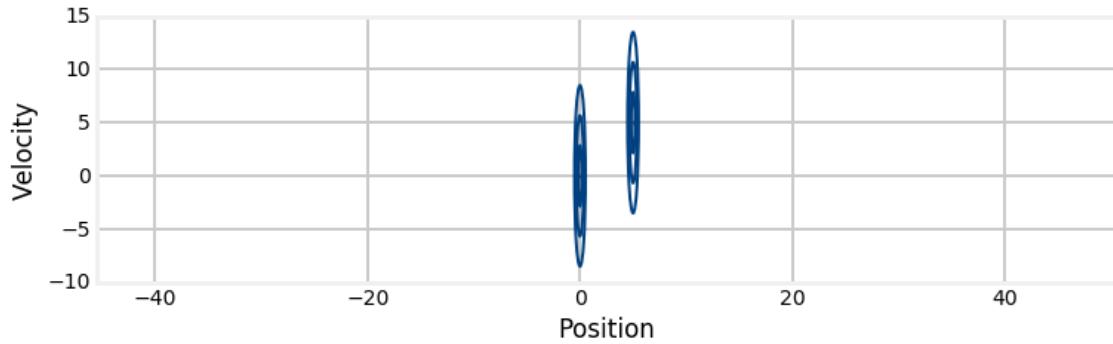
At time 1 we might be fairly certain about the position ( $x=0$ ) but have no idea about the velocity. We can plot that with a covariance matrix like this. The narrow width expresses our relative certainty about position, and the tall height expresses our lack of knowledge about velocity.

In [41]: `mkf_internal.show_x_error_chart(1)`



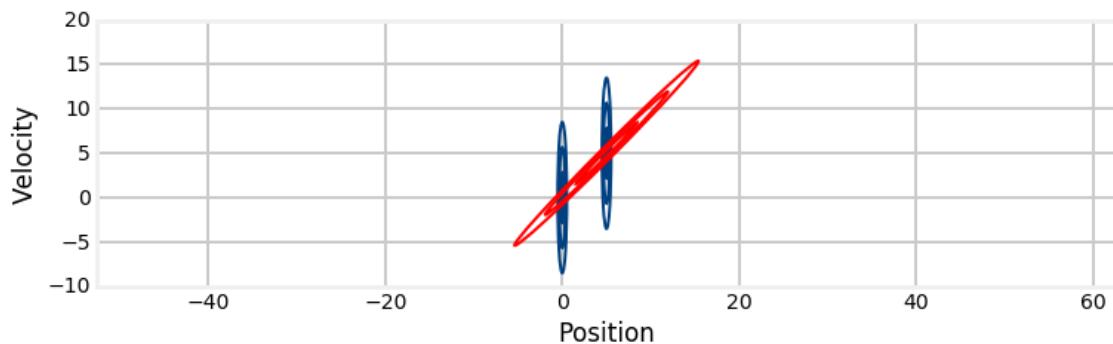
Now after one second we get a position update of  $x=5$ .

In [42]: `mkf_internal.show_x_error_chart(2)`



This implies that our velocity is roughly 5 m/s. But of course position and velocity are correlated. If the velocity is 5 m/s the position would be 5, but if the velocity was 10 m/s the position would be 10. So let's draw a covariance matrix in red showing the relationship between the position and velocity.

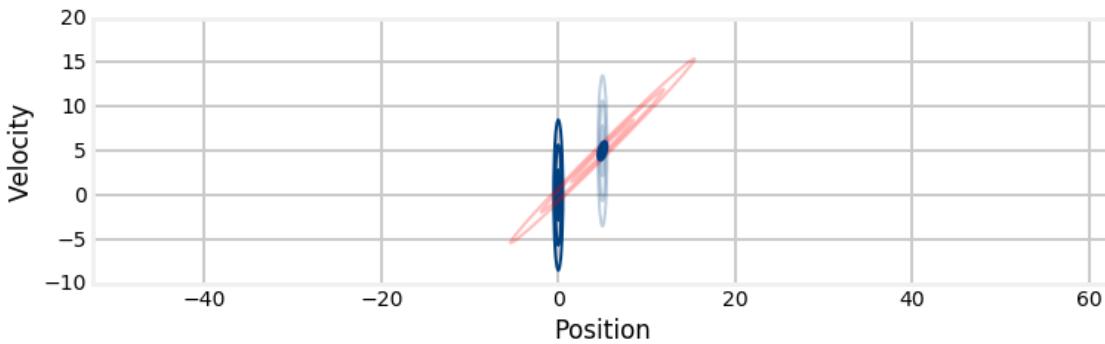
In [43]: `mkf_internal.show_x_error_chart(3)`



It won't be clear until the next chapter how I calculate this covariance. Ignore the calculation, and think about what this implies. We have no easy way to say where the object really is because we are so uncertain about the velocity. Hence the ellipse stretches very far in the x-axis. Our uncertainty in velocity of course means it is also very spread in the y-axis. But as I said in the last paragraph, position and velocity is correlated. If the velocity is 5 m/s the next position would be 5, and if the velocity is 10 the next position would be 10. They are extremely correlated, so the ellipse must be very narrow.

This superposition of the two covariances is where the magic happens. The only reasonable estimate at time  $t=1$  (where position=5) is roughly the intersection between the two covariance matrices! More exactly, we can use the math from the last section and multiply the two covariances together. From a Bayesian point of view we multiply the prior with the probability of the evidence (the likelihood) to get the posterior. If we multiply the position covariance with the velocity covariance using the Bayesian equations we get the result shown in the next chart.

```
In [44]: mkf_internal.show_x_error_chart(4)
```



We can see that the new covariance (the posterior) lies at the intersection of the position covariance and the velocity covariance. It is slightly tilted, showing that there is some correlation between the position and velocity. Far more importantly, it is much smaller than either the position or velocity covariances. In the previous chapter our variance would get smaller each time we performed an `update()` because the previous estimate was multiplied by the new measurement. The same thing happens here. However, the amount by which the covariance got smaller by is much larger in this chapter. This is because we are using two different kinds of information which are nevertheless correlated. Knowing the velocity approximately and the position approximately allows us to very quickly hone in on the correct answer.

This is a key point in Kalman filters, so read carefully! Our sensor is only detecting the position of the aircraft (how doesn't matter). This is called an observed variable. It does not have a sensor that provides velocity. But based on the position estimates we can compute velocity. We call the velocity a hidden variable. Hidden means what it sounds like - there is no sensor that is measuring velocity, thus its value is hidden from us. We are able to use the correlation between position and velocity to infer its value very accurately.

To round out the terminology there are also unobserved variables. For example, the aircraft's state includes things such as as heading, engine RPM, weight, color, the first name of the pilot, and so on. We cannot sense these directly using the position sensor so they are not observed. There is no way to infer them from the sensor measurements and correlations (red planes don't go faster than white planes), so they are not hidden. Instead, they are unobservable. If you include an unobserved variable in your filter state the estimate for that variable will be nonsense.

What makes this possible? Imagine for a moment that we superimposed the velocity from a different airplane over the position graph. Clearly the two are not related, and there is no way that combining the two could possibly yield any additional information. In contrast, the velocity of this airplane tells us something very

important - the direction and speed of travel. So long as the aircraft does not alter its velocity the velocity allows us to predict where the next position is. After a relatively small amount of error in velocity the probability that it is a good match with the position is very small. Think about it - if you suddenly change direction your position is also going to change a lot. If the measurement of the position is not in the direction of the velocity change it is very unlikely to be true. The two are correlated, so if the velocity changes so must the position, and in a predictable way.

It is important to understand that we are taking advantage of the fact that velocity and position are correlated. We get a rough estimate of velocity from the distance and time between two measurement, and use Bayes theorem to and produce very accurate estimates after only a few observations. Please reread this section if you have any doubts. If you do not understand this you will quickly find it impossible to reason about what you will learn in the following chapters.

## 5.8 Summary

We have taken advantage of the geometry and correlations of the system to produce a very accurate estimate. The math does not care whether we are working with two positions, or a position and a correlated velocity, or if these are spatial dimensions. If floor space is correlated to house price you can write a Kalman filter to track house prices. If age is correlated to disease incidence you can write a Kalman filter to track diseases. If the zombie population is inversely correlated with the number of shotguns then you can write a Kalman filter to track zombie populations. I showed you this in terms of geometry and talked about triangulation. That was just to build your intuition. You can write a Kalman filter for state variables that have no geometric representation, such as filters for stock prices or milk production of cows (I received an email from someone tracking milk production!) Get used to thinking of these as Gaussians with correlations. If we can express our uncertainties as a multidimensional Gaussian we can then multiply the prior with the likelihood and get a much more accurate result.

## 5.9 References

- [1] <http://docs.scipy.org/doc/scipy/reference/tutorial/stats.html>
- [2] FilterPy library. Roger Labbe. <https://github.com/rlabbe/filterpy>



# Chapter 6

## Multivariate Kalman Filters - Working with Multiple State Variables

### 6.1 Introduction

In this chapter I am purposefully glossing over many aspects of the mathematics behind Kalman filters. Some things I show you will only work for special cases, others will be ‘magical’ - it will not be clear how I derived a certain result. I prefer that you develop an intuition for how these filters work through several worked examples. If I started with rigorous, generalized equations you would be left scratching your head about what all these terms mean and how you might apply them to your problem. In later chapters I will provide a more rigorous mathematical foundation, and at that time I will have to either correct approximations that I made in this chapter or provide additional information that I did not cover here.

To make this possible we will restrict ourselves to a subset of problems which we can describe with Newton’s equations of motion. These filters are called discretized continuous-time kinematic filters. In the **Kalman Filter Math** chapter we will develop the math required for solving any kind of dynamic system.

### 6.2 Newton’s Equations of Motion

Newton’s equations of motion are: given a constant velocity  $v$  of a system we can compute its position  $x$  after time  $t$  with:

$$x = vt + x_0$$

For example, if we start at position 13, our velocity is 10 m/s, and we travel for 12 seconds our final position is  $133 = (10 \times 12) + 13$ .

We can incorporate constant acceleration with this equation

$$x = \frac{1}{2}at^2 + v_0t + x_0$$

And if we assume constant jerk we get

$$x = \frac{1}{6}jt^3 + \frac{1}{2}a_0t^2 + v_0t + x_0$$

As a reminder, we generate these equations by integrating a differential equation. Given a constant velocity  $v$  we can compute the distance traveled over time with the equation

$$\begin{aligned} v &= \frac{dx}{dt} \\ dx &= v dt \\ \int_{x_0}^x dx &= \int_0^t v dt \\ x - x_0 &= vt - 0 \\ x &= vt + x_0 \end{aligned}$$

When you design a Kalman filter you start with a system of differential equations that describe the dynamics of the system. Most systems of differential equations do not easily integrate in this way. We start with Newton's equation because we can integrate and get a closed form solution, which makes the Kalman filter easier to design. An added benefit is that Newton's equations are the right equations to use to track objects, which is one of the main uses of Kalman filters.

## 6.3 Kalman Filter Algorithm

Now that we have that under our belts we can show how the Kalman filter works. The algorithm is the same Bayesian filter algorithm that we have used in every chapter. The update step is slightly more complicated, but I will explain why when we get to it.

### Initialization

1. Initialize the state of the filter
2. Initialize our belief in the state

### Predict

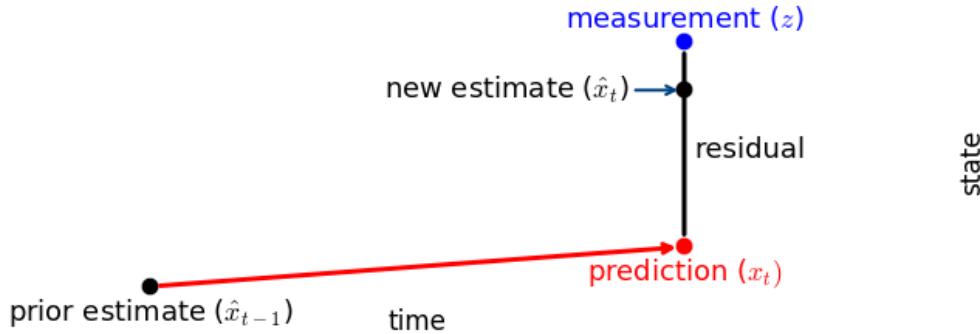
1. Use system behaviors to predict its state at the next time step
2. Adjust belief to account for the uncertainty in prediction

### Update

1. Get a measurement and associated belief about its accuracy
2. Compute residual between estimated state and measurement
3. Compute scaling factor based on whether the measurement or prediction is more accurate
4. Set state between the prediction and measurement based on scaling factor
5. Update belief in the state based on how certain we are in the measurement

As a reminder, here is a graphical depiction of the algorithm:

```
In [2]: from mkf_internal import show_residual_chart
show_residual_chart()
```



## 6.4 Tracking a Dog

Let's go back to our tried and true problem of tracking our dog. This time we will include the fundamental insight of the previous chapter and use hidden variables to improve our estimates. I could start with the math, but instead let's implement a filter, learning as we go. On the surface the math is different and perhaps more complicated than the previous chapters, but the ideas are all the same - we are just multiplying and adding Gaussians.

We start by writing a simulation for the dog. The simulation will run for `count` steps, moving the dog forward approximately 1 meter for each step. At each step the velocity will vary according to the process variance `process_var`. After updating the position we compute a measurement with an assumed sensor variance of `z_var`. The function returns an NumPy array of the positions and another of the measurements.

```
In [3]: import math
        import numpy as np
        from numpy.random import randn

def compute_dog_data(z_var, process_var, count=1, dt=1.):
    x, vel = 0., 1.
    z_std = math.sqrt(z_var)
    p_std = math.sqrt(process_var)
    xs, zs = [], []
    for _ in range(count):
        v = vel + (randn() * p_std * dt)
        x += v*dt
        xs.append(x)
        zs.append(x + randn() * z_std)
    return np.array([xs, zs]).T
```

I have programmed the equations of the Kalman filter into the `KalmanFilter` class in FilterPy. You will import it with:

```
from filterpy.kalman import KalmanFilter
```

I need to introduce a lot of new material to implement the tracking filter, and it is easy to get lost in the details. We are using multivariate Gaussians so we need to use matrices and linear algebra to implement the filter. Keep in mind that we are doing the same things as in the preceding chapters:

Use a Gaussian to represent our estimate of the state and error

Use a Gaussian to represent our sensor and its error

Use a process model to predict the next state

Incorporate the measurement into the prediction using Bayes Theorem.

### 6.4.1 Initialization - Choose State Variables and Set Initial Conditions

#### Step 1: Design State Variable as a Multivariate Gaussian

We previously tracked a dog in one dimension by using a Gaussian. The mean ( $\mu$ ) represented the most likely position, and the variance ( $\sigma^2$ ) represented the probability distribution of the position. The position is the state of the system, and we call  $\mu$  the state variable.

In this problem we will be tracking both the position and velocity of the dog. It is important to understand that this is a design choice with implications and assumptions that we are not yet prepared to explore. For example, we could also track acceleration, or even jerk.

State variables can either be observed variables - directly measured by a sensor, or hidden variables - inferred from the observed variables. For our dog tracking problem the sensor reads position, so position is observed and velocity is hidden.

In the univariate chapter we represented the dog's position with a scalar value (e.g.  $\mu = 3.27$ ). In the last chapter we learned to use a multivariate Gaussian. For example, if we wanted to specify a position of 10.0 m and a velocity of 4.5 m/s, we would write:

$$\mu = \begin{bmatrix} 10.0 \\ 4.5 \end{bmatrix}$$

The Kalman filter is implemented using linear algebra. We use an  $n \times 1$  matrix (called a vector) to store  $n$  state variables. For the dog tracking problem, we use  $x$  to denote position, and the first derivative of  $x$ ,  $\dot{x}$ , for velocity. Kalman filter equations use  $\mathbf{x}$  for the state, so we define  $\mathbf{x}$  as:

$$\mathbf{x} = \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

We use  $\mathbf{x}$  instead of  $\mu$ , but recognize this is the mean of the multivariate Gaussian.

Another way to write this is  $\mathbf{x} = [x \ \dot{x}]^\top$  because the transpose of a row vector is a column vector. This notation is easier to use in text because it takes less vertical space.

$\mathbf{x}$  and the position  $x$  coincidentally have the same name. If we were tracking the dog in the y-axis we would write  $\mathbf{x} = [y \ \dot{y}]^\top$ , not  $\mathbf{y} = [y \ \dot{y}]^\top$ .  $\mathbf{x}$  is the standard name for the state variable used in the Kalman filter literature and we will not vary it to give it a more meaningful name. This consistency in naming allows us to communicate with our peers.

Let's code this. FilterPy's `KalmanFilter` class contains a variable `x` for the state variable. Initialization of `x` is as simple as

```
In [4]: from filterpy.kalman import KalmanFilter
kf = KalmanFilter(dim_x=2, dim_z=1)

kf.x = np.array([[10.0],
                 [4.5]])
print(kf.x)
```

```
[[ 10. ]
 [ 4.5]]
```

I often use the transpose trick in my code to turn a row matrix into a vector:

```
In [5]: kf.x = np.array([[10.0, 4.5]]).T
        print(kf.x)
```

```
[[ 10. ]
 [ 4.5]]
```

However, NumPy recognizes 1D arrays as vectors, so I can simplify this line to use a 1D array.

```
In [6]: kf.x = np.array([10.0, 4.5])
```

The other multivariate Gaussian variable is the covariance matrix  $\Sigma$ . Kalman filter equations typically use the symbol  $\mathbf{P}$ . In the univariate Kalman filter we specified an initial value for  $\sigma^2$ , and then the filter took care of updating its value as measurements were added to the filter. The same thing happens in the multidimensional Kalman filter. We specify an initial value for  $\mathbf{P}$  and the filter updates it during each epoch.

We need to set the variances to reasonable values. For example, we may choose  $\sigma_{\text{pos}}^2 = 500m^2$  if we are quite uncertain about the initial position. Top speed for a dog is around 21 m/s, so in the absence of any other information about the velocity we can set  $3\sigma_{\text{vel}} = 21$ , or  $\sigma_{\text{vel}}^2 = 49$ .

In the last chapter we showed that the position and velocities are correlated. But how correlated are they for a dog? I have no idea. As we will see the filter computes this for us, so I initialize the covariances to zero. Of course if you know the covariances you should use them.

Recall that the diagonals of the covariance matrix contains the variance of each variable, and the off-diagonal elements contains the covariances. Thus we have:

$$\mathbf{P} = \begin{bmatrix} 500 & 0 \\ 0 & 49 \end{bmatrix}$$

I can code this using the function `numpy.diag` which creates a diagonal matrix from the values for the diagonal. Recall from linear algebra that a diagonal matrix is one with zeros in the off-diagonal elements.

```
In [7]: kf.P = np.diag([500, 49])
        kf.P
```

```
Out[7]: array([[ 500.,    0.],
               [    0.,   49.]])
```

I could have been explicit. This creates the same value:

```
In [8]: kf.P = np.array([[500., 0.],
                      [0., 49.]])
        kf.P
```

```
Out[8]: array([[ 500.,    0.],
               [    0.,   49.]])
```

We are done. We've expressed the state of the filter as a multivariate Gaussian and implemented it in code.

### 6.4.2 Predict Step

The next step is designing the process model which tells the filter how to predict the state (mean and covariance) of the system for the next time step. We do this with a set of equations that describe the physical model of the system.

#### Step 2: Design the State Transition Function

Previously we modeled the dog's motion with

$$x = v\Delta t + x_0$$

and implemented it with

```
def predict(pos, variance, movement, movement_variance):
    return (pos + movement, variance + movement_variance)
```

We will do the same thing in this chapter, using multivariate Gaussians instead of univariate Gaussians. You might imagine this sort of implementation:

$$\mathbf{x} = \begin{bmatrix} 5.4 \\ 4.2 \end{bmatrix}, \dot{\mathbf{x}} = \begin{bmatrix} 1.1 \\ 0. \end{bmatrix} \mathbf{x} = \mathbf{x} + \dot{\mathbf{x}}$$

But we need to generalize this. The Kalman filter equations work with any linear system, not just Newtonian ones. Maybe the system you are filtering is the plumbing system in a chemical plant, and the amount of flow in a given pipe is determined by a linear combination of the settings of different valves.

$$\begin{aligned} \text{pipe}_1 &= 0.134(\text{valve}_1) + 0.41(\text{valve}_2 - \text{valve}_3) + 1.34 \\ \text{pipe}_2 &= 0.210(\text{valve}_2) - 0.62(\text{valve}_1 - \text{valve}_5) + 1.86 \end{aligned}$$

Linear algebra has a powerful way to express systems of equations. Take this system

$$\begin{cases} 2x + 3y = 8 \\ 3x - y = 1 \end{cases}$$

We can put this in matrix form by writing:

$$\begin{bmatrix} 2 & 3 \\ 3 & -1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 8 \\ 1 \end{bmatrix}$$

If you perform the matrix multiplication in this equation the result will be the two equations above. In linear algebra we would write this as  $\mathbf{Ax} = \mathbf{B}$ , where

$$\mathbf{A} = \begin{bmatrix} 2 & 3 \\ 3 & -1 \end{bmatrix}, \mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}, \mathbf{B} = \begin{bmatrix} 8 \\ 1 \end{bmatrix}$$

We call the set of equations that describe how the system behaves the process model. We use the process model to perform the innovation, because the equations tell us what the next state will be given the current state. Kalman filters implement this using the linear equation, where  $\bar{\mathbf{x}}$  is the prior, or predicted state:

$$\bar{\mathbf{x}} = \mathbf{Fx}$$

Our job as Kalman filters designers is to specify  $\mathbf{F}$  such that  $\bar{\mathbf{x}} = \mathbf{Fx}$  performs the innovation (prediction) for our system. To do this we need one equation for each state variable. In our problem  $\mathbf{x} = [x \ \dot{x}]^T$ , so we need one equation to compute  $\bar{x}$  and another to compute  $\bar{\dot{x}}$ . We already know the equation for the position innovation:

$$\bar{x} = \dot{x}\Delta t + x$$

What is our equation for velocity? We have no predictive model for how our dog's velocity will change over time. In this case we assume that it remains constant between innovations. Of course this is not exactly true, but so long as the velocity doesn't change too much over each innovation you will see that the filter performs very well. So we say

$$\bar{\dot{x}} = \dot{x}$$

This gives us the process model for our system

$$\begin{cases} \bar{x} = \dot{x}\Delta t + x \\ \bar{\dot{x}} = \dot{x} \end{cases}$$

We need to express this set of equations in the form  $\bar{\mathbf{x}} = \mathbf{Fx}$ . Rearranging terms makes it easier to see what to do.

$$\begin{cases} \bar{x} = 1x + \Delta t \dot{x} \\ \bar{\dot{x}} = 0x + 1\dot{x} \end{cases}$$

We can rewrite this in matrix form as

$$\begin{bmatrix} \bar{x} \\ \bar{\dot{x}} \end{bmatrix} = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

$$\bar{\mathbf{x}} = \mathbf{Fx}$$

$\mathbf{F}$  is often called the state transition function. In the `KalmanFilter` class we implement the state transition function with

```
In [9]: dt = 0.1
        kf.F = np.array([[1, dt],
                          [0, 1]])
        print(kf.F)

[[ 1.   0.1]
 [ 0.   1. ]]
```

Let's test this! `KalmanFilter` has a `predict` method that performs the prediction by computing  $\bar{\mathbf{x}} = \mathbf{Fx}$ . Let's call it and see what happens. We've set the position to 10.0 and the velocity to 4.5 meter/sec. We've defined `dt = 0.1`, which means the time step is 0.1 seconds, so we expect the new position to be 10.45 meters after the innovation. The velocity should be unchanged.

```
In [10]: kf.x = np.array([10.0, 4.5])
        kf.P = np.diag([500, 49])

# Q is amount of noise in the process,
```

```
# we will learn about it later
kf.Q = 0
kf.predict()
print(kf.x)
```

```
[ 10.45  4.5 ]
```

This worked. Note that the code does not distinguish between the prior and posterior in the variable names, so after calling `predict()` the prior  $\bar{x}$  is stored in `KalmanFilter.x`. If we call `predict()` several times in a row the value will be updated each time.

```
In [11]: kf.predict()
        print('x =', kf.x)
        kf.predict()
        print('x =', kf.x)

x = [ 10.9   4.5]
x = [ 11.35   4.5 ]
```

### Step 3: Design the Process Noise Matrix

`KalmanFilter.predict()` computes both the mean and covariance of the innovation. This is the value of  $\mathbf{P}$  after three innovations (predictions), which we denote  $\bar{\mathbf{P}}$  in the Kalman filter equations.

```
In [12]: kf.P
```

```
Out[12]: array([[ 504.41,   14.7 ],
               [ 14.7 ,   49. ]])
```

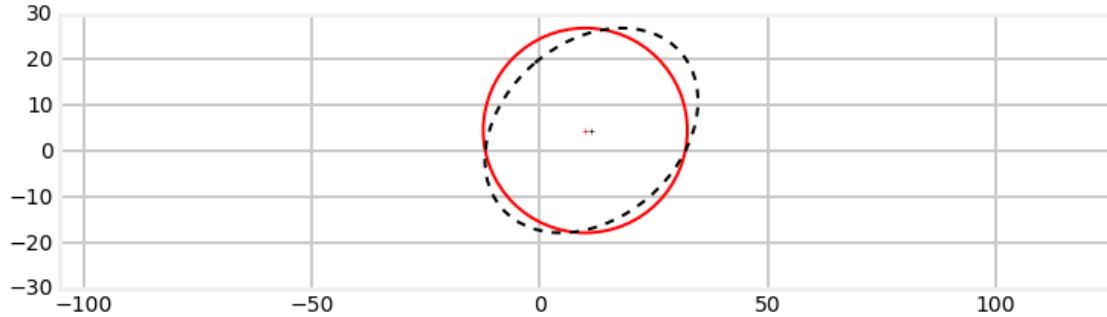
Inspecting the diagonals shows us that the position variance got larger. We've performed three prediction steps with no measurements, and our uncertainty grew. The off-diagonal elements became non-zero - the Kalman filter detected a correlation between position and velocity! The variance of the velocity did not change.

Here I plot the covariance before and after the prediction. The initial value is in solid red, and the prior (prediction) is in dashed black. I've altered the covariance and time step to better illustrate the change.

```
In [13]: from filterpy.stats import plot_covariance_ellipse

dt = 0.3
kf.F = np.array([[1, dt], [0, 1]])
kf.x = np.array([10.0, 4.5])
kf.P = np.diag([500, 500])
plot_covariance_ellipse(kf.x, kf.P, edgecolor='r')

kf.predict()
plot_covariance_ellipse(kf.x, kf.P, edgecolor='k', ls='dashed')
```



You can see that the center of the ellipse shifted by a small amount (from 10 to 10.90) because the position changed. The ellipse also elongated, showing the correlation between position and velocity. How does the filter compute new values for  $\mathbf{P}$ , and what is it based on? It's a little too early to discuss this, but recall that in every filter so far the predict step entailed a loss of information. The same is true here. I will give you the details once we have covered a bit more ground.

### Process Noise

A quick review on process noise. A car is driving along the road with the cruise control on; it should travel at a constant speed. We model this with  $x_k = \dot{x}_k \Delta t + x_{k-1}$ . However, it is affected by a number of unknown factors. The cruise control cannot perfectly maintain a constant velocity. Winds affect the car, as do hills and potholes. Passengers roll down windows, changing the drag profile of the car.

We can model this system with the differential equation

$$\dot{\mathbf{x}} = f(\mathbf{x}) + w$$

where  $f(\mathbf{x})$  models the state transition and  $w$  is white process noise.

We will learn how to go from a differential equation to the Kalman filter matrices in the **Kalman Filter Math** chapter. For now you just need to know that to account for the noise in the system the filter adds a process noise covariance matrix  $\mathbf{Q}$  to the covariance  $\mathbf{P}$ . We do not add anything to  $\mathbf{x}$  because the noise is white - which means that the mean of the noise will be 0. If the mean is 0,  $\mathbf{x}$  will not change.

The univariate Kalman filter used `variance = variance + process_noise` to compute the variance for the variance of the prediction step. The multivariate Kalman filter does the same, essentially  $\mathbf{P} = \mathbf{P} + \mathbf{Q}$ . I say ‘essentially’ because there are other terms unrelated to noise in the covariance equation that we will see later.

Deriving the process noise matrix can be quite demanding, and we will put it off until the Kalman math chapter. For now know that  $\mathbf{Q}$  equals the expected value of the white noise  $w$ , computed as  $\mathbf{Q} = \mathbb{E}[ww^T]$ . In this chapter we will focus on building an intuitive understanding on how modifying this matrix alters the behavior of the filter.

FilterPy provides functions which compute  $\mathbf{Q}$  for the kinematic problems of this chapter. `Q_discrete_white_noise` takes 3 parameters. `dim`, which specifies the dimension of the matrix, `dt`, which is the time step in seconds, and `var`, the variance in the noise. Briefly, it discretizes the noise over the given time period under assumptions that we will discuss later. This code computes  $\mathbf{Q}$  for white noise with a variance of 2.35 and a time step of 0.1 seconds:

```
In [14]: from filterpy.common import Q_discrete_white_noise
kf.Q = Q_discrete_white_noise(dim=2, dt=0.1, var=2.35)
kf.Q
```

```
Out[14]: array([[ 0.    ,  0.001],
   [ 0.001,  0.024]])
```

#### Step 4: Design the Control Function

The Kalman filter does not just filter data, it allows us to incorporate control inputs for systems like robots and airplanes. Suppose we are controlling a robot. At each time step we would send control signals to the robot based on its current position vs desired position. Kalman filter equations incorporate that knowledge into the filter equations, creating a predicted position based both on current velocity and control inputs to the drive motors. Remember, we never throw information away.

For a linear system the effect of control inputs can be described as a set of linear equations, which we can express with linear algebra as

$$\mathbf{B}\mathbf{u}$$

Here  $\mathbf{u}$  is the control input, and  $\mathbf{B}$  is the control input model or control function. For example,  $\mathbf{u}$  might be a voltage controlling how fast the wheel's motor turns, and multiplying by  $\mathbf{B}$  yields  $[\frac{x}{\dot{x}}]$ .

Therefore the complete Kalman filter equation for the prior mean is

$$\bar{\mathbf{x}} = \mathbf{F}\mathbf{x} + \mathbf{B}\mathbf{u}$$

and this is the equation that is computed when you call `KalmanFilter.predict()`.

Your dog may be trained to respond to voice commands. All available evidence suggests that my dog has no control inputs whatsoever, so I set  $\mathbf{B}$  to zero. In Python we write

```
In [15]: kf.B = 0. # my dog doesn't listen to me!
```

#### Prediction: Summary

Your job as a designer is to specify the matrices for

- $\mathbf{x}$ : the state variables and initial value
- $\mathbf{P}$ : the covariance matrix initial value
- $\mathbf{F}$ : the state transition function
- $\mathbf{B}, \mathbf{u}$ : Optionally, the control input and function
- $\mathbf{Q}$ : the process noise matrix

### 6.4.3 Update Step

Now we can implement the update step of the filter. The good news is that you only have to supply two more matrices, and they are easy to understand. The bad news is that the Kalman filter equations themselves are hard to understand. For now I'll leave the equations unexplained so we can get the full Kalman filter working. After that we'll look at the equations.

#### Step 5: Design the Measurement Function

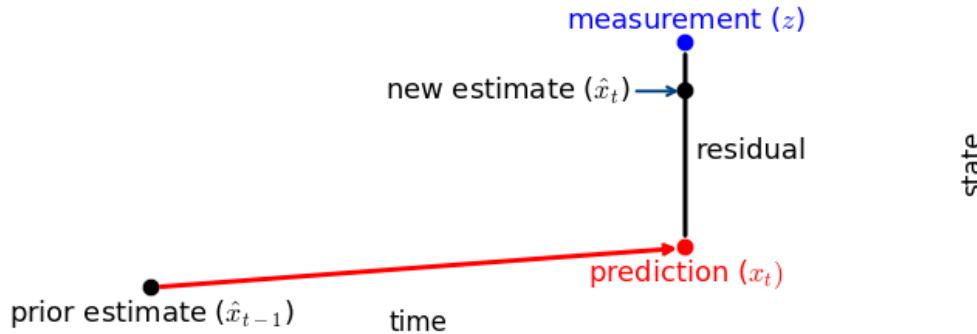
The Kalman filter computes the update step in what we call measurement space. We mostly ignored this issue in the univariate chapter because of the complication it adds. We tracked our dog's position using a

sensor that reported his position. Computing the residual was easy - subtract the filter's predicted position from the measurement:

$$\text{residual} = \text{measured position} - \text{estimated position}$$

We need to compute the residual because we scale it by a scaling factor to get the new estimate.

```
In [16]: from mkf_internal import show_residual_chart
show_residual_chart()
```



What would happen if we were trying to track temperature using a thermometer that outputs a voltage corresponding to the temperature reading? The equation for the residual computation would be nonsense; you can't subtract a temperature from a voltage.

$$\text{residual} = \text{voltage} - \text{temperature} \quad (\text{NONSENSE!})$$

We need to convert the temperature into a voltage so we can perform the subtraction. For the thermometer we might write:

```
CELSIUS_TO_VOLTS = 0.21475
residual = measurement - (CELSIUS_TO_VOLTS * predicted_state)
```

The Kalman filter generalizes this problem by having you supply a measurement function that converts a state into a measurement.

Why are we working in measurement space? Why not work in state space by converting the voltage into a temperature, allowing the residual to be a difference in temperature?

We cannot do that because most measurements are not invertible. The state for the tracking problem contains the hidden variable  $\dot{x}$ . There is no way to convert a measurement of position into a state containing velocity. On the other hand, it is trivial to convert a state containing position and measurement into a equivalent "measurement" containing only position. We have to work in measurement space to make the computation of the residual possible.

Both the measurement  $\mathbf{z}$  and state  $\mathbf{x}$  are vectors so we need to use a matrix to perform the conversion. The Kalman filter equation that performs this step is:

$$\mathbf{y} = \mathbf{z} - \mathbf{H}\bar{\mathbf{x}}$$

where  $\mathbf{y}$  is the residual,  $\bar{\mathbf{x}}$  is the prior,  $\mathbf{z}$  is the measurement, and  $\mathbf{H}$  is the measurement function. So we take the prior, convert it to a measurement by multiplying it with  $\mathbf{H}$ , and subtract that from the measurement. This gives us the difference between our prediction and measurement in measurement space!

We need to design  $\mathbf{H}$  so that  $\mathbf{H}\bar{\mathbf{x}}$  yields a measurement. For this problem we have a sensor that measures position, so  $\mathbf{z}$  will be a one variable vector:

$$\mathbf{z} = [z]$$

The residual equation will have the form

$$\begin{aligned}\mathbf{y} &= \mathbf{z} - \mathbf{H}\bar{\mathbf{x}} \\ [y] &= [z] - [? \ ?] \begin{bmatrix} x \\ \dot{x} \end{bmatrix}\end{aligned}$$

$\mathbf{H}$  has to be a 1x2 matrix for  $\mathbf{H}\bar{\mathbf{x}}$  to be 1x1. Recall that multiplying matrices  $m \times n$  by  $n \times p$  yields a  $m \times p$  matrix.

We will want to multiply the position  $x$  by 1 to get the corresponding measurement of the position. We do not need to use velocity to find the corresponding measurement so we multiply  $\dot{x}$  by 0.

$$\mathbf{y} = \mathbf{z} - [1 \ 0] \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

And so, for our Kalman filter we set

$$\mathbf{H} = [1 \ 0]$$

In [17]: `kf.H = np.array([[1., 0.]])`

We have designed the majority of our Kalman filter. All that is left is to model the noise in the sensors.

### Step 6: Design the Measurement Noise Matrix

The measurement noise matrix models the noise in our sensors as a covariance matrix. In practice this can be difficult. A complicated system may have many sensors, the correlation between them might not be clear, and usually their noise is not a pure Gaussian. For example, a sensor might be biased to read high if the temperature is high, and so the noise is not distributed equally on both sides of the mean. We will learn to deal with these problems later.

The Kalman filter equations uses a covariance matrix  $\mathbf{R}$  for the measurement noise. The matrix will have dimension  $m \times m$ , where  $m$  is the number of sensors. It is a covariance matrix to account for correlations between the sensors. We have only 1 sensor so  $\mathbf{R}$  is:

$$\mathbf{R} = [\sigma_z^2]$$

If  $\sigma_z^2$  is 5 meters squared we'd have  $\mathbf{R} = [5]$ .

If we had two position sensors, the first with a variance of 5  $\text{m}^2$ , the second with a variance of 3  $\text{m}^2$ , we would write

$$\mathbf{R} = \begin{bmatrix} 5 & 0 \\ 0 & 3 \end{bmatrix}$$

We put the variances on the diagonal because this is a covariance matrix, where the variances lie on the diagonal, and the covariances, if any, lie in the off-diagonal elements. Here we assume there is no correlation in the noise between the two sensors, so the covariances are 0.

For our problem we only have one sensor, so we can implement this as

```
In [18]: kf.R = np.array([[5.]])
```

We have a few shortcuts available. `KalmanFilter.R` is initialized to the identity matrix. For this simple problem we could write:

```
kf.R *= 5.
```

Finally, since this is a 1x1 matrix you can just use a scalar.

```
kf.R = 5.
```

#### 6.4.4 Implementing the Kalman Filter

I've given you all of the code for the filter, but now let's collect it in one place. First we construct an `KalmanFilter` object. We have to specify the number of variable in the state with the `dim_x` parameter, and the number of sensors/measurements with `dim_z`. We have two random variables in the state and one measurement, so we write

```
In [19]: from filterpy.kalman import KalmanFilter
dog_filter = KalmanFilter(dim_x=2, dim_z=1)
```

Now we initialize the filter's matrices and vectors.

```
In [20]: from filterpy.common import Q_discrete_white_noise
```

```
dt = .1
dog_filter.x = np.array([0., 0.])      # state (location and velocity)
dog_filter.F = np.array([[1., dt],       # state transition function
                       [0., 1.]])
dog_filter.H = np.array([[1., 0.]])      # Measurement function
dog_filter.R = 5                         # measurement noise
dog_filter.Q = Q_discrete_white_noise(   # process noise
    2, dt=dt, var=0.1)
dog_filter.P = np.diag([500., 49.])      # covariance matrix
```

Let's look at this line by line.

1: We assign the initial value for our state. Here we initialize both the position and velocity to zero.

2: We set  $\mathbf{F} = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix}$ , as in design step 2 above.

3: We set  $\mathbf{H} = [1 \ 0]$ , as in design step 4 above.

4: We set  $\mathbf{R} = [5]$ .

5 We use the `Q_discrete_white_noise()` method to set  $\mathbf{Q}$ 's variance to 0.1.

6: We set  $\mathbf{P} = \begin{bmatrix} 500 & 0 \\ 0 & 49 \end{bmatrix}$

I've put this in a function to allow you to specify different initial values for  $\mathbf{R}$ ,  $\mathbf{P}$ , and  $\mathbf{Q}$  and put it in a helper function. We will be creating and running many of these filters, and this saves us a lot of headaches.

```
In [21]: def pos_vel_filter(x, P, R, Q=0., dt=1.0):
    """ Returns a KalmanFilter which implements a
    constant velocity model for a state [x dx].T
    """
    kf = KalmanFilter(dim_x=2, dim_z=1)
    kf.x = np.array([x[0], x[1]]) # location and velocity
    kf.F = np.array([[1, dt],
                    [0, 1]])      # state transition matrix
    kf.H = np.array([[1, 0]])      # Measurement function
    kf.R *= R                      # measurement uncertainty
    if np.isscalar(P):
        kf.P *= P                  # covariance matrix
    else:
        kf.P[:] = P
    if np.isscalar(Q):
        kf.Q = Q_discrete_white_noise(dim=2, dt=dt, var=Q)
    else:
        kf.Q = Q
    return kf
```

All that is left is to write the code to run the Kalman filter.

This is the complete code for the filter, and most of it is boilerplate. I've made it flexible enough to support several uses in this chapter, so it is a bit verbose. Lets work through it line by line.

The first lines checks to see if you provided it with measurement data in `data`. If not, it creates the data using the `compute_dog_data` function we wrote earlier.

The next lines uses our helper function to create a Kalman filter.

```
# create the Kalman filter
kf = pos_vel_filter(P=P, R=R, Q=Q)
```

The next line initialize empty lists to store the state and covariance so we can plot them later.

```
xs, cov = [], []
```

Finally we get to the filter. All we need to do is perform the update and predict steps of the Kalman filter for each measurement. The `KalmanFilter` class provides the two methods `update()` and `predict()` for this purpose. `update()` performs the measurement update step of the Kalman filter, and so it takes a variable containing the sensor measurement.

Absent the work of storing the results, the loop reads:

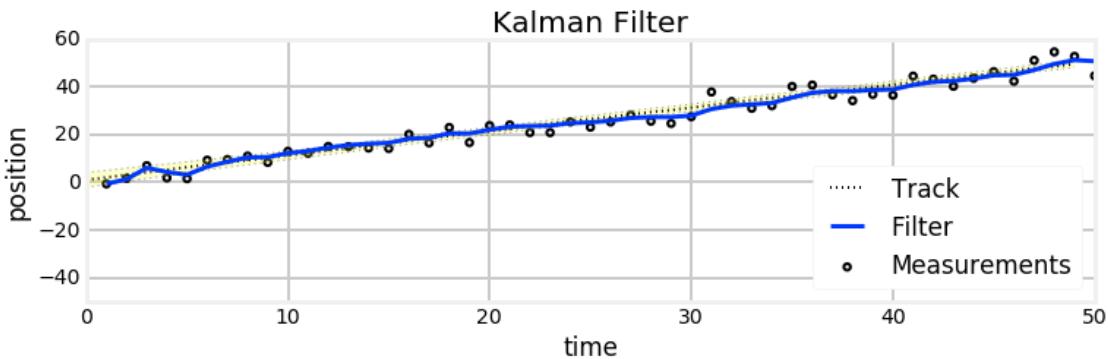
```
for z in zs:
    kf.predict()
    kf.update(z)
```

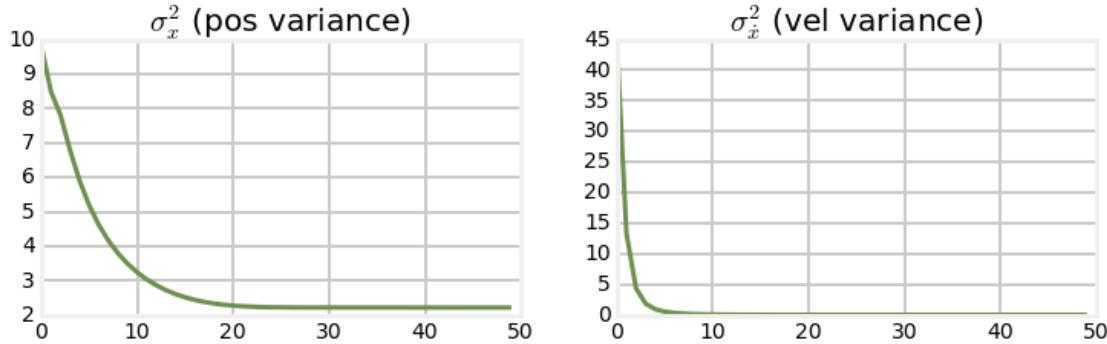
It really cannot get much simpler than that. As we tackle more complicated problems this code will remain largely the same; all of the work goes into setting up the `KalmanFilter` matrices; executing the filter is trivial.

The rest of the code optionally plots the results and then returns the saved states and covariances.

Let's run it. I start by filtering 50 measurements with a noise variance of 10 and a process variance of 0.01.

```
In [23]: P = np.diag([500., 49.])
Ms, Ps = run(count=50, R=10, Q=0.01, P=P)
```





There is still a lot to learn, but we have implemented our first, full Kalman filter using the same theory and equations as published by Rudolf Kalman! Code very much like this runs inside of your GPS, inside every airliner, inside of robots, and so on.

The first plot plots the output of the Kalman filter against the measurements and the actual position of our dog (the ‘Track’ line). After the initial settling in period the filter should track the dog’s position very closely. The yellow shaded portion between the black dotted lines shows 1 standard deviations of the filter’s variance, which I explain in the next paragraph.

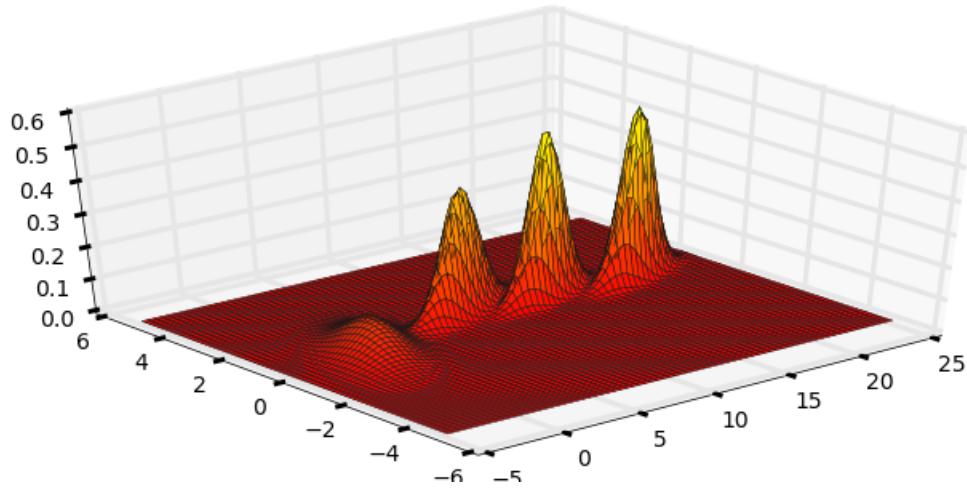
The next two plots show the variance of  $x$  and of  $\dot{x}$ . If you look at the code, you will see that I have plotted the diagonals of  $\mathbf{P}$  over time. Recall that the diagonal of a covariance matrix contains the variance of each state variable. So  $\mathbf{P}[0, 0]$  is the variance of  $x$ , and  $\mathbf{P}[1, 1]$  is the variance of  $\dot{x}$ . You can see that we quickly converge to small variances for both.

The covariance matrix  $\mathbf{P}$  tells us the theoretical performance of the filter assuming everything we tell it is true. Recall from the Gaussian chapter that the standard deviation is the square root of the variance, and that approximately 68% of a Gaussian distribution occurs within one standard deviation. Therefore, if at least 68% of the filter output is within one standard deviation we can be sure that the filter is performing well. In the top chart I have displayed the one standard deviation as the yellow shaded area between the two dotted lines. To my eye it looks like perhaps the filter is slightly exceeding that bounds, so the filter probably needs some tuning.

In the univariate chapter we filtered very noisy signals with much simpler code than the code above. However, realize that right now we are working with a very simple example - an object moving through 1-D space and one sensor. That is about the limit of what we can compute with the code in the last chapter. In contrast, we can implement very complicated, multidimensional filter with this code merely by altering our assignments to the filter’s variables. Perhaps we want to track 100 dimensions in financial models. Or we have an aircraft with a GPS, INS, TACAN, radar altimeter, baro altimeter, and airspeed indicator, and we want to integrate all those sensors into a model that predicts position, velocity, and accelerations in 3D (which requires 9 state variables). We can do that with the code in this chapter.

I want you to get a better feel for how the Gaussians change over time, so here is a 3D plot showing the Gaussians every 7th epoch (time step). Every 7th separates them enough so can see each one independently.

```
In [24]: from book_format import set_figsize, figsize
        from nonlinear_plots import plot_gaussians
        P = np.diag([3., 1.])
        np.random.seed(3)
        Ms, Ps = run(count=25, R=10, Q=0.01, P=P, do_plot=False)
        with figsize(y=5):
            plot_gaussians(Ms[::7], Ps[::7], (-5,25), (-5, 5), 75)
```



## 6.5 The Kalman Filter Equations

We are now ready to learn how `predict()` and `update()` perform their computations.

A word about my notation. I'm a programmer, and I am used to code reading  $x = x + 1$ . That is not an equation as the sides are not equal, but an assignment. If we wanted to write this in correct mathematical notation we'd write

$$x_k = x_{k-1} + 1$$

The Kalman filter equations are littered with subscripts and superscripts to keep the equations mathematically consistent. This makes them extremely hard to read. In most of the book I opt for subscriptless statements. As a programmer you should understand that I am showing you assignments which implement an algorithm that is to be executed step by step. I'll elaborate on this once we have a concrete example.

### 6.5.1 Prediction Equations

The Kalman filter uses these equations to compute the prior - the predicted next state of the system. They compute the mean ( $\bar{\mathbf{x}}$ ) and covariance ( $\bar{\mathbf{P}}$ ) of the system.

$$\begin{aligned}\bar{\mathbf{x}} &= \mathbf{F}\mathbf{x} + \mathbf{B}\mathbf{u} \\ \bar{\mathbf{P}} &= \mathbf{F}\mathbf{P}\mathbf{F}^T + \mathbf{Q}\end{aligned}$$

#### Mean

$$\bar{\mathbf{x}} = \mathbf{F}\mathbf{x} + \mathbf{B}\mathbf{u}$$

We are starting out easy - you were already exposed to the first equation while designing the state transition function  $\mathbf{F}$  and control function  $\mathbf{B}$ .

As a reminder, the linear equation  $\mathbf{Ax} = \mathbf{B}$  represents a system of equations, where  $\mathbf{A}$  holds the coefficients of the state  $\mathbf{x}$ . Performing the multiplication  $\mathbf{Ax}$  computes the right hand side values for that set of equations.

If  $\mathbf{F}$  contains the state transition for a given time step, then the product  $\mathbf{Fx}$  computes the state after that transition. Easy! Likewise,  $\mathbf{B}$  is the control function,  $\mathbf{u}$  is the control input, so  $\mathbf{Bu}$  computes the contribution

of the controls to the state after the transition. Thus, the prior of the state ( $\bar{\mathbf{x}}$ ) is computed as the sum of  $\mathbf{F}\mathbf{x}$  and  $\mathbf{B}\mathbf{u}$ .

### Covariance

$$\bar{\mathbf{P}} = \mathbf{F}\mathbf{P}\mathbf{F}^T + \mathbf{Q}$$

This equation is not as easy to understand so we will spend more time on it.

In the univariate chapter when we added Gaussians in the predict step we did it this way:

$$\mu = \mu + \mu_{move}\sigma^2 = \sigma^2 + \sigma_{move}^2$$

We added the variance of the movement to the variance of our estimate to reflect the loss of knowledge. We need to do the same thing here, except it isn't quite that easy with multivariate Gaussians. We can't simply write  $\bar{\mathbf{P}} = \mathbf{P} + \mathbf{Q}$ .

In a multivariate Gaussians the state variables are correlated. What does this imply? Our knowledge of the velocity is imperfect, but we are adding it to the position with

$$\bar{x} = \dot{x}\Delta t + x$$

Since we do not have perfect knowledge of the value of  $\dot{x}$  the sum  $\bar{x} = \dot{x}\Delta t + x$  gains uncertainty. Because the positions and velocities are correlated we cannot simply add the covariance matrices. The correct equation is

$$\bar{\mathbf{P}} = \mathbf{F}\mathbf{P}\mathbf{F}^T$$

Expressions in the form  $\mathbf{ABA}^T$  are common in linear algebra. You can think of it as projecting the middle term by the outer term. We will be using this many times in the rest of the book.

I admit this may be a ‘magical’ equation to you.

When we initialize  $\mathbf{P}$  with

$$\mathbf{P} = \begin{bmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_v^2 \end{bmatrix}$$

the value for  $\mathbf{F}\mathbf{P}\mathbf{F}^T$  is:

$$\mathbf{F}\mathbf{P}\mathbf{F}^T = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_v^2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \Delta t & 1 \end{bmatrix} = \begin{bmatrix} \sigma_x^2 & \sigma_v^2 \Delta t \\ 0 & \sigma_v^2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \Delta t & 1 \end{bmatrix} = \begin{bmatrix} \sigma_x^2 + \sigma_v^2 \Delta t^2 & \sigma_v^2 \Delta t \\ \sigma_v^2 \Delta t & \sigma_v^2 \end{bmatrix}$$

The initial value for  $\mathbf{P}$  had no covariance between the position and velocity. Position is computed as  $\dot{x}\Delta t + x$ , so there is a correlation between the position and velocity. The multiplication  $\mathbf{F}\mathbf{P}\mathbf{F}^T$  computes a covariance of  $\sigma_v^2 \Delta t$ . The exact value is not important; you just need to recognize that  $\mathbf{F}\mathbf{P}\mathbf{F}^T$  uses the process model to automatically compute the covariance between the position and velocity!

If you have some experience with linear algebra and statistics, this may help. The covariance due to the prediction can be modeled as the expected value of the error in the prediction step, given by this equation.

$$\begin{aligned} \bar{\mathbf{P}} &= \mathbb{E}[(\mathbf{F}\mathbf{x})(\mathbf{F}\mathbf{x})^T] \\ &= \mathbb{E}[\mathbf{F}\mathbf{x}\mathbf{x}^T\mathbf{F}^T] \\ &= \mathbf{F} \mathbb{E}[\mathbf{x}\mathbf{x}^T] \mathbf{F}^T \end{aligned}$$

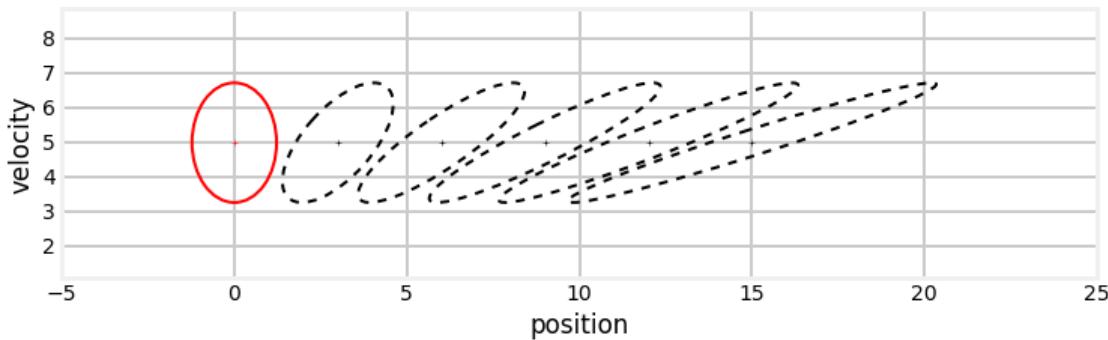
Of course,  $\mathbb{E}[\mathbf{x}\mathbf{x}^T]$  is just  $\mathbf{P}$ , giving us

$$\bar{\mathbf{P}} = \mathbf{F}\mathbf{P}\mathbf{F}^T$$

For now we will look at its effect. Here I use  $\mathbf{F}$  from our filter and project the state forward 6/10ths of a second. I do this five times so you can see how  $\bar{\mathbf{P}}$  continues to change.

```
In [25]: import matplotlib.pyplot as plt
import book_plots
dt = 0.6
x = np.array([0., 5.])
F = np.array([[1., dt], [0, 1.]])
P = np.array([[1.5, 0], [0, 3.]])
plot_covariance_ellipse(x, P, edgecolor='r')

for _ in range(5):
    x = np.dot(F, x)
    P = np.dot(F, P).dot(F.T)
    plot_covariance_ellipse(x, P, edgecolor='k', ls='dashed')
book_plots.set_labels(x='position', y='velocity')
```



You can see that with a velocity of 5 the position correctly moves 3 units in each 6/10ths of a second step. At each step the width of the ellipse is larger, indicating that we have lost information about the position due to adding  $\dot{x}\Delta t$  to  $x$  at each step. The height has not changed - our system model says the velocity does not change, so the belief we have about the velocity cannot change. As time continues you can see that the ellipse becomes more and more tilted. Recall that a tilt indicates correlation.  $\mathbf{F}$  linearly correlates  $x$  with  $\dot{x}$  with the expression  $\bar{x} = \dot{x}\Delta t + x$ . The  $\mathbf{F}\mathbf{P}\mathbf{F}^T$  computation correctly incorporates this correlation into the covariance matrix.

Here is an animation of this equation that allows you to change the design of  $\mathbf{F}$  to see how it affects shape of  $\mathbf{P}$ . The `F00` slider affects the value of  $\mathbf{F}[0, 0]$ . `covar` sets the intial covariance between the position and velocity( $\sigma_x\sigma_{\dot{x}}$ ). I recommend answering these questions at a minimum

- what if  $x$  is not correlated to  $\dot{x}$ ? (set  $\mathbf{F}[0, 0]$  to 0, the rest at defaults)
- what if  $x = 2\dot{x}\Delta t + x_0$ ? (set  $\mathbf{F}[0, 0]$  to 2, the rest at defaults)
- what if  $x = \dot{x}\Delta t + 2x_0$ ? (set  $\mathbf{F}[0, 0]$  to 2, the rest at defaults)
- what if  $x = \dot{x}\Delta t$ ? (set  $\mathbf{F}[0, 0]$  to 0, the rest at defaults)

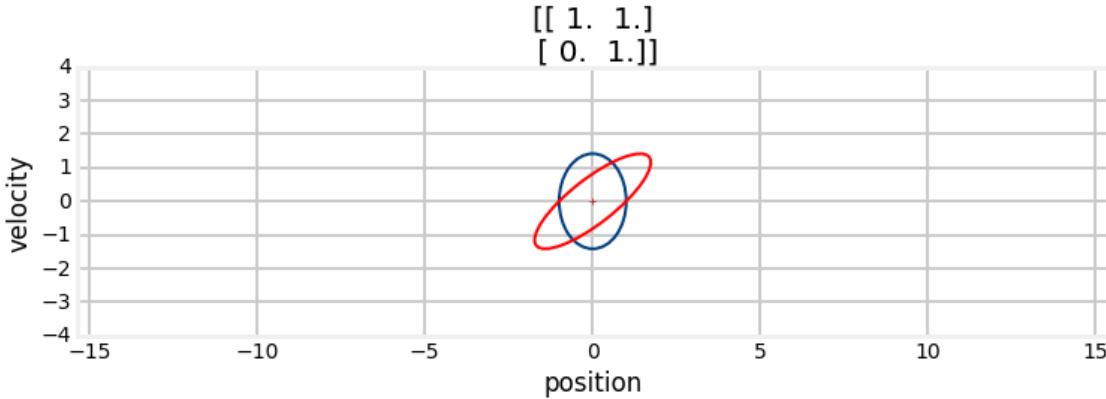
```
In [26]: from IPython.html.widgets import interact, interactive, fixed
from IPython.html.widgets import IntSlider, FloatSlider
```

```

def plot_FPFT(F00, F01, F10, F11, covar):
    dt = 1.
    x = np.array((0, 0.))
    P = np.array(((1, covar), (covar, 2)))
    F = np.array(((F00, F01), (F10, F11)))
    plot_covariance_ellipse(x, P)
    plot_covariance_ellipse(x, np.dot(F, P).dot(F.T), ec='r')
    plt.gca().set_aspect('equal')
    plt.xlim(-4, 4)
    plt.ylim(-4, 4)
    plt.title(str(F))
    plt.xlabel('position')
    plt.ylabel('velocity')

interact(plot_FPFT,
         F00=IntSlider(value=1, min=0, max=2.),
         F01=FloatSlider(value=1, min=0., max=2.,
                         description='F01(dt)'),
         F10=FloatSlider(value=0, min=0., max=2.),
         F11=FloatSlider(value=1, min=0., max=2.),
         covar=FloatSlider(value=0, min=0, max=1.));

```



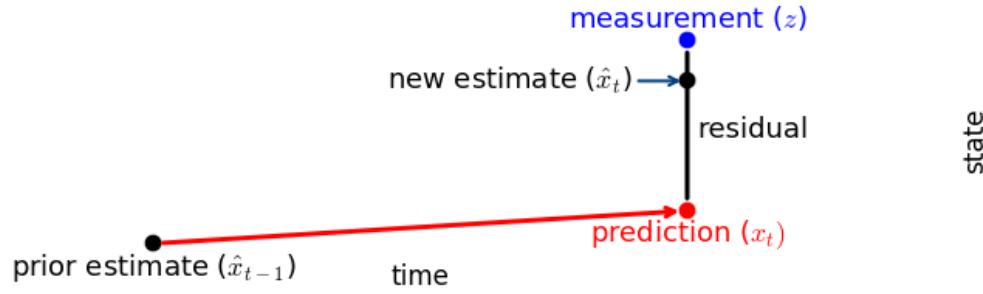
(If you are reading this in a static form: instructions to run this online are here: <http://git.io/vmv6e>)

### 6.5.2 Update Equations

The update equations look messier than the predict equations, but that is mostly due to the Kalman filter computing the update in measurement space. This is because measurements are not invertible. For example, later on we will be using sensors that only give you the range to a target. It is impossible to convert a range into a position - an infinite number of positions (in a circle) will yield the same range. On the other hand, we can always compute the range (measurement) given a position (state).

Before I continue, recall that we are trying to do something very simple: choose a new estimate chosen somewhere between a measurement and a prediction, as in this chart:

```
In [27]: with figsize(y=2.5):
    show_residual_chart()
```



The equations will be complicated because the state has multiple dimensions, but this operations is what we are doing. Don't let the equations distract you from the simplicity of this idea.

### System Uncertainty

$$\mathbf{S} = \mathbf{H}\bar{\mathbf{P}}\mathbf{H}^T + \mathbf{R}$$

To work in measurement space the Kalman filter has to project the covariance matrix into measurement space. The math for this is  $\mathbf{H}\bar{\mathbf{P}}\mathbf{H}^T$ , where  $\bar{\mathbf{P}}$  is the prior covariance and  $\mathbf{H}$  is the measurement function.

You should recognize this  $\mathbf{ABA}^T$  form - the prediction step used  $\mathbf{FPF}^T$  to update  $\mathbf{P}$  with the state transition function. Here, we use the same form to update it with the measurement function. In a real sense the linear algebra is changing the coordinate system for us.

Once the covariance is in measurement space we need to account for the sensor noise. This is very easy - we just add matrices. The result is variously called the system uncertainty or innovation covariance.

Compare the equations for the system uncertainty and the covariance

$$\begin{aligned}\mathbf{S} &= \mathbf{H}\bar{\mathbf{P}}\mathbf{H}^T + \mathbf{R} \\ \bar{\mathbf{P}} &= \mathbf{FPF}^T + \mathbf{Q}\end{aligned}$$

In each equation  $\mathbf{P}$  is put into a different space with either the function  $\mathbf{H}$  or  $\mathbf{F}$ . Once that is done we add the noise matrix associated with that space.

### Kalman Gain

$$\mathbf{K} = \bar{\mathbf{P}}\mathbf{H}^T \mathbf{S}^{-1}$$

Look back at the residual diagram. Once we have a prediction and a measurement we need to select an estimate somewhere between the two. If we have more certainty about the measurement the estimate will be closer to it. If instead we have more certainty about the prediction then the estimate will be closer to it.

In the univariate chapter we scaled the mean using this equation

$$\mu = \frac{\sigma_{\text{prior}}^2 \mu_2 + \sigma_z^2 \mu_{\text{prior}}}{\sigma_{\text{prior}}^2 + \sigma_z^2}$$

which we simplified to

$$\mu = (1 - K)\mu_{\text{prior}} + K\mu_z$$

$K$  is the Kalman gain, and it is a real number between 0 and 1. Ensure you understand how it selects a mean somewhere between the prediction and measurement. The Kalman gain is a percentage or ratio - if  $K$  is .9 it takes 90% of the measurement and 10% of the prediction.

For the multivariate Kalman filter  $\mathbf{K}$  is a vector, not a scalar. Here is the equation again:  $\mathbf{K} = \bar{\mathbf{P}}\mathbf{H}^T\mathbf{S}^{-1}$ . Is this a ratio? We can think of the inverse of a matrix as linear algebra's way of finding the reciprocal. Division is not defined for matrices, but it is useful to think of it in this way. So we can read the equation for  $\mathbf{K}$  as meaning

$$\mathbf{K} \approx \frac{\bar{\mathbf{P}}\mathbf{H}^T}{\mathbf{S}}$$

$$\mathbf{K} \approx \frac{\text{uncertainty}_{\text{prediction}}}{\text{uncertainty}_{\text{measurement}}} \mathbf{H}^T$$

The Kalman gain equation computes a ratio based on how much we trust the prediction vs the measurement. We did the same thing in every prior chapter. The equation is complicated because we are doing this in multiple dimensions via matrices, but the concept is simple. The  $\mathbf{H}^T$  term is less clear, I'll explain it soon.

### Residual

$$\mathbf{y} = \mathbf{z} - \mathbf{H}\mathbf{x}$$

This is an easy one as we've covered this equation while designing the measurement function  $\mathbf{H}$ . Recall that the measurement function converts a state into a measurement. So  $\mathbf{H}\mathbf{x}$  converts  $\mathbf{x}$  into an equivalent measurement. Once that is done, we can subtract it from the measurement  $\mathbf{z}$  to get the residual - the difference between the measurement and prediction.

### State Update

$$\mathbf{x} = \bar{\mathbf{x}} + \mathbf{K}\mathbf{y}$$

We select our new state to be along the residual, scaled by the Kalman gain. The scaling is performed by  $\mathbf{K}\mathbf{y}$ , which both scales the residual and converts it back into state space with the  $\mathbf{H}^T$  term which is in  $\mathbf{K}$ . This is added to the prior, yielding the equation:  $\mathbf{x} = \bar{\mathbf{x}} + \mathbf{K}\mathbf{y}$ . Let me write out  $\mathbf{K}$  so we can see the entire computation:

$$\begin{aligned} \mathbf{x} &= \bar{\mathbf{x}} + \mathbf{K}\mathbf{y} \\ &= \bar{\mathbf{x}} + \bar{\mathbf{P}}\mathbf{H}^T\mathbf{S}^{-1}\mathbf{y} \\ &\approx \bar{\mathbf{x}} + \frac{\text{uncertainty}_{\text{prediction}}}{\text{uncertainty}_{\text{measurement}}} \mathbf{H}^T\mathbf{y} \end{aligned}$$

### Covariance Update

$$\mathbf{P} = (\mathbf{I} - \mathbf{K}\mathbf{H})\bar{\mathbf{P}}$$

$\mathbf{I}$  is the identity matrix, and is the way we represent 1 in multiple dimensions.  $\mathbf{H}$  is our measurement function, and is a constant. We can think of the equation as  $\mathbf{P} = (1 - c\mathbf{K})\mathbf{P}$ .  $\mathbf{K}$  is our ratio of how much prediction vs measurement we use. If  $\mathbf{K}$  is large then  $(1 - c\mathbf{K})$  is small, and  $\mathbf{P}$  will be made smaller than it was. If  $\mathbf{K}$  is small, then  $(1 - c\mathbf{K})$  is large, and  $\mathbf{P}$  will be relatively larger. This means that we adjust the size of our uncertainty by some factor of the Kalman gain.

This equation can be numerically unstable and I don't use it in FilterPy. Later I'll share more complicated but numerically stable forms of this equation.

### 6.5.3 An Example not using FilterPy

FilterPy hides the details of the implementation from us. Normally you will appreciate this, but let's implement the last filter without FilterPy. To do so we need to define our matrices as variables, and then implement the Kalman filter equations explicitly.

Here we initialize our matrices:

```
In [28]: dt = 1.
R_var = 10
Q_var = 0.01
x = np.array([[10.0, 4.5]]).T
P = np.diag([500, 49])
F = np.array([[1, dt],
              [0, 1]])
H = np.array([[1., 0.]])
R = np.array([[R_var]])
Q = Q_discrete_white_noise(dim=2, dt=dt, var=Q_var)
```

```
In [29]: from numpy import dot
from scipy.linalg import inv

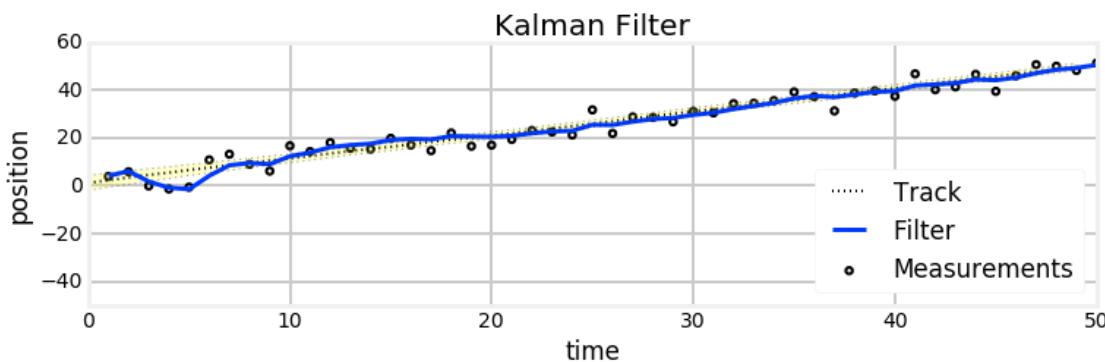
count = 50
data = compute_dog_data(R_var, Q_var, count)
xs, cov = [], []
for z in data[:, 1]:
    # predict
    x = dot(F, x)
    P = dot(F, P).dot(F.T) + Q

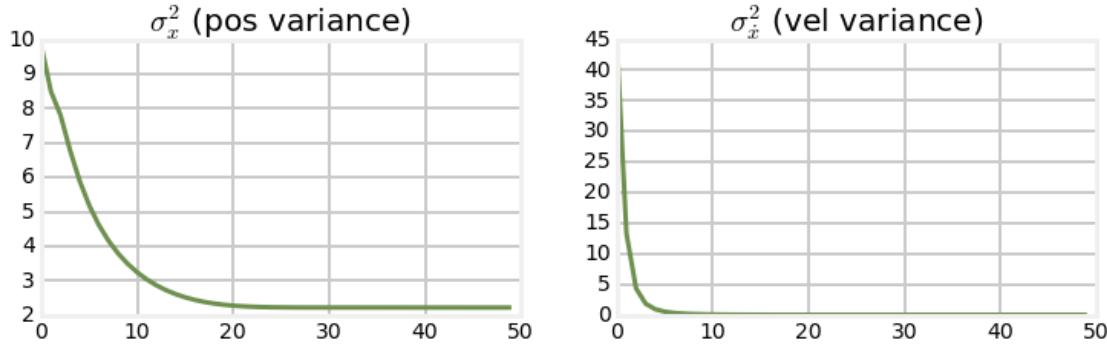
    #update
    S = dot(H, P).dot(H.T) + R
    K = dot(P, H.T).dot(inv(S))
    y = z - dot(H, x)
    x += dot(K, y)
    P = P - dot(K, H).dot(P)

    xs.append(x)
    cov.append(P)

xs, cov = np.array(xs), np.array(cov)

plot_track(xs[:, 0], data[:, 0], data[:, 1], cov, dt=dt)
```





The results are identical to the FilterPy version. Which you prefer is up to you. If you do not like object oriented programming FilterPy might be off-putting. I prefer not polluting my namespace with variables such as `x`, `P`, and so on; `dog_filter.x` is, to me, more readable.

More importantly, this example requires you to remember and program the equations for the Kalman filter. Sooner or later you will make a mistake. FilterPy's version ensures that your code will be correct. On the other hand, if you make a mistake in your definitions, such as making `H` a column vector instead of a row vector, FilterPy's error message will be harder to debug than this explicit code.

FilterPy's `KalmanFilter` class provides additional functionality such as smoothing, batch processing, faded memory filtering, computation of the maximum likelihood function, and more. You get all of this functionality without having to explicitly program it.

Finally, there are multiple forms of some of the Kalman filter equations, some which we will present later. It is cumbersome to change the equations using FilterPy, whereas it is trivial if you explicitly code the equations yourself.

In the rest of the book I will continue using FilterPy - I want you to focus on issues of design, not implementation of linear algebra. If the linear algebra interests you I encourage you to rewrite the examples without FilterPy.

#### 6.5.4 Summary

We have learned the Kalman filter equations. Here they are all together for your review. There was a lot to learn, but I hope that as you went through each you recognized its kinship with the equations in the univariate filter. In the [Kalman Math](#) chapter I will show you that if we set the dimension of `x` to one that these equations revert back to the equations for the univariate filter. This is not “like” the univariate filter - it is a multiple dimension implementation of it.

Predict Step

$$\begin{aligned}\bar{\mathbf{x}} &= \mathbf{Fx} + \mathbf{Bu} \\ \bar{\mathbf{P}} &= \mathbf{FPF}^T + \mathbf{Q}\end{aligned}$$

Update Step

$$\begin{aligned}\mathbf{S} &= \mathbf{H}\bar{\mathbf{P}}\mathbf{H}^T + \mathbf{R} \\ \mathbf{K} &= \bar{\mathbf{P}}\mathbf{H}^T \mathbf{S}^{-1} \\ \mathbf{y} &= \mathbf{z} - \mathbf{H}\bar{\mathbf{x}} \\ \mathbf{x} &= \bar{\mathbf{x}} + \mathbf{Ky} \\ \mathbf{P} &= (\mathbf{I} - \mathbf{Kh})\bar{\mathbf{P}}\end{aligned}$$

I want to share a form of the equations that you will see in the literature. There are many different notation systems used, but this gives you an idea of what to expect.

$$\begin{aligned}\hat{\mathbf{x}}_{k|k-1} &= \mathbf{F}_k \hat{\mathbf{x}}_{k-1|k-1} + \mathbf{B}_k \mathbf{u}_k \\ \mathbf{P}_{k|k-1} &= \mathbf{F}_k \mathbf{P}_{k-1|k-1} \mathbf{F}_k^T + \mathbf{Q}_k \\ \tilde{\mathbf{y}}_k &= \mathbf{z}_k - \mathbf{H}_k \hat{\mathbf{x}}_{k|k-1} \\ \mathbf{S}_k &= \mathbf{H}_k \mathbf{P}_{k|k-1} \mathbf{H}_k^T + \mathbf{R}_k \\ \mathbf{K}_k &= \mathbf{P}_{k|k-1} \mathbf{H}_k^T \mathbf{S}_k^{-1} \\ \hat{\mathbf{x}}_{k|k} &= \hat{\mathbf{x}}_{k|k-1} + \mathbf{K}_k \tilde{\mathbf{y}}_k \\ \mathbf{P}_{k|k} &= (I - \mathbf{K}_k \mathbf{H}_k) \mathbf{P}_{k|k-1}\end{aligned}$$

This notation uses the Bayesian  $a | b$  notation, which means  $a$  given the evidence of  $b$ . The hat means estimate. Thus  $\hat{\mathbf{x}}_{k|k}$  means the estimate of the state  $\mathbf{x}$  at step  $k$  (the first  $k$ ) given the evidence from step  $k$  (the second  $k$ ). The posterior, in other words.  $\hat{\mathbf{x}}_{k|k-1}$  means the estimate for the state  $\mathbf{x}$  at step  $k$  given the estimate from step  $k-1$ . The prior, in other words.

This notation, copied from [Wikipedia](#) [1], allows a mathematician to express himself exactly. In formal publications presenting new results this precision is necessary. As a programmer I find it fairly unreadable. I am used to thinking about variables changing state as a program runs, and do not use a different variable name for each new computation. There is no agreed upon format in the literature, so each author makes different choices. I find it challenging to switch quickly between books and papers, and so have adopted my admittedly less precise notation. Mathematicians may write scathing emails to me, but I hope programmers and students will rejoice at my simplified notation.

The **Symbology** Appendix lists the notation used by various authors. This brings up another difficulty. Different authors use different variable names.  $\mathbf{x}$  is fairly universal, but after that it is anybody's guess. For example, it is common to use  $\mathbf{A}$  for what I call  $\mathbf{F}$ . You must read carefully, and hope that the author defines their variables (they often do not).

If you are a programmer trying to understand a paper's math equations, I suggest starting by removing all of the superscripts, subscripts, and diacriticals, replacing them with a single letter. If you work with equations like this every day this is superfluous advice, but when I read I am usually trying to understand the flow of computation. To me it is far more understandable to remember that  $P$  in this step represents the updated value of  $P$  computed in the last step, as opposed to trying to remember what  $P_{k-1}(+)$  denotes, and what its relation to  $P_k(-)$  is, if any, and how any of that relates to the completely different notation used in the paper I read 5 minutes ago.

## 6.6 Exercise: Show Effect of Hidden Variables

In our filter velocity is a hidden variable. How would a filter perform if we did not use velocity in the state?

Write a Kalman filter that uses the state  $\mathbf{x} = [x]$  and compare it against a filter that uses  $\mathbf{x} = [x \ \dot{x}]^\top$ .

In [30]: # your code here

### 6.6.1 Solution

We've already implemented a Kalman filter for position and velocity, so I will provide the code without much comment, and then plot the result.

```
In [31]: from math import sqrt
        from numpy.random import randn

        def univariate_filter(x0, P, R, Q):
            f = KalmanFilter(dim_x=1, dim_z=1, dim_u=1)
            f.x = np.array([[x0]])
            f.P *= P
            f.H = np.array([[1.]])
            f.F = np.array([[1.]])
            f.B = np.array([[1.]])
            f.Q *= Q
            f.R *= R
            return f

        def plot_1d_2d(xs, xs1d, xs2d):
            plt.plot(xs1d, label='1D Filter')
            plt.scatter(range(len(xs2d)), xs2d, c='r', label='2D Filter')
            plt.plot(xs, ls='--', color='k', lw=1, label='track')
            plt.title('State')
            plt.legend(loc=4)
            plt.show()

        def compare_1D_2D(x0, P, R, Q, vel, u=0):
            # storage for filter output
            xs, xs1, xs2 = [], [], []

            # 1d KalmanFilter
            f1D = univariate_filter(x0, P, R, Q)

            # 2D Kalman filter
            f2D = pos_vel_filter(x=(x0, vel), P=P, R=R, Q=0)

            pos = 0 # true position
            for i in range(100):
                pos += vel
                xs.append(pos)

                # control input u - discussed below
                f1D.predict(u=u)
                f2D.predict()
```

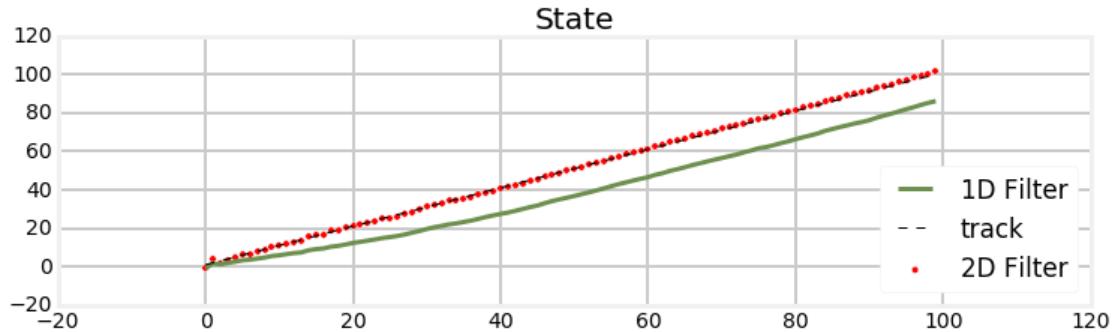
```

z = pos + randn()*sqrt(R) # measurement
f1D.update(z)
f2D.update(z)

xs1.append(f1D.x[0])
xs2.append(f2D.x[0])
plot_1d_2d(xs, xs1, xs2)

compare_1D_2D(x0=0, P=50., R=5., Q=.02, vel=1.)

```



### 6.6.2 Discussion

The filter that incorporates velocity into the state produces much better estimates than the filter that only tracks position. The univariate filter has no way to estimate the velocity or change in position, so it lags the tracked object.

In the univariate Kalman filter chapter we had a control input  $u$  to the predict equation:

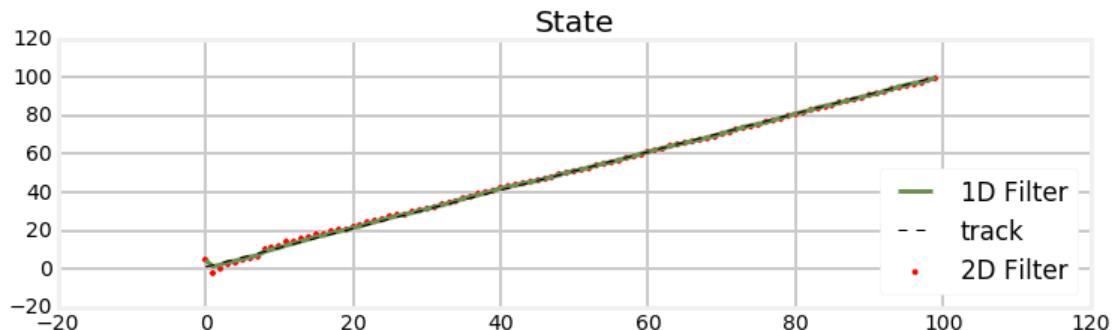
```

def predict(self, u=0.0):
    self.x += u
    self.P += self.Q

```

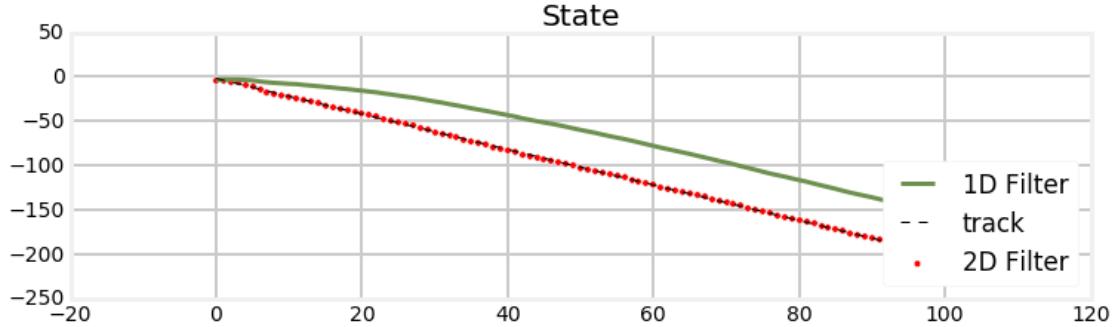
It may strike you as ‘cheating’ that I did not use it here. Let’s try that:

In [32]: `compare_1D_2D(x0=0, P=50., R=5., Q=.02, vel=1., u=1.)`



Here the performance of the two filters are very close, and perhaps the univariate filter is tracking more closely. But let's see what happens when the actual velocity  $\text{vel}$  is different from the control input  $u$ :

In [33]: `compare_1D_2D(x0=0, P=50., R=5., Q=.02, vel=-2., u=1.)`



If we are tracking a robot which we are also controlling the univariate filter can do a very good job because the control input allows the filter to make an accurate prediction. But if we are tracking passively the control input is not much help unless we can make an accurate apriori guess as to the velocity. This is rarely possible.

## 6.7 How Velocity is Calculated

I haven't really explained how the filter computes the velocity, or any hidden variable. If we plug in the values we calculated for each of the filter's matrices we can see what happens.

First we need to compute the system uncertainty.

$$\begin{aligned} \mathbf{S} &= \mathbf{H}\bar{\mathbf{P}}\mathbf{H}^T + \mathbf{R} \\ &= [1 \quad 0] \begin{bmatrix} \sigma_x^2 & \sigma_{xv} \\ \sigma_{xv} & \sigma_v^2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} + [\sigma_z^2] \\ &= [\sigma_x^2 \quad \sigma_{xv}] \begin{bmatrix} 1 \\ 0 \end{bmatrix} + [\sigma_z^2] \\ &= [\sigma_x^2 + \sigma_z^2] \end{aligned}$$

Now that we have  $\mathbf{S}$  we can find the value for the Kalman gain:

$$\begin{aligned} \mathbf{K} &= \bar{\mathbf{P}}\mathbf{H}^T \mathbf{S}^{-1} \\ &= \begin{bmatrix} \sigma_x^2 & \sigma_{xv} \\ \sigma_{xv} & \sigma_v^2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \left[ \frac{1}{\sigma_x^2 + \sigma_z^2} \right] \\ &= \begin{bmatrix} \sigma_x^2 \\ \sigma_{xv} \end{bmatrix} \left[ \frac{1}{\sigma_x^2 + \sigma_z^2} \right] \\ &= \begin{bmatrix} \sigma_x^2 / (\sigma_x^2 + \sigma_z^2) \\ \sigma_{xv} / (\sigma_x^2 + \sigma_z^2) \end{bmatrix} \end{aligned}$$

In other words, the Kalman gain for  $x$  is

$$K_x = \frac{VAR(x)}{VAR(x) + VAR(z)}$$

This should be very familiar to you from the univariate case.

The Kalman gain for  $\dot{x}$  is

$$K_{\dot{x}} = \frac{COV(x, \dot{x})}{VAR(x) + VAR(z)}$$

What is the effect of this? Recall that we compute the state as

$$\begin{aligned} \mathbf{x} &= \bar{\mathbf{x}} + \mathbf{K}(z - \mathbf{H}\mathbf{x}) \\ &= \bar{\mathbf{x}} + \mathbf{K}y \end{aligned}$$

Here the residual  $y$  is a scalar. Therefore it is multiplied into each element of  $\mathbf{K}$ . Therefore we have

$$\begin{bmatrix} x \\ \dot{x} \end{bmatrix} = \begin{bmatrix} \bar{x} \\ \dot{\bar{x}} \end{bmatrix} + \begin{bmatrix} K_x \\ K_{\dot{x}} \end{bmatrix} y$$

Which gives this system of equations:

$$x = \bar{x} + yK_x \quad \dot{x} = \dot{\bar{x}} + yK_{\dot{x}}$$

The prediction  $\bar{x}$  was computed as  $x + \bar{x}\Delta t$ . If the prediction was perfect then the residual will be  $y = 0$  (ignoring noise in the measurement) and the velocity estimate will be unchanged. On the other hand, if the velocity estimate was very bad then the prediction will be very bad, and the residual will be large:  $y \gg 0$ . In this case we update the velocity estimate with  $yK_{\dot{x}}$ .  $K_{\dot{x}}$  is proportional to  $COV(x, \dot{x})$ . Therefore the velocity is updated by the error in the position times the a value proportional to the covariance between the position and velocity. The higher the correlation the larger the correction.

To bring this full circle,  $COV(x, \dot{x})$  is the off-diagonal elements of  $\mathbf{P}$ . Recall that those values were computed with  $\mathbf{F}\mathbf{P}\mathbf{F}^T$ . So the covariance of position and velocity is computed during the predict step. The Kalman gain for the velocity is proportional to this covariance, and we adjust the velocity estimate based on how inaccurate it was during the last epoch times a value proportional to this covariance.

In summary, these linear algebra equations may be very unfamiliar to you, but computation is actually very simple. It is essentially the same computation that we performed in the g-h filter. Our constants are different in this chapter because we take the noise in the process model and sensors into account, but the math is the same.

## 6.8 Adjusting the Filter

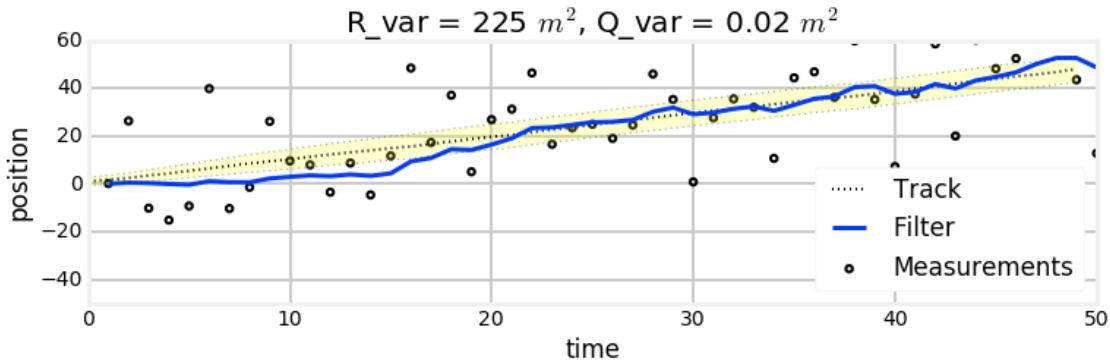
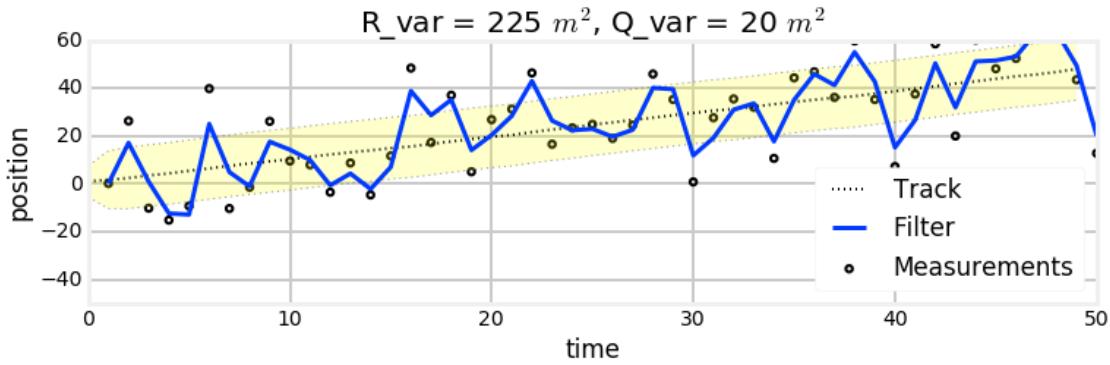
Let's start varying our parameters to see the effect of various changes. This is a very normal thing to be doing with Kalman filters. It is difficult, and often impossible to exactly model our sensors. An imperfect model means imperfect output from our filter. Engineers spend a lot of time tuning Kalman filters so that they perform well with real world sensors. We will spend time now to learn the effect of these changes. As you learn the effect of each change you will develop an intuition for how to design a Kalman filter. Designing a Kalman filter is as much art as science. We are modeling a physical system using math, and models are imperfect.

Let's look at the effects of the measurement noise  $\mathbf{R}$  and process noise  $\mathbf{Q}$ . We will want to see the effect of different settings for  $\mathbf{R}$  and  $\mathbf{Q}$ , so I have hard coded our measurements in `zs` based on a variance of 225

meters squared. That is very large, but it magnifies the effects of various design choices on the graphs, making it easier to recognize what is happening. Our first experiment holds  $\mathbf{R}$  constant while varying  $\mathbf{Q}$ .

```
In [34]: from numpy.random import seed
seed(2)
data = compute_dog_data(z_var=225, process_var=.02, count=50)

run(data=data, R=225, Q=200, P=P, plot_P=False,
     title='R_var = 225 $m^2$, Q_var = 20 $m^2$')
run(data=data, R=225, Q=.02, P=P, plot_P=False,
     title='R_var = 225 $m^2$, Q_var = 0.02 $m^2$');
```



The filter in the first plot should follow the noisy measurement closely. In the second plot the filter should vary from the measurement quite a bit, and be much closer to a straight line than in the first graph. Why does  $\mathbf{Q}$  affect the plots this way?

Let's remind ourselves of what the term process uncertainty means. Consider the problem of tracking a ball. We can accurately model its behavior in a vacuum with math, but with wind, varying air density, temperature, and an spinning ball with an imperfect surface our model will diverge from reality.

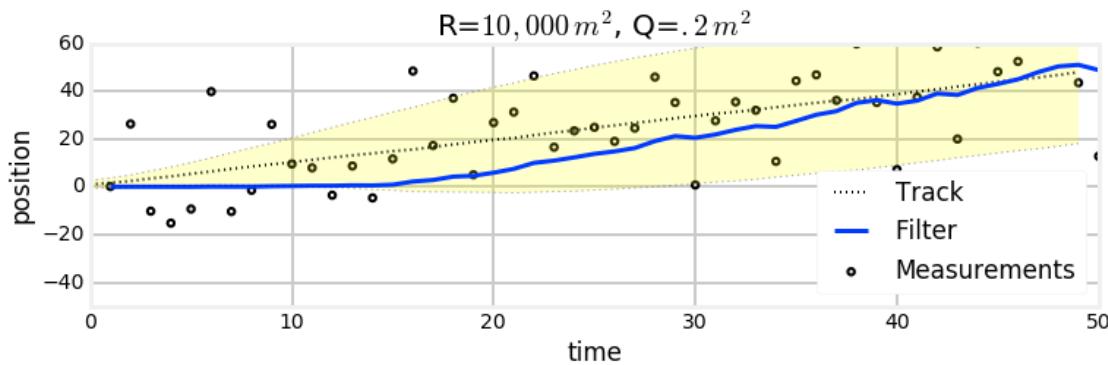
In the first case we set  $\mathbf{Q}_{var}=20 \text{ m}^2$ , which is quite large. In physical terms this is telling the filter "I don't trust my motion prediction step" as we are saying that the variance in the velocity is 20. Strictly speaking, we are telling the filter there is a lot of external noise that we are not modeling with  $\mathbf{F}$ , but the upshot of

that is to not trust the motion prediction step. The filter will be computing velocity ( $\dot{x}$ ), but then mostly ignoring it because we are telling the filter that the computation is extremely suspect. Therefore the filter has nothing to trust but the measurements, and thus it follows the measurements closely.

In the second case we set  $Q_{var}=0.02 \text{ m}^2$ , which is quite small. In physical terms we are telling the filter “trust the prediction, it is really good!”. More strictly this actually says there is very small amounts of process noise (variance  $0.02 \text{ m}^2$ ), so the process model is very accurate. So the filter ends up ignoring some of the measurement as it jumps up and down, because the variation in the measurement does not match our trustworthy velocity prediction.

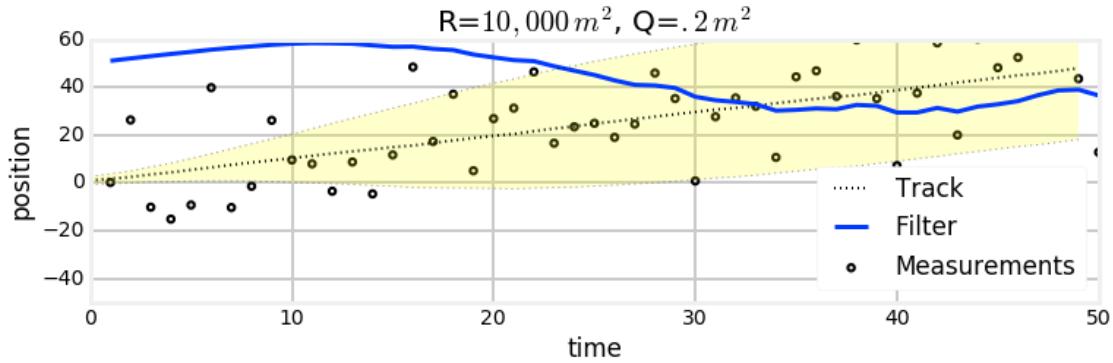
Now let's set  $Q_{var}$  to  $0.2 \text{ m}^2$ , and bump  $R_{var}$  up to  $10,000 \text{ m}^2$ . This is telling the filter that the measurement noise is very large.

```
In [35]: run(data=data, R=10000, Q=.2, P=P, plot_P=False,
           title='R=$10,000\, m^2, Q=$.2\, m^2');
```



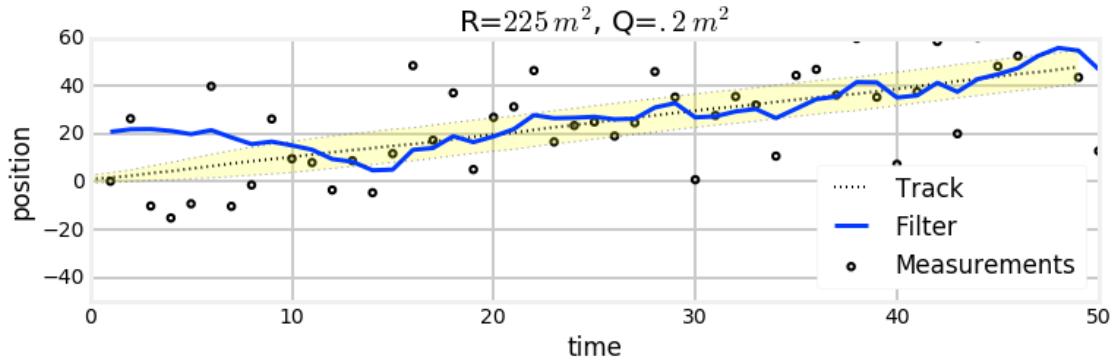
The effect of this can be subtle. We have created an suboptimal filter because the actual measurement noise variance is  $225 \text{ m}^2$ , not  $10,000 \text{ m}^2$ . By setting the filter's noise variance so high we force the filter to favor the prediction over the measurement. This can lead to apparently very smooth and good looking results. In the chart above the track may look extremely good to you since it follows the ideal path very closely. But, the ‘great’ behavior at the start should give you pause - the filter has not converged yet ( $P$  is still large) so it should not be able to be so close to the actual position. We can see that  $P$  has not converged because the entire chart is colored with the yellow background denoting the size of  $P$ . Let's see the result of a bad initial guess for the position by guessing the initial position to be 50 m and the initial velocity to be 1 m/s.

```
In [36]: run(data=data, R=10000, Q=.2, P=P, plot_P=False,
           x0=np.array([50., 1.]),
           title='R=$10,000\, m^2, Q=$.2\, m^2');
```



Here we can see that the filter cannot acquire the track. This happens because even though the filter is getting reasonably good measurements it assumes that the measurements are bad, and eventually predicts forward from a bad position at each step. If you think that perhaps that bad initial position would give similar results for a smaller measurement noise, let's set it back to the correct value of  $225 \text{ m}^2$ .

```
In [37]: run(data=data, R=225, Q=.2, P=P, plot_P=False,
           x0=np.array([20., 1.]),
           title='R=$225\text{ m}^2, Q=$.2\text{ m}^2');
```



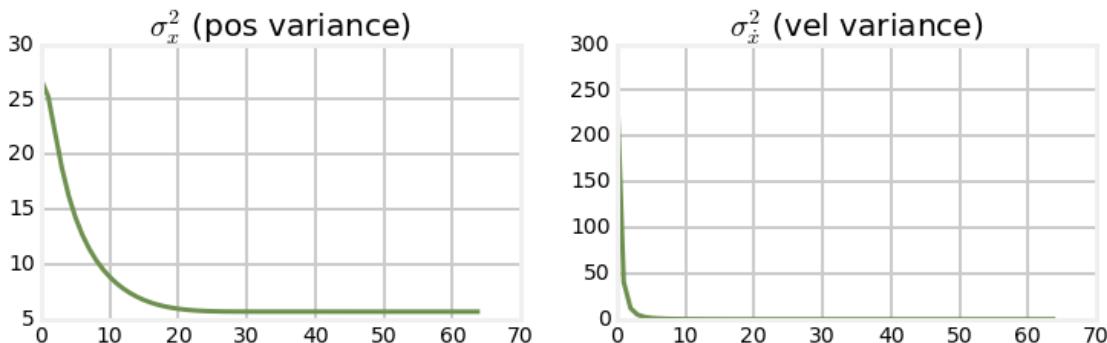
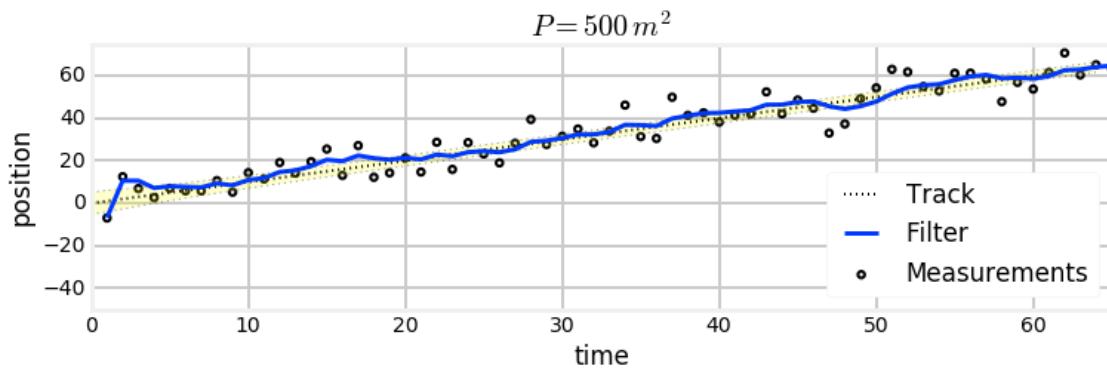
Here we see that the filter initially struggles for several iterations to acquire the track, but then it accurately tracks our dog. In fact, this is nearly optimum - we have not designed  $\mathbf{Q}$  optimally, but  $\mathbf{R}$  is optimal. A rule of thumb for  $\mathbf{Q}$  is to set it between  $\frac{1}{2}\Delta a$  to  $\Delta a$ , where  $\Delta a$  is the maximum amount that the acceleration will change between sample period. This only applies for the assumption we are making in this chapter - that acceleration is constant and uncorrelated between each time period. In the Kalman Math chapter we will discuss several different ways of designing  $\mathbf{Q}$ .

To some extent you can get similar looking output by varying either  $\mathbf{R}$  or  $\mathbf{Q}$ , but I urge you to not ‘magically’ alter these until you get output that you like. Always think about the physical implications of these assignments, and vary  $\mathbf{R}$  and/or  $\mathbf{Q}$  based on your knowledge of the system you are filtering. Back that up with extensive simulations and/or trial runs of real data.

## 6.9 A Detailed Examination of the Covariance Matrix

Let's start by revisiting plotting a track. I will hard code the data and noise to avoid being at the mercy of the random number generator, which might generate data that does not illustrate what I want to talk about. I will start with  $P=500$ .

```
In [38]: import mkf_internal
var = 27.5
data = mkf_internal.zs_var_275()
run(data=data, R=var, Q=.02, P=500., plot_P=True,
     title='P=500 m^2');
```



Looking at the output we see a very large spike in the filter output at the beginning. We set  $P = 500 \mathbf{I}_2$  (this is shorthand notation for a  $2 \times 2$  diagonal matrix with 500 in the diagonal). We now have enough information to understand what this means, and how the Kalman filter treats it. The 500 in the upper left hand corner corresponds to  $\sigma_x^2$ ; therefore we are saying the standard deviation of  $x$  is  $\sqrt{500}$ , or roughly 22.36 m. Roughly 99% of the samples occur within  $3\sigma$ , therefore  $\sigma_x^2 = 500$  is telling the Kalman filter that the prediction (the prior) could be up to 67 meters off. That is a large error, so when the measurement spikes the Kalman filter distrusts its own estimate and jumps wildly to try to incorporate the measurement. Then, as the filter evolves  $\mathbf{P}$  quickly converges to a more realistic value.

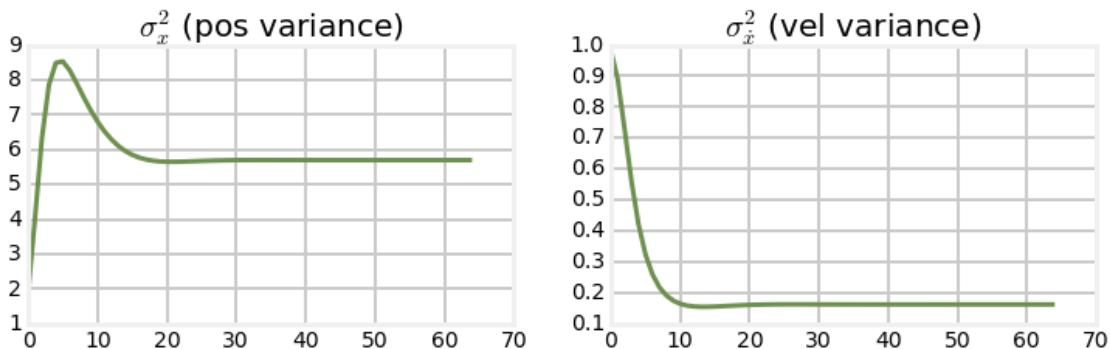
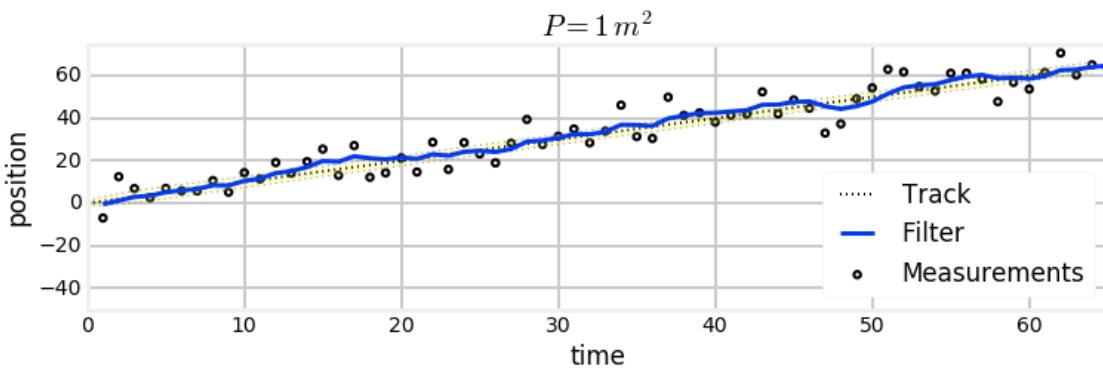
Let's look at the math behind this. The equation for the Kalman gain is

$$\mathbf{K} = \bar{\mathbf{P}}\mathbf{H}^T\mathbf{S}^{-1} \approx \frac{\bar{\mathbf{P}}\mathbf{H}^T}{\mathbf{S}} \approx \frac{\text{uncertainty}_{\text{prediction}}}{\text{uncertainty}_{\text{measurement}}} \mathbf{H}^T$$

It is a ratio of the uncertainty of the prediction vs measurement. Here the uncertainty in the prediction is large, so  $\mathbf{K}$  is large (near 1 if this was a scalar).  $\mathbf{K}$  is multiplied by the residual  $\mathbf{y} = \mathbf{z} - \mathbf{H}\bar{\mathbf{x}}$ , which is the measurement minus the prediction, so a large  $\mathbf{K}$  favors the measurement. Therefore if  $\mathbf{P}$  is large relative to the sensor uncertainty  $\mathbf{R}$  the filter will form most of the estimate from the measurement.

Now let us see the effect of a smaller initial value of  $\mathbf{P} = 1.0 \mathbf{I}_2$ .

```
In [39]: run(data=data, R=var, Q=.02, P=1., plot_P=True,
          title='\$P=1\\, m^2$');
```



This looks good at first blush. The plot does not have the spike that the former plot did; the filter starts tracking the measurements and doesn't take any time to settle to the signal. However, if we look at the plots for  $\mathbf{P}$  you can see that there is an initial spike for the variance in position, and that it never really converges. Poor design leads to a long convergence time, and suboptimal results.

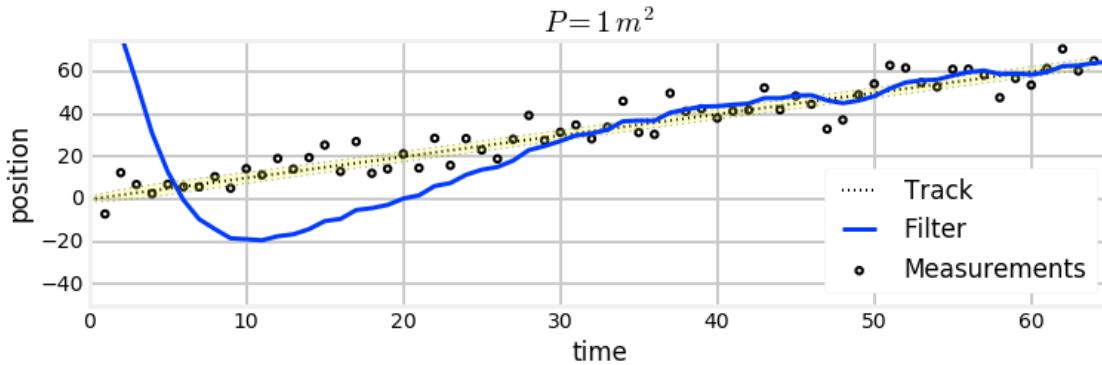
So despite the filter tracking very close to the actual signal we cannot conclude that the ‘magic’ is to use a small  $\mathbf{P}$ . Yes, this will avoid having the Kalman filter take time to accurately track the signal, but if we are truly uncertain about the initial measurements this can cause the filter to generate very bad results. If we are tracking a living object we are probably very uncertain about where it is before we start tracking it. On the other hand, if we are filtering the output of a thermometer, we are as certain about the first

measurement as the 1000th. For your Kalman filter to perform well you must set  $\mathbf{P}$  to a value that truly reflects your knowledge about the data.

Let's see the result of a bad initial estimate coupled with a very small  $\mathbf{P}$ . We will set our initial estimate at  $x = 100$  m (whereas the dog actually starts at 0m), but set  $P=1$  m<sup>2</sup>. This is clearly an incorrect value for  $\mathbf{P}$  as the estimate is off by 100 m but we tell the filter that it the  $3\sigma$  error is 3 m.

```
In [40]: x = np.array([100., 0.]).T
```

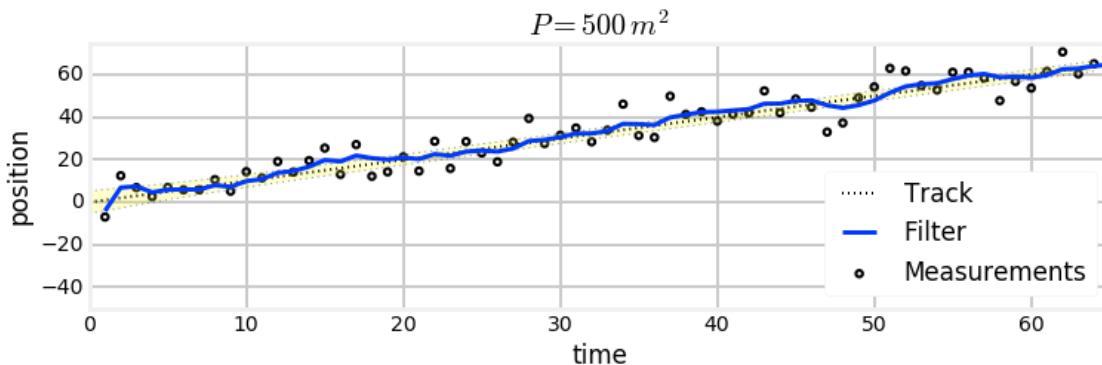
```
run(data=data, R=var, Q=.02, P=1., x0=x,
plot_P=False, title='P=1\, m^2');
```



We can see that the initial estimates are terrible and that it takes the filter a long time to start converging onto the signal . This is because we told the Kalman filter that we strongly believe in our initial estimate of 100 m and were incorrect in that belief.

Now, let's provide a more reasonable value for  $\mathbf{P}$  and see the difference.

```
In [41]: x = np.array([100., 0.]).T
run(data=data, R=var, Q=.02, P=500., x0=x,
plot_P=False, title='P=500\, m^2');
```



In this case the Kalman filter is very uncertain about the initial state, so it converges onto the signal much faster. It is producing good output after only 5 to 6 epochs. With the theory we have developed so far this

is about as good as we can do. However, this scenario is a bit artificial; if we do not know where the object is when we start tracking we do not initialize the filter to some arbitrary value, such as 0 m or 100 m. I address this in the **Filter Initialization** section below.

Lets do another Kalman filter for our dog, and this time plot the covariance ellipses on the same plot as the position.

```
In [42]: from mkf_internal import plot_track_ellipses

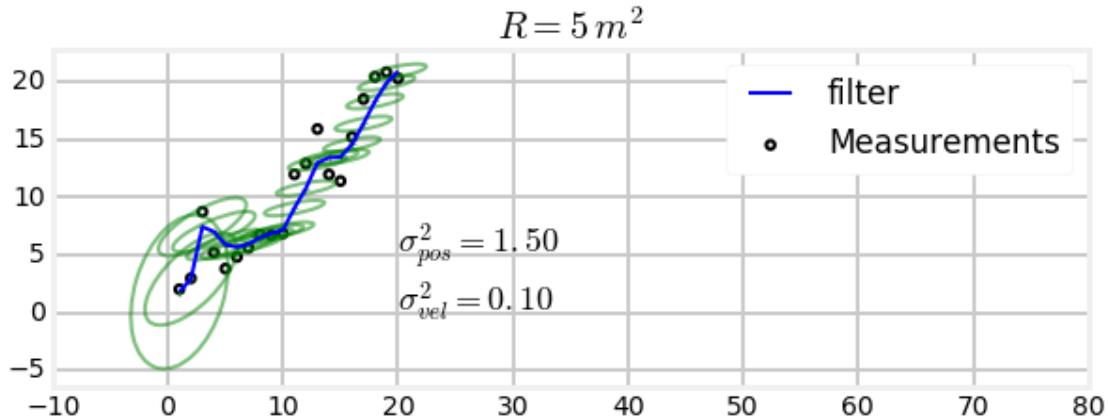
def plot_covariances(count, R, Q=0, P=20., title=''):
    data = compute_dog_data(R, Q, count)
    f = pos_vel_filter(x=(0., 0.), R=R, Q=Q, P=P)

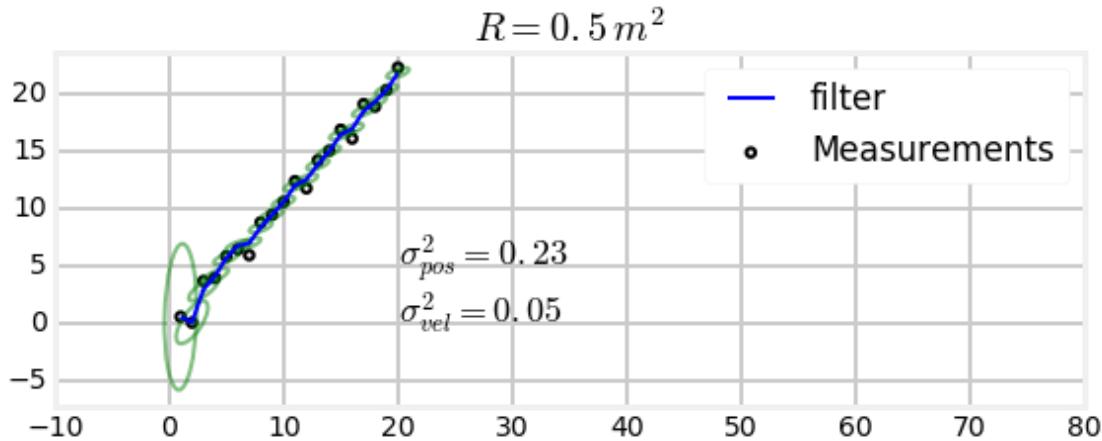
    xs, cov = [], []
    for z in data[:, 1]:
        f.predict()
        f.update(z)

    xs.append(f.x[0])
    cov.append(f.P)

    plot_track_ellipses(count, data[:, 1], xs, cov, title)

plot_covariances(R=5, Q=.02, count=20, title='$R = 5\text{ m}^2$')
plot_covariances(R=.5, Q=.02, count=20, title='$R = 0.5\text{ m}^2$')
```





If you are viewing this in Jupyter Notebook or on the web, here is an animation of the filter filtering the data. I've tuned the filter parameters such that it is easy to see a change in  $\mathbf{P}$  as the filter progresses.

The output on these is a bit messy, but you should be able to see what is happening. In both plots we are drawing the covariance matrix for each point. We start with the covariance  $\mathbf{P} = \begin{pmatrix} 20 & 0 \\ 0 & 20 \end{pmatrix}$ , which signifies a lot of uncertainty about our initial belief. After we receive the first measurement the Kalman filter updates this belief, and so the variance is no longer as large. In the top plot the first ellipse (the one on the far left) should be a slightly squashed ellipse. As the filter continues processing the measurements the covariance ellipse quickly shifts shape until it settles down to being a long, narrow ellipse tilted in the direction of movement.

Think about what this means physically. The x-axis of the ellipse denotes our uncertainty in position, and the y-axis our uncertainty in velocity. So, an ellipse that is taller than it is wide signifies that we are more uncertain about the velocity than the position. Conversely, a wide, narrow ellipse shows high uncertainty in position and low uncertainty in velocity. Finally, the amount of tilt shows the amount of correlation between the two variables.

The first plot, with  $R = 5\text{m}^2$ , finishes up with an ellipse that is wider than it is tall. If that is not clear I have printed out the variances for the last ellipse in the lower right hand corner.

In contrast, the second plot, with  $R=0.5 \text{ m}^2$ , has a final ellipse that is taller than wide. The ellipses in the second plot are all much smaller than the ellipses in the first plot. This stands to reason because a small  $\mathbf{R}$  implies a small amount of noise in our measurements. Small noise means accurate predictions, and thus a strong belief in our position.

## 6.10 Question: Explain Ellipse Differences

Why are the ellipses for  $\mathbf{R} = 5\text{m}^2$  shorter, and more tilted than the ellipses for  $\mathbf{R} = 0.5\text{m}^2$ . Hint: think about this in the context of what these ellipses mean physically, not in terms of the math. If you aren't sure about the answer, change  $\mathbf{R}$  to truly large and small numbers such as  $100 \text{ m}^2$  and  $0.1 \text{ m}^2$ , observe the changes, and think about what this means.

### 6.10.1 Solution

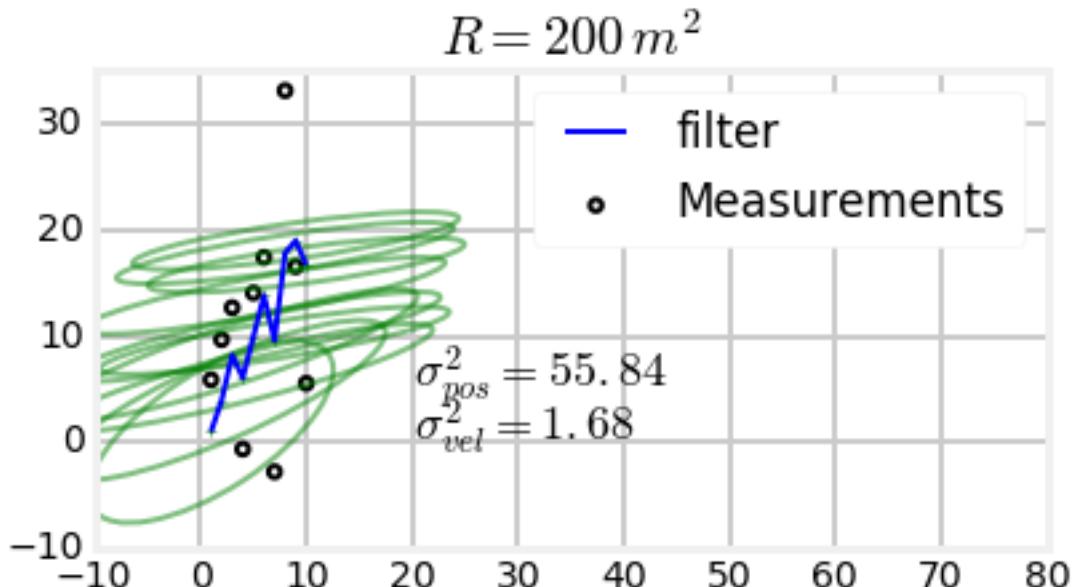
The x-axis is for position, and x-axis is velocity. An ellipse that is vertical, or nearly so, says there is no correlation between position and velocity, and an ellipse that is diagonal says that there is a lot of correlation.

Phrased that way, the results sound unlikely. The tilt of the ellipse changes, but the correlation shouldn't be changing over time. But this is a measure of the output of the filter, not a description of the actual, physical world. When  $\mathbf{R}$  is very large we are telling the filter that there is a lot of noise in the measurements. In that case the Kalman gain  $\mathbf{K}$  is set to favor the prediction over the measurement, and the prediction comes from the velocity state variable. Thus there is a large correlation between  $x$  and  $\dot{x}$ . Conversely, if  $\mathbf{R}$  is small, we are telling the filter that the measurement is very trustworthy, and  $\mathbf{K}$  is set to favor the measurement over the prediction. Why would the filter want to use the prediction if the measurement is nearly perfect? If the filter is not using much from the prediction there will be very little correlation reported.

**This is a critical point to understand!** The Kalman filter is a mathematical model for a real world system. A report of little correlation does not mean there is no correlation in the physical system, just that there was no linear correlation in the mathematical model. It's a report of how much measurement vs prediction was incorporated into the model.

Let's bring that point home with a truly large measurement error. We will set  $\mathbf{R} = 200 \text{ m}^2$ . Think about what the plot will look like before scrolling down.

```
In [43]: plot_covariances(R=200., Q=.2, count=10, title='$R = 200\, \text{m}^2$')
```



I hope the result was what you were expecting. The ellipse quickly became very wide and not very tall. It did this because the Kalman filter mostly used the prediction vs the measurement to produce the filtered result. We can also see how the filter output is slow to acquire the track. The Kalman filter assumes that the measurements are extremely noisy, and so it is very slow to update its estimate for  $\dot{x}$ .

Keep looking at these plots until you grasp how to interpret the covariance matrix  $\mathbf{P}$ . When you work with a  $9 \times 9$  matrix it may seem overwhelming - there are 81 numbers to interpret. Just break it down - the diagonal contains the variance for each state variable, and all off diagonal elements are the product of two variances and a scaling factor  $p$ . You cannot plot a  $9 \times 9$  matrix on the screen so you have to develop your intuition and understanding in this simple, 2-D case.

When plotting covariance ellipses, make sure to always use `ax.set_aspect('equal')` or `plt.axis('equal')` in your code (the former lets you set the `xlim` and `ylim` values). If the axis

use different scales the ellipses will be drawn distorted. For example, the ellipse may be drawn as being taller than it is wide, but it may actually be wider than tall.

## 6.11 Filter Initialization

There are many schemes for initializing the filter (i.e. choosing the initial values for  $\mathbf{x}$  and  $\mathbf{P}$ ). I will share a common approach that performs well in most situations. In this scheme you do not initialize the filter until you get the first measurement,  $\mathbf{z}_0$ . From this you can compute the initial value for  $\mathbf{x}$  with  $\mathbf{x}_0 = \mathbf{z}_0$ . If  $\mathbf{z}$  is not of the same size, type, and units as  $\mathbf{x}$ , which is usually the case, we can use our measurement function as follow.

We know

$$\mathbf{z} = \mathbf{H}\mathbf{x}$$

Hence,

$$\begin{aligned}\mathbf{H}^{-1}\mathbf{H}\mathbf{x} &= \mathbf{H}^{-1}\mathbf{z} \\ \mathbf{x} &= \mathbf{H}^{-1}\mathbf{z}\end{aligned}$$

Matrix inversion requires a square matrix, but  $\mathbf{H}$  is rarely square. SciPy will compute the Moore-Penrose pseudo-inverse of a matrix with `scipy.linalg.pinv`, so your code might look like

```
In [44]: import scipy.linalg as la

H = np.array([[1, 0.]])
z0 = 3.2

x = np.dot(la.pinv(H), z0)
print(x)

[[ 3.2]
 [ 0. ]]
```

Specialized knowledge of your problem domain may lead you to a different computation, but this is one way to do it.

Now we need to compute a value for  $\mathbf{P}$ . This will vary by problem, but in general you will use the measurement error  $\mathbf{R}$  for identical terms, and maximum values for the rest of the terms. Maybe that isn't clear. In this chapter we have been tracking an object using position and velocity as the state, and the measurements have been positions. In that case we would initialize  $\mathbf{P}$  with

$$\mathbf{P} = \begin{bmatrix} \mathbf{R}_0 & 0 \\ 0 & vel_{max}^2 \end{bmatrix}$$

The diagonal of  $\mathbf{P}$  contains the variance of each state variable, so we populate it with reasonable values.  $\mathbf{R}$  is a reasonable variance for the position, and the maximum velocity squared is a reasonable variance for the velocity. It is squared because variance is squared:  $\sigma^2$ .

You really need to understand the domain in which you are working and initialize your filter on the best available information. For example, suppose we were trying to track horses in a horse race. The initial measurements might be very bad, and provide you with a position far from the starting gate. We know that the horse must start at the starting gate; initializing the filter to the initial measurement would lead to suboptimal results. In this scenario we would want to always initialize the Kalman filter with the starting gate position of the horse.

## 6.12 Batch Processing

The Kalman filter is designed as a recursive algorithm - as new measurements come in we immediately create a new estimate. But it is very common to have a set of data that have been already collected which we want to filter. Kalman filters can always be run in a `batch` mode, where all of the measurements are filtered at once. We have implemented this in `KalmanFilter.batch_filter()`. Internally, all the function does is loop over the measurements and collect the resulting state and covariance estimates in arrays. It simplifies your logic and conveniently gathers all of the outputs into arrays.

First collect your measurements into an array or list. Maybe it is in a CSV file, for example.

```
zs = read_altitude_from_csv('altitude_data.csv')
```

Or maybe you will generate it using a generator:

```
zs = [some_func(i) for i in range(1000)]
```

Then call the `batch_filter()` method.

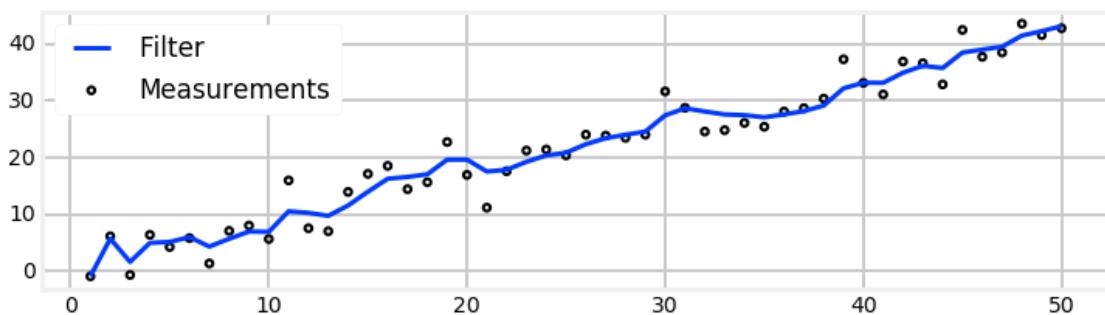
```
Xs, Ps, Xs_pred, Ps_pred = kf.filter.batch_filter(zs)
```

The function takes the list/array of measurements, filters it, and returns an NumPy array of state estimates ( $X_s$ ), covariance matrices ( $P_s$ ), and the predictions for the same ( $X_{s\_pred}$ ,  $P_{s\_pred}$ ).

Here is a complete example.

```
In [45]: count = 50
data = compute_dog_data(10, .2, count)
zs = data[:, 1]
P = np.diag([500., 49.])
f = pos_vel_filter(x=(0., 0.), R=3., Q=.02, P=P)
xs, _, _, _ = f.batch_filter(zs)

book_plots.plot_measurements(range(1, count + 1), zs)
book_plots.plot_filter(range(1, count + 1), xs[:, 0])
plt.legend(loc='best');
```



## 6.13 Smoothing the Results

I have an entire chapter on using the Kalman filter to smooth data; I will not repeat the chapter's information here. However, it is so easy to use, and offers such a profoundly improved output that I will tease you will a few examples. The smoothing chapter is not especially difficult; you are sufficiently prepared to read it now.

Let's assume that we are tracking a car that has been traveling in a straight line. We get a measurement that implies that the car is starting to turn to the left. The Kalman filter moves the state estimate somewhat towards the measurement, but it cannot judge whether this is a particularly noisy measurement or the true start of a turn.

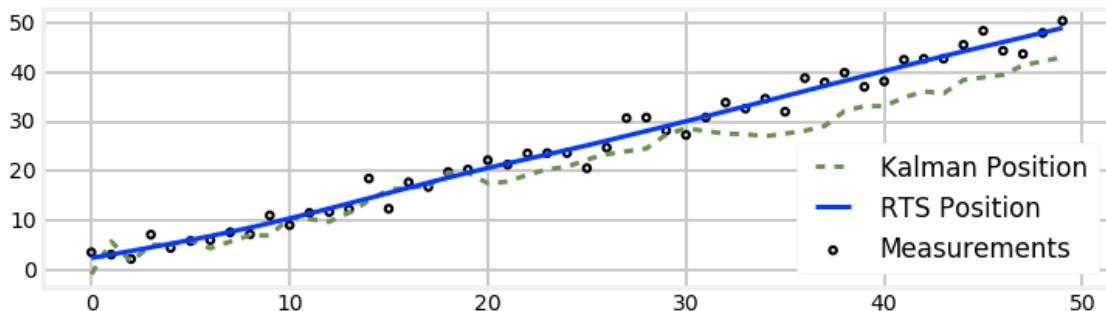
However, if we have future measurements we can decide if a turn was made. Suppose the subsequent measurements all continue turning left. We can then be sure that that a turn was initiated. On the other hand, if the subsequent measurements continued on in a straight line we would know that the measurement was noisy and should be mostly ignored. Instead of making an estimate part way between the measurement and prediction the estimate will either fully incorporate the measurement or ignore it, depending on what the future measurements imply about the object's movement.

`KalmanFilter` implements a form of this algorithm which is called an RTS smoother, named after the inventors of the algorithm: Rauch, Tung, and Striebel. The method is `rts_smoother()`. To use it pass in the means and covariances computed from the `batch_filter` step, and receive back the smoothed means, covariances, and Kalman gain.

```
In [46]: from numpy.random import seed
count = 50
seed(8923)

P = np.diag([500., 49.])
f = pos_vel_filter(x=(0., 0.), R=3., Q=.02, P=P)
data = compute_dog_data(3., .02, count)
zs = data[:, 1]
Xs, Covs, _, _ = f.batch_filter(zs)
Ms, Ps, _ = f.rts_smoothen(Xs, Covs)

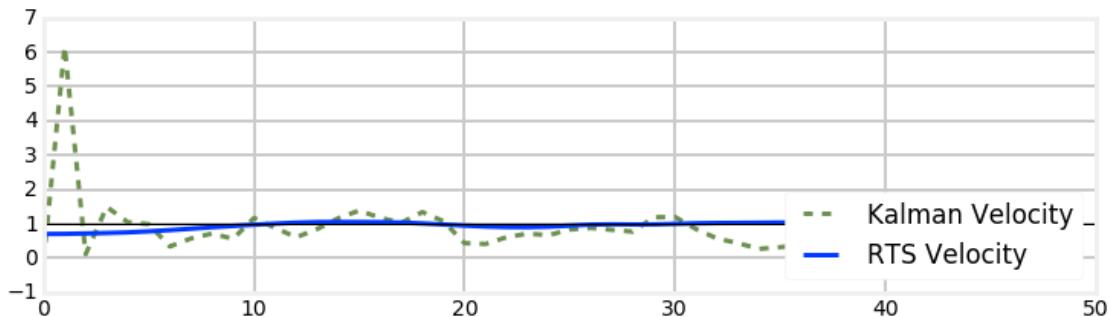
book_plots.plot_measurements(zs)
plt.plot(xs[:, 0], ls='--', label='Kalman Position')
plt.plot(Ms[:, 0], label='RTS Position')
plt.legend(loc=4);
```



This output is fantastic! Two things are very apparent to me in this chart. First, the RTS output is much, much smoother than the KF output. Second, it is almost always more accurate than the KF output (we

will examine this claim in detail in the **Smoothing** chapter). The improvement in the velocity, which is an hidden variable, is even more dramatic:

```
In [47]: plt.plot(xs[:, 1], ls='--', label='Kalman Velocity')
    plt.plot(Ms[:, 1], label='RTS Velocity')
    plt.legend(loc=4)
    plt.gca().axhline(1, lw=1, c='k');
```



We will explore why this is so in the next exercise.

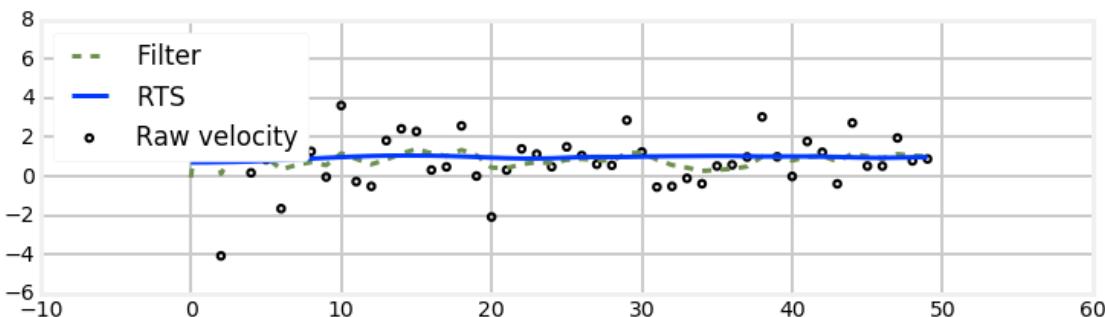
## 6.14 Exercise: Compare Velocities

Since we are plotting velocities let's look at what the the ‘raw’ velocity is, which we can compute by subtracting subsequent measurements. i.e the velocity at time 1 can be approximated by  $xs[1] - xs[0]$ . Plot the raw value against the values estimated by the Kalman filter and by the RTS filter. Discuss what you see.

```
In [48]: # your code here
```

### 6.14.1 Solution

```
In [49]: dx = np.diff(xs[:, 0], axis=0)
    plt.scatter(range(1, count), dx, facecolor='none', edgecolor='k',
                lw=2, label='Raw velocity')
    plt.plot(xs[:, 1], ls='--', label='Filter')
    plt.plot(Ms[:, 1], label='RTS')
    plt.legend(loc='best');
```



We see that the noise swamps the signal, causing the raw values to be essentially worthless. The filter is maintaining a separate estimate for velocity. The Kalman gain  $\mathbf{K}$  is multidimensional. For example, it might have the value  $\mathbf{K} = [0.1274, 0.843]^\top$ . the first value is used to scale the residual of the position, and the second value will scale the residual of the velocity. The covariance matrix tells the filter how correlated the position and velocity are, and each will be optimally filtered.

I show this to reiterate the importance of using Kalman filters to compute velocities, accelerations, and even higher order values. I use a Kalman filter even when my measurements are so accurate that I am willing to use them unfiltered because it allows me accurate estimates for velocities and accelerations.

## 6.15 Discussion and Summary

We extended Gaussians into multiple dimensions to allow us to simultaneously handle multiple dimensions, both spacial and others (velocity, etc). This led us to use linear algebra. We made a key insight: hidden variables have the ability to significantly increase the accuracy of the filter. This counterintuitive result is because the hidden variables are correlated with the observed variables. Intuitively, I think of it as triangulating onto the solution - only a small number of states can be true for a given combination of a velocity Gaussian and position Gaussian, roughly described by their intersection. But of course this works for non-spatial problems as well.

I gave an intuitive definition of observability. Observability was invented by Dr. Kalman for linear systems, and there is a fair amount of theory behind it. It answers the question of whether a system state can be determined by observing the system's output. For our problems this has been easy to determine, but more complicated systems may require rigorous analysis. Wikipedia's [Observability](#) article has an overview; if you need to learn the topic [Grewal2008] is a good source.

There is one important caveat about hidden variables. It is easy to construct a filter that produces estimates for hidden variables. I could write a filter that estimates the color of a tracked car. But there is no way to compute car color from positions, so the estimate for the color will be nonsense. The designer must verify that these variables are being estimated correctly. If you do not have a velocity sensor and yet are estimating velocity, you will need to test that the velocity estimates are correct.; do not trust that they are. For example, suppose the velocity has a periodic component to it - it looks like a sine wave. If your sample time is less than 2 times the frequency you will not be able to accurately estimate the velocity (due to Nyquist's Theorem). Imagine that the sample period is equal to the frequency of the velocity. The filter will report that the velocity is constant because each time it will be sampling the system at the same point on the sin wave.

Initialization poses a particularly difficult problem for hidden variables. If you start with a bad initialization the filter can usually recover the observed variables, but may struggle and fail with the hidden one. Estimating hidden variables is a powerful tool, but a dangerous one.

I established a series of steps for designing a Kalman filter. These are not a usual part of the Kalman filter literature, and are only meant as a guide, not a prescription. Designing for a hard problem is an iterative process. You make a guess at the state vector, work out what your measurement and state models are, run some tests, and then alter the design as necessary.

The design of  $\mathbf{R}$  and  $\mathbf{Q}$  is often quite challenging. I've made it appear to be quite scientific. Your sensor has Gaussian noise of  $\mathcal{N}(0, \sigma^2)$ , so set  $\mathbf{R} = \sigma^2$ . Easy! But I told you a dirty lie. Sensors are not Gaussian. We started the book with a bathroom scale. Suppose  $\sigma = 1$  kg, and you try to weigh something that weighs 0.5 kg. Theory tells us we will get negative measurements, but of course the scale will never report weights less than zero. Real world sensors typically have fat tails (known as kurtosis) and skew. In some cases, such as with the scale, one or both tails are truncated.

The case with  $\mathbf{Q}$  is more dire. I hope you were skeptical when I blithely assigned a noise matrix to my

prediction about the movements of a dog. Who can say what a dog will do next? The Kalman filter in my GPS doesn't know about hills, the outside winds, or my terrible driving skills. Yet the filter requires a precise number to encapsulate all of that information, and it needs to work while I drive off-road in the desert, and when a Formula One champion drives on a closed circuit track.

These problems led some researchers and engineers to derogatorily call the Kalman filter a 'ball of mud'. In other words, it doesn't always hold together so well. Another term to know - Kalman filters can become smug. Their estimates are based solely on what you tell it the noises are. Those values can lead to overly confident estimates.  $\mathbf{P}$  gets smaller and smaller while the filter is actually becoming more and more inaccurate! In the worst case the filter diverges. We will see a lot of that when we start studying nonlinear filters.

The Kalman filter is a mathematical model of the world. The output is only as accurate as that model. To make the math tractable we had to make some assumptions. We assume that the sensors and motion model have Gaussian noise. We assume that everything is linear. If that is true, the Kalman filter is optimal in a least squares sense. This means that there is no way to make a better estimate than what the filter gives us. However, these assumption are almost never true, and hence the model is necessarily limited, and a working filter is rarely optimal.

In later chapters we will deal with the problem of nonlinearity. For now I want you to understand that designing the matrices of a linear filter is an experimental procedure more than a mathematical one. Use math to establish the initial values, but then you need to experiment. If there is a lot of unaccounted noise in the world (wind, etc) you may have to make  $\mathbf{Q}$  larger. If you make it too large the filter fails to respond quickly to changes. In the **Adaptive Filters** chapter you will learn some alternative techniques that allow you to change the filter design in real time in response to the inputs and performance, but for now you need to find one set of values that works for the conditions your filter will encounter. Noise matrices for an acrobatic plane might be different if the pilot is a student than if the pilot is an expert as the dynamics will be quite different.

## 6.16 References

- [1] 'Kalman Filters'. Wikipedia [https://en.wikipedia.org/wiki/Kalman\\_filter#Details](https://en.wikipedia.org/wiki/Kalman_filter#Details)
- [Grewal2008] Grewal, Mohinder S., Andrews, Angus P. Kalman Filtering: Theory and Practice Using MATLAB. Third Edition. John Wiley & Sons. 2008.

# Chapter 7

## Kalman Filter Math

If you've gotten this far I hope that you are thinking that the Kalman filter's fearsome reputation is somewhat undeserved. Sure, I hand waved some equations away, but I hope implementation has been fairly straightforward for you. The underlying concept is quite straightforward - take two measurements, or a measurement and a prediction, and choose the output to be somewhere between the two. If you believe the measurement more your guess will be closer to the measurement, and if you believe the prediction is more accurate your guess will lie closer to it. That's not rocket science (little joke - it is exactly this math that got Apollo to the moon and back!).

To be honest I have been choosing my problems carefully. For an arbitrary problem designing the Kalman filter matrices can be extremely difficult. I haven't been too tricky, though. Equations like Newton's equations of motion can be trivially computed for Kalman filter applications, and they make up the bulk of the kind of problems that we want to solve.

I have illustrated the concepts with code and reasoning, not math. But there are topics that do require more mathematics than I have used so far. This chapter presents the math that you will need for the rest of the book.

### 7.1 Modeling a Dynamic System

A dynamic system is a physical system whose state (position, temperature, etc) evolves over time. Calculus is the math of changing values, so we use differential equations to model dynamic systems. Some systems cannot be modeled with differential equations, but we will not encounter those in this book.

Modeling dynamic systems is properly the topic of several college courses. To an extent there is no substitute for a few semesters of ordinary and partial differential equations followed by a graduate course in control system theory. If you are a hobbyist, or trying to solve one very specific filtering problem at work you probably do not have the time and/or inclination to devote a year or more to that education.

Fortunately, I can present enough of the theory to allow us to create the system equations for many different Kalman filters. My goal is to get you to the stage where you can read a publication and understand it well enough to implement the algorithms. The background math is deep, but in practice we end up using a few simple techniques over and over again.

This is the longest section of pure math in this book. You will need to master everything in this section to understand the Extended Kalman filter (EKF), the workhorse nonlinear filter. I do cover more modern filters that do not require as much of this math. You can choose to skim now, and come back to this if you decide to learn the EKF.

We need to start by understanding the underlying equations and assumptions that the Kalman filter uses. We are trying to model real world phenomena, so what do we have to consider?

Each physical system has a process. For example, a car traveling at a certain velocity goes so far in a fixed amount of time, and its velocity varies as a function of its acceleration. We describe that behavior with the well known Newtonian equations that we learned in high school.

$$\begin{aligned} v &= at \\ x &= \frac{1}{2}at^2 + v_0t + x_0 \end{aligned}$$

Once we learned calculus we saw them in this form:

$$\mathbf{v} = \frac{d\mathbf{x}}{dt}, \quad \mathbf{a} = \frac{d\mathbf{v}}{dt} = \frac{d^2\mathbf{x}}{dt^2}$$

A typical automobile tracking problem would have you compute the distance traveled given a constant velocity or acceleration as we did in previous chapters. But, of course we know this is not all that is happening. No car travels on a perfect road. There are bumps that cause the car to slow down, there is wind drag, there are hills that raise and lower the speed. The suspension is a mechanical system with friction and imperfect springs. Gusts of wind alter the car's state.

Accurately modeling a system with linear equations is impossible except for the most trivial problems. So control theory is forced to make a simplification. At any time  $t$  we say that the true state (such as the position of our car) is the predicted value from the imperfect model plus some unknown process noise:

$$x(t) = x_{pred}(t) + noise(t)$$

This is not meant to imply that  $noise(t)$  is a function that we can derive analytically. It is merely a statement of fact - we can always describe the true value as the predicted value plus the process noise. "Noise" does not imply random events. If we are tracking a thrown ball in the atmosphere, and our model assumes the ball is in a vacuum, then the effect of air drag is process noise in this context.

In the next section we will learn techniques to convert a set of differential equations into a set of first-order differential equations. Assuming that, we can say that our model of the system without noise is:

$$\dot{\mathbf{x}} = \mathbf{Ax}$$

**A** is known as the systems dynamics matrix as it describes the dynamics of the system. Now we need to model the noise. We will call that **w**, and add it to the equation.

$$\dot{\mathbf{x}} = \mathbf{Ax} + \mathbf{w}$$

**w** may strike you as a poor choice for the name, but you will soon see that the Kalman filter assumes white noise.

Finally, we need to consider any inputs into the system. We assume an input **u**, and that there exists a linear model that defines how that input changes the system. For example, pressing the accelerator in your car makes it accelerate, and gravity causes balls to fall. Both are control inputs. We will need a matrix **B** to convert **u** into the effect on the system. We add that into our equation:

$$\dot{\mathbf{x}} = \mathbf{Ax} + \mathbf{Bu} + \mathbf{w}$$

And that's it. That is one of the equations that Dr. Kalman set out to solve, and he found an optimal estimator if we assume certain properties of **w**.

## 7.2 State-Space Representation of Dynamic Systems

In the last section we derived the equation

$$\dot{\mathbf{x}} = \mathbf{Ax} + \mathbf{Bu} + \mathbf{w}$$

However, for our filters we are not interested in the derivative of  $\mathbf{x}$ , but in  $\mathbf{x}$  itself. Ignoring the noise for a moment, we want an equation that recursively finds the value of  $\mathbf{x}$  at time  $t_k$  in terms of  $\mathbf{x}$  at time  $t_{k-1}$ :

$$\mathbf{x}(t_k) = \mathbf{F}(\Delta t)\mathbf{x}(t_{k-1}) + \mathbf{B}(t_k) + \mathbf{u}(t_k)$$

Convention allows us to write  $\mathbf{x}(t_k)$  as  $\mathbf{x}_k$ , which means the the value of  $\mathbf{x}$  at the  $k^{th}$  value of  $t$ .

$$\mathbf{x}_k = \mathbf{Fx}_{k-1} + \mathbf{B}_k \mathbf{u}_k$$

$\mathbf{F}$  is the familiar state transition matrix, named due to its ability to transition the state from the previous time step to the current time step. It is very similar to the system dynamics matrix  $\mathbf{A}$ . The difference is that the system dynamics matrix  $\mathbf{A}$  models a set of linear differential equations, and is continuous.  $\mathbf{F}$  is discrete, and represents a set of linear equations (not differential equations) which step transition  $\mathbf{x}_{k-1}$  to  $\mathbf{x}_k$  over a discrete time step  $\Delta t$ .

Normally finding this equation is quite difficult. The equation  $\dot{x} = v$  is the simplest possible differential equation and we trivially integrate it as:

$$\int_{x_{k-1}}^{x_k} dx = \int_0^{\Delta t} v dt x_k - x_0 = v \Delta t x_k = v \Delta t + x_0$$

This equation is recursive: we compute the value of  $x$  at time  $t$  based on its value at time  $t-1$ . This recursive form enables us to represent the system (process model) in the form required by the Kalman filter:

$$\begin{aligned} \mathbf{x}_k &= \mathbf{Fx}_{k-1} \\ &= \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_{k-1} \\ \dot{x}_{k-1} \end{bmatrix} \end{aligned}$$

We can do that only because  $\dot{x} = v$  is simplest differential equation possible. Almost all other in physical systems result in more complicated differential equation which do not yield to this approach.

State-space methods became popular around the time of the Apollo missions, largely due to the work of Dr. Kalman. The idea is simple. Model a system with a set of  $n^{th}$ -order differential equations. Convert them into an equivalent set of first-order differential equations. Put them into the vector-matrix form used in the previous section:  $\dot{\mathbf{x}} = \mathbf{Ax} + \mathbf{Bu}$ . Once in this form we use of of several techniques to convert these linear differential equations into the recursive equation:

$$\mathbf{x}_k = \mathbf{Fx}_{k-1} + \mathbf{B}_k \mathbf{u}_k$$

Some books call the state transition matrix the fundamental matrix. Many use  $\Phi$  instead of  $\mathbf{F}$ . Sources based heavily on control theory tend to use these forms.

These are called state-space methods because we are expressing the solution of the differential equations in terms of the system state.

### 7.2.1 Forming First Order Equations from Higher Order Equations

Many models of physical systems require second or higher order differential equations with control input  $u$ :

$$a_n \frac{d^n y}{dt^n} + a_{n-1} \frac{d^{n-1} y}{dt^{n-1}} + \cdots + a_2 \frac{d^2 y}{dt^2} + a_1 \frac{dy}{dt} + a_0 = u$$

State-space methods require first-order equations. Any higher order system of equations can be reduced to first-order by defining extra variables for the derivatives and then solving.

Let's do an example. Given the system  $\ddot{x} - 6\dot{x} + 9x = t$  find the equivalent first order equations. I've used the dot notation for the time derivatives for clarity.

The first step is to isolate the highest order term onto one side of the equation.

$$\ddot{x} = 6\dot{x} - 9x + t$$

We define two new variables:

$$x_1(t) = x, \quad x_2(t) = \dot{x}$$

Now we will substitute these into the original equation and solve. The solution yields a set of first-order equations in terms of these new variables. It is conventional to drop the  $(t)$  for notational convenience.

We know that  $\dot{x}_1 = x_2$  and that  $\dot{x}_2 = \ddot{x}$ . Therefore

$$\begin{aligned}\dot{x}_2 &= \ddot{x} \\ &= 6\dot{x} - 9x + t \\ &= 6x_2 - 9x_1 + t\end{aligned}$$

Therefore our first-order system of equations is

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= 6x_2 - 9x_1 + t\end{aligned}$$

If you practice this a bit you will become adept at it. Isolate the highest term, define a new variable and its derivatives, and then substitute.

### 7.2.2 First Order Differential Equations In State-Space Form

Substituting the newly defined variables from the previous section:

$$\frac{dx_1}{dt} = x_2, \quad \frac{dx_2}{dt} = x_3, \quad \dots, \quad \frac{dx_{n-1}}{dt} = x_n$$

into the first order equations yields:

$$\frac{dx_n}{dt} = \frac{1}{a_n} \sum_{i=0}^{n-1} a_i x_{i+1} + \frac{1}{a_n} u$$

Using vector-matrix notation we have:

$$\begin{bmatrix} \frac{dx_1}{dt} \\ \frac{dx_2}{dt} \\ \vdots \\ \frac{dx_n}{dt} \end{bmatrix} = \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \vdots \\ \dot{x}_n \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -\frac{a_0}{a_n} & -\frac{a_1}{a_n} & -\frac{a_2}{a_n} & \cdots & -\frac{a_{n-1}}{a_n} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ \frac{1}{a_n} \end{bmatrix} u$$

which we then write as  $\dot{\mathbf{x}} = \mathbf{Ax} + \mathbf{Bu}$ .

### 7.2.3 Finding the Fundamental Matrix for Time Invariant Systems

We express the system equations in state-space form with

$$\dot{\mathbf{x}} = \mathbf{Ax}$$

where  $\mathbf{A}$  is the system dynamics matrix, and want to find the fundamental matrix  $\mathbf{F}$  that propagates the state  $\mathbf{x}$  over the interval  $\Delta t$  with the equation

$$\mathbf{x}(t_k) = \mathbf{F}(\Delta t)\mathbf{x}(t_{k-1})$$

In other words,  $\mathbf{A}$  is a set of continuous differential equations, and we need  $\mathbf{F}$  to be a set of discrete linear equations that computes the change in  $\mathbf{A}$  over a discrete time step.

It is conventional to drop the  $t_k$  and  $(\Delta t)$  and use the notation

$$\mathbf{x}_k = \mathbf{F}\mathbf{x}_{k-1}$$

$\mathbf{x}_k$  does not mean the  $k^{th}$  value of  $\mathbf{x}$ , but the value of  $\mathbf{x}$  at the  $k^{th}$  value of  $t$ . In this book I go even further and drop the suffixes entirely in favor of an overline to denote the prediction:  $\bar{\mathbf{x}} = \mathbf{Fx}$ .

Broadly speaking there are three common ways to find this matrix for Kalman filters. The technique most often used with Kalman filters is to use a matrix exponential. Linear Time Invariant Theory, also known as LTI System Theory, is a second technique. Finally, there are numerical techniques. You may know of others, but these three are what you will most likely encounter in the Kalman filter literature and praxis.

### 7.2.4 The Matrix Exponential

The solution to the equation  $\frac{dx}{dt} = kx$  can be found by:

$$\frac{dx}{dt} = kx \frac{dx}{x} = k dt \int \frac{1}{x} dx = \int k dt \log x = kt + cx = e^{kt+c}x = e^c e^{kt}x = c_0 e^{kt}$$

Using similar math, the solution to the first-order equation

$$\dot{\mathbf{x}} = \mathbf{Ax}, \quad \mathbf{x}(0) = \mathbf{x}_0$$

where  $\mathbf{A}$  is a constant matrix, is

$$\mathbf{x} = e^{\mathbf{At}}\mathbf{x}_0$$

Substituting  $F = e^{\mathbf{At}}$ , we can write

$$\mathbf{x}_k = \mathbf{F}\mathbf{x}_{k-1}$$

which is the form we are looking for! We have reduced the problem of finding the fundamental matrix to one of finding the value for  $e^{\mathbf{A}t}$ .

$e^{\mathbf{A}t}$  is known as the matrix exponential. It can be computed with this power series:

$$e^{\mathbf{A}t} = \mathbf{I} + \mathbf{A}t + \frac{(\mathbf{A}t)^2}{2!} + \frac{(\mathbf{A}t)^3}{3!} + \dots$$

That series is found by doing a Taylor series expansion of  $e^{\mathbf{A}t}$ , which I will not cover here.

Let's use this to find the solution to Newton's equations. Using  $v$  as an substitution for  $\dot{x}$ , and assuming constant velocity we get the linear matrix-vector form

$$\begin{bmatrix} \dot{x} \\ \dot{v} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ v \end{bmatrix}$$

This is a first order differential equation, so we can set  $\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$  and solve the following equation. I have substituted the interval  $\Delta t$  for  $t$  to emphasize that the fundamental matrix is discrete:

$$\mathbf{F} = e^{\mathbf{A}\Delta t} = \mathbf{I} + \mathbf{A}\Delta t + \frac{(\mathbf{A}\Delta t)^2}{2!} + \frac{(\mathbf{A}\Delta t)^3}{3!} + \dots$$

If you perform the multiplication you will find that  $\mathbf{A}^2 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$ , which means that all higher powers of  $\mathbf{A}$  are also  $\mathbf{0}$ . Thus we get an exact answer without an infinite number of terms:

$$\begin{aligned} \mathbf{F} &= \mathbf{I} + \mathbf{A}\Delta t + \mathbf{0} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \Delta t \\ &= \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix} \end{aligned}$$

We plug this into  $\mathbf{x}_k = \mathbf{F}\mathbf{x}_{k-1}$  to get

$$\mathbf{x}_k = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix} \mathbf{x}_{k-1}$$

You will recognize this as the matrix we derived analytically for the constant velocity Kalman filter in the **Multivariate Kalman Filter** chapter.

### 7.2.5 Time Invariance

If the behavior of the system depends on time we can say that a dynamic system is described by the first-order differential equation

$$g(t) = \dot{x}$$

However, if the system is time invariant the equation is of the form:

$$f(x) = \dot{x}$$

What does time invariant mean? Consider a home stereo. If you input a signal  $x$  into it at time  $t$ , it will output some signal  $f(x)$ . If you instead perform the input at time  $t + \Delta t$  the output signal will be the same  $f(x)$ , shifted in time.

A counter-example is  $x(t) = \sin(t)$ , with the system  $f(x) = t x(t) = t \sin(t)$ . This is not time invariant; the value will be different at different times due to the multiplication by  $t$ . An aircraft is not time invariant. If you make a control input to the aircraft at a later time its behavior will be different because it will have burned fuel and thus lost weight. Lower weight results in different behavior.

We can solve these equations by integrating each side. I demonstrated integrating the time invariants system  $v = \dot{x}$  above. However, integrating the time invariant equation  $\dot{x} = f(x)$  is not so straightforward. Using the separation of variables techniques we divide by  $f(x)$  and move the  $dt$  term to the right so we can integrate each side:

$$\frac{dx}{dt} = f(x) \int_{x_0}^x \frac{1}{f(x)} dx = \int_{t_0}^t dt$$

If we let  $F(x) = \int \frac{1}{f(x)} dx$  we get

$$F(x) - F(x_0) = t - t_0$$

We then solve for  $x$  with

$$F(x) = t - t_0 + F(x_0) = F^{-1}[t - t_0 + F(x_0)]$$

In other words, we need to find the inverse of  $F$ . This is not trivial, and a significant amount of coursework in a STEM education is devoted to finding tricky, analytic solutions to this problem.

However, they are tricks, and many simple forms of  $f(x)$  either have no closed form solution or pose extreme difficulties. Instead, the practicing engineer turns to state-space methods to find approximate solutions.

The advantage of the matrix exponential is that we can use it for any arbitrary set of differential equations which are time invariant. However, we often use this technique even when the equations are not time invariant. As an aircraft flies it burns fuel and loses weight. However, the weight loss over one second is negligible, and so the system is nearly linear over that time step. Our answers will still be reasonably accurate so long as the time step is short.

### Example: Mass-Spring-Damper Model

Suppose we wanted to track the motion of a weight on a spring and connected to a damper, such as an automobile's suspension. The equation for the motion with  $m$  being the mass,  $k$  the spring constant, and  $c$  the damping force, under some input  $u$  is

$$m \frac{d^2x}{dt^2} + c \frac{dx}{dt} + kx = u$$

For notational convenience I will write that as

$$m\ddot{x} + c\dot{x} + kx = u$$

I can turn this into a system of first order equations by setting  $x_1(t) = x(t)$ , and then substituting as follows:

$$\begin{aligned} x_1 &= x \\ x_2 &= \dot{x}_1 \\ \dot{x}_2 &= \dot{x}_1 = \ddot{x} \end{aligned}$$

As is common I dropped the  $(t)$  for notational convenience. This gives the equation

$$m\dot{x}_2 + cx_2 + kx_1 = u$$

Solving for  $\dot{x}_2$  we get a first order equation:

$$\dot{x}_2 = -\frac{c}{m}x_2 - \frac{k}{m}x_1 + \frac{1}{m}u$$

We put this into matrix form:

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -k/m & -c/m \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1/m \end{bmatrix} u$$

Now we use the matrix exponential to find the state transition matrix:

$$\Phi(t) = e^{\mathbf{A}t} = \mathbf{I} + \mathbf{A}t + \frac{(\mathbf{A}t)^2}{2!} + \frac{(\mathbf{A}t)^3}{3!} + \dots$$

The first two terms give us

$$\mathbf{F} = \begin{bmatrix} 1 & t \\ -(k/m)t & 1 - (c/m)t \end{bmatrix}$$

This may or may not give you enough precision. You can easily check this by computing  $\frac{(\mathbf{A}t)^2}{2!}$  for your constants and seeing how much this matrix contributes to the results.

### 7.2.6 Linear Time Invariant Theory

Linear Time Invariant Theory, also known as LTI System Theory, gives us a way to find  $\Phi$  using the inverse Laplace transform. You are either nodding your head now, or completely lost. Don't worry, I will not be using the Laplace transform in this book except in this paragraph, as the computation can be quite difficult. LTI system theory tells us that

$$\Phi(t) = \mathcal{L}^{-1}[(s\mathbf{I} - \mathbf{F})^{-1}]$$

I have no intention of going into this other than to say that the Laplace transform  $\mathcal{L}$  converts a signal into a space  $s$  that excludes time, but finding a solution to the equation above is non-trivial. If you are interested, the Wikipedia article on LTI system theory provides an introduction [2]. I mention LTI because you will find some literature using it to design the Kalman filter matrices for difficult problems.

### 7.2.7 Numerical Solutions

Finally, there are numerous numerical techniques to find  $\mathbf{F}$ . As filters get larger finding analytical solutions becomes very tedious (though packages like SymPy make it easier). C. F. van Loan [3] has developed a technique that finds both  $\Phi$  and  $\mathbf{Q}$  numerically. Given the continuous model

$$\dot{x} = Ax + Gw$$

where  $w$  is the unity white noise, van Loan's method computes both  $\mathbf{F}_k$  and  $\mathbf{Q}_k$ .

I have implemented van Loan's method in `FilterPy`. You may use it as follows:

```
from filterpy.common import van_loan_discretization

A = np.array([[0., 1.], [-1., 0.]])
G = np.array([[0.], [2.]]) # white noise scaling
F, Q = van_loan_discretization(A, G, dt=0.1)
```

In the section [Numeric Integration of Differential Equations](#) I present alternative methods which are very commonly used in Kalman filtering.

## 7.3 Design of the Process Noise Matrix

In general the design of the **Q** matrix is among the most difficult aspects of Kalman filter design. This is due to several factors. First, the math requires a good foundation in signal theory. Second, we are trying to model the noise in something for which we have little information. Consider trying to model the process noise for a thrown baseball. We can model it as a sphere moving through the air, but that leave many unknown factors - the wind, ball rotation and spin decay, the coefficient of drag of a ball with stitches, the effects of wind and air density, and so on. We develop the equations for an exact mathematical solution for a given process model, but since the process model is incomplete the result for **Q** will also be incomplete. This has a lot of ramifications for the behavior of the Kalman filter. If **Q** is too small then the filter will be overconfident in its prediction model and will diverge from the actual solution. If **Q** is too large than the filter will be unduly influenced by the noise in the measurements and perform sub-optimally. In practice we spend a lot of time running simulations and evaluating collected data to try to select an appropriate value for **Q**. But let's start by looking at the math.

Let's assume a kinematic system - some system that can be modeled using Newton's equations of motion. We can make a few different assumptions about this process.

We have been using a process model of

$$\dot{\mathbf{x}} = \mathbf{Ax} + \mathbf{Bu} + \mathbf{w}$$

where **w** is the process noise. Kinematic systems are continuous - their inputs and outputs can vary at any arbitrary point in time. However, our Kalman filters are discrete. We sample the system at regular intervals. Therefore we must find the discrete representation for the noise term in the equation above. This depends on what assumptions we make about the behavior of the noise. We will consider two different models for the noise.

### 7.3.1 Continuous White Noise Model

We model kinematic systems using Newton's equations. We have either used position and velocity, or position, velocity, and acceleration as the models for our systems. There is nothing stopping us from going further - we can model jerk, jounce, snap, and so on. We don't do that normally because adding terms beyond the dynamics of the real system degrades the estimate.

Let's say that we need to model the position, velocity, and acceleration. We can then assume that acceleration is constant for each discrete time step. Of course, there is process noise in the system and so the acceleration is not actually constant. The tracked object will alter the acceleration over time due to external, unmodeled forces. In this section we will assume that the acceleration changes by a continuous time zero-mean white noise  $w(t)$ . In other words, we are assuming that velocity is acceleration changing by small amounts that over time average to 0 (zero-mean).

Since the noise is changing continuously we will need to integrate to get the discrete noise for the discretization interval that we have chosen. We will not prove it here, but the equation for the discretization of the noise is

$$\mathbf{Q} = \int_0^{\Delta t} \mathbf{F}(t) \mathbf{Q}_c \mathbf{F}^\top(t) dt$$

where  $\mathbf{Q}_c$  is the continuous noise. This gives us

$$\Phi = \begin{bmatrix} 1 & \Delta t & \Delta t^2/2 \\ 0 & 1 & \Delta t \\ 0 & 0 & 1 \end{bmatrix}$$

for the fundamental matrix, and

$$\mathbf{Q}_c = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \Phi_s$$

for the continuous process noise matrix, where  $\Phi_s$  is the spectral density of the white noise.

We could carry out these computations ourselves, but I prefer using SymPy to solve the equation.

$$\mathbf{Q}_c = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \Phi_s$$

```
In [2]: import sympy
from sympy import (init_printing, Matrix, MatMul,
                    integrate, symbols)

init_printing(use_latex='mathjax')
dt, phi = symbols('Delta{t} Phi_s')
F_k = Matrix([[1, dt, dt**2/2],
              [0, 1, dt],
              [0, 0, 1]])
Q_c = Matrix([[0, 0, 0],
              [0, 0, 0],
              [0, 0, 1]]) * phi

Q=sympy.integrate(F_k * Q_c * F_k.T, (dt, 0, dt))

# factor phi out of the matrix to make it more readable
Q = Q / phi
sympy.MatMul(Q, phi)
```

Out[2]:

$$\begin{bmatrix} \frac{\Delta t^5}{20} & \frac{\Delta t^4}{8} & \frac{\Delta t^3}{6} \\ \frac{\Delta t^4}{8} & \frac{\Delta t^3}{3} & \frac{\Delta t^2}{2} \\ \frac{\Delta t^3}{6} & \frac{\Delta t^2}{2} & \Delta t \end{bmatrix} \Phi_s$$

For completeness, let us compute the equations for the 0th order and 1st order equations.

```
In [3]: F_k = sympy.Matrix([[1]])
Q_c = sympy.Matrix([[phi]])
```

```
print('0th order discrete process noise')
sympy.integrate(F_k*Q_c*F_k.T, (dt, 0, dt))
```

0th order discrete process noise

Out[3] :

$$[\Delta t \Phi_s]$$

```
In [4]: F_k = sympy.Matrix([[1, dt],
                           [0, 1]])
Q_c = sympy.Matrix([[0, 0],
                    [0, 1]]) * phi

Q = sympy.integrate(F_k * Q_c * F_k.T, (dt, 0, dt))

print('1st order discrete process noise')
# factor phi out of the matrix to make it more readable
Q = Q / phi
sympy.MatMul(Q, phi)
```

1st order discrete process noise

Out[4] :

$$\begin{bmatrix} \frac{\Delta t^3}{3} & \frac{\Delta t^2}{2} \\ \frac{\Delta t^2}{2} & \Delta t \end{bmatrix} \Phi_s$$

### 7.3.2 Piecewise White Noise Model

Another model for the noise assumes that the highest order term (say, acceleration) is constant for the duration of each time period, but differs for each time period, and each of these is uncorrelated between time periods. In other words there is a discontinuous jump in acceleration at each time step. This is subtly different than the model above, where we assumed that the last term had a continuously varying noisy signal applied to it.

We will model this as

$$f(x) = Fx + \Gamma w$$

where  $\Gamma$  is the noise gain of the system, and  $w$  is the constant piecewise acceleration (or velocity, or jerk, etc).

Lets start by looking at a first order system. In this case we have the state transition function

$$\mathbf{F} = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix}$$

In one time period, the change in velocity will be  $w(t)\Delta t$ , and the change in position will be  $w(t)\Delta t^2/2$ , giving us

$$\Gamma = \begin{bmatrix} \frac{1}{2}\Delta t^2 \\ \Delta t \end{bmatrix}$$

The covariance of the process noise is then

$$Q = \mathbb{E}[\Gamma w(t)w(t)\Gamma^T] = \Gamma\sigma_v^2\Gamma^T$$

We can compute that with SymPy as follows

```
In [5]: var=symbols('sigma^2_v')
v = Matrix([[dt**2 / 2], [dt]])

Q = v * var * v.T

# factor variance out of the matrix to make it more readable
Q = Q / var
sympy.MatMul(Q, var)
```

Out[5] :

$$\begin{bmatrix} \frac{\Delta t^4}{4} & \frac{\Delta t^3}{2} \\ \frac{\Delta t^3}{2} & \Delta t^2 \end{bmatrix} \sigma_v^2$$

The second order system proceeds with the same math.

$$\mathbf{F} = \begin{bmatrix} 1 & \Delta t & \Delta t^2/2 \\ 0 & 1 & \Delta t \\ 0 & 0 & 1 \end{bmatrix}$$

Here we will assume that the white noise is a discrete time Wiener process. This gives us

$$\Gamma = \begin{bmatrix} \frac{1}{2}\Delta t^2 \\ \Delta t \\ 1 \end{bmatrix}$$

There is no ‘truth’ to this model, it is just convenient and provides good results. For example, we could assume that the noise is applied to the jerk at the cost of a more complicated equation.

The covariance of the process noise is then

$$Q = \mathbb{E}[\Gamma w(t)w(t)\Gamma^T] = \Gamma\sigma_v^2\Gamma^T$$

We can compute that with SymPy as follows

```
In [6]: var=symbols('sigma^2_v')
v = Matrix([[dt**2 / 2], [dt], [1]])

Q = v * var * v.T

# factor variance out of the matrix to make it more readable
Q = Q / var
sympy.MatMul(Q, var)
```

**Out [6] :**

$$\begin{bmatrix} \frac{\Delta t^4}{4} & \frac{\Delta t^3}{2} & \frac{\Delta t^2}{2} \\ \frac{\Delta t^3}{2} & \Delta t^2 & \Delta t \\ \frac{\Delta t^2}{2} & \Delta t & 1 \end{bmatrix} \sigma_v^2$$

We cannot say that this model is more or less correct than the continuous model - both are approximations to what is happening to the actual object. Only experience and experiments can guide you to the appropriate model. In practice you will usually find that either model provides reasonable results, but typically one will perform better than the other.

The advantage of the second model is that we can model the noise in terms of  $\sigma^2$  which we can describe in terms of the motion and the amount of error we expect. The first model requires us to specify the spectral density, which is not very intuitive, but it handles varying time samples much more easily since the noise is integrated across the time period. However, these are not fixed rules - use whichever model (or a model of your own devising) based on testing how the filter performs and/or your knowledge of the behavior of the physical model.

A good rule of thumb is to set  $\sigma$  somewhere from  $\frac{1}{2}\Delta a$  to  $\Delta a$ , where  $\Delta a$  is the maximum amount that the acceleration will change between sample periods. In practice we pick a number, run simulations on data, and choose a value that works well.

### 7.3.3 Using FilterPy to Compute Q

FilterPy offers several routines to compute the **Q** matrix. The function `Q_continuous_white_noise()` computes **Q** for a given value for  $\Delta t$  and the spectral density.

```
In [7]: from filterpy.common import Q_continuous_white_noise
        from filterpy.common import Q_discrete_white_noise
```

```
Q = Q_continuous_white_noise(dim=2, dt=1, spectral_density=1)
print(Q)
```

```
[[ 0.333  0.5 ]
 [ 0.5     1.    ]]
```

```
In [8]: Q = Q_continuous_white_noise(dim=3, dt=1, spectral_density=1)
print(Q)
```

```
[[ 0.05   0.125  0.167]
 [ 0.125  0.333  0.5   ]
 [ 0.167  0.5     1.    ]]
```

The function `Q_discrete_white_noise()` computes **Q** assuming a piecewise model for the noise.

```
In [9]: Q = Q_discrete_white_noise(2, var=1.)
print(Q)
```

```
[[ 0.25  0.5 ]
 [ 0.5   1.   ]]
```

```
In [10]: Q = Q_discrete_white_noise(3, var=1.)
print(Q)
```

```
[[ 0.25  0.5   0.5 ]
 [ 0.5   1.    1.   ]
 [ 0.5   1.    1.   ]]
```

### 7.3.4 Simplification of $\mathbf{Q}$

Many treatments use a much simpler form for  $\mathbf{Q}$ , setting it to zero except for a noise term in the lower rightmost element. Is this justified? Well, consider the value of  $\mathbf{Q}$  for a small  $\Delta t$

In [11]: `import numpy as np`

```
np.set_printoptions(precision=8)
Q = Q_continuous_white_noise(
    dim=3, dt=0.05, spectral_density=1)
print(Q)
np.set_printoptions(precision=3)

[[ 0.00000002  0.00000078  0.00002083]
 [ 0.00000078  0.00004167  0.00125    ]
 [ 0.00002083  0.00125     0.05      ]]
```

We can see that most of the terms are very small. Recall that the only Kalman filter using this matrix is

$$\mathbf{P} = \mathbf{F}\mathbf{P}\mathbf{F}^T + \mathbf{Q}$$

If the values for  $\mathbf{Q}$  are small relative to  $\mathbf{P}$  than it will be contributing almost nothing to the computation of  $\mathbf{P}$ . Setting  $\mathbf{Q}$  to the zero matrix except for the lower right term

$$\mathbf{Q} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & \sigma^2 \end{bmatrix}$$

while not correct, is often a useful approximation. If you do this you will have to perform quite a few studies to guarantee that your filter works in a variety of situations.

If you do this, ‘lower right term’ means the most rapidly changing term for each variable. If the state is  $x = [x \ \dot{x} \ \ddot{x} \ y \ \dot{y} \ \ddot{y}]^T$  Then  $\mathbf{Q}$  will be 6x6; the elements for both  $\ddot{x}$  and  $\ddot{y}$  will have to be set to non-zero in  $\mathbf{Q}$ .

## 7.4 Numeric Integration of Differential Equations

We’ve been exposed to several numerical techniques to solve linear differential equations. These include state-space methods, the Laplace transform, and van Loan’s method.

These work well for linear ordinary differential equations (ODEs), but do not work well for nonlinear equations. For example, consider trying to predict the position of a rapidly turning car. Cars maneuver by turning the front wheels. This makes them pivot around their rear axle as it moves forward. Therefore the path will be continuously varying and a linear prediction will necessarily produce an incorrect value. If the change in the system is small enough relative to  $\Delta t$  this can often produce adequate results, but that will rarely be the case with the nonlinear Kalman filters we will be studying in subsequent chapters.

For these reasons we need to know how to numerically integrate ODEs. This can be a vast topic that requires several books. If you need to explore this topic in depth [Computational Physics in Python](#) by Dr. Eric Ayars is excellent, and available for free here:

<http://phys.csuchico.edu/ayars/312/Handouts/comp-phys-python.pdf>

However, I will cover a few simple techniques which will work for a majority of the problems you encounter.

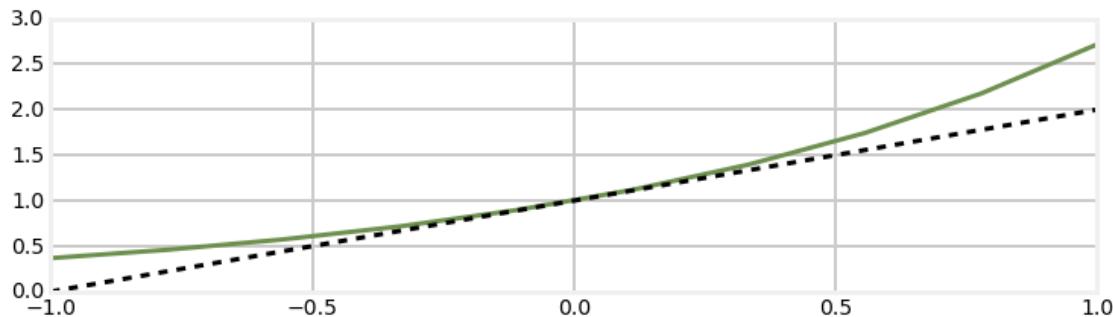
### 7.4.1 Euler's Method

Let's say we have the initial condition problem of

$$y' = y, y(0) = 1$$

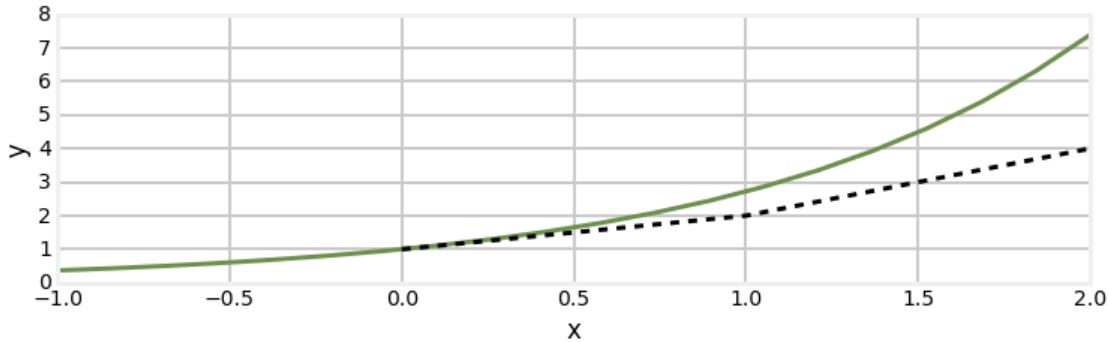
We happen to know the exact answer is  $y = e^t$  because we solved it earlier, but for an arbitrary ODE we will not know the exact solution. In general all we know is the derivative of the equation, which is equal to the slope. We also know the initial value: at  $t = 0$ ,  $y = 1$ . If we know these two pieces of information we can predict the value at  $y(t = 1)$  using the slope at  $t = 0$  and the value of  $y(0)$ . I've plotted this below.

```
In [12]: import matplotlib.pyplot as plt
t = np.linspace(-1, 1, 10)
plt.plot(t, np.exp(t))
t = np.linspace(-1, 1, 2)
plt.plot(t,t+1, ls='--', c='k');
```



You can see that the slope is very close to the curve at  $t = 0.1$ , but far from it at  $t = 1$ . But let's continue with a step size of 1 for a moment. We can see that at  $t = 1$  the estimated value of  $y$  is 2. Now we can compute the value at  $t = 2$  by taking the slope of the curve at  $t = 1$  and adding it to our initial estimate. The slope is computed with  $y' = y$ , so the slope is 2.

```
In [13]: import book_plots
t = np.linspace(-1, 2, 20)
plt.plot(t, np.exp(t))
t = np.linspace(0, 1, 2)
plt.plot([1, 2, 4], ls='--', c='k')
book_plots.set_labels(x='x', y='y');
```



Here we see the next estimate for  $y$  is 4. The errors are getting large quickly, and you might be unimpressed. But 1 is a very large step size. Let's put this algorithm in code, and verify that it works by using a small step size.

```
In [14]: def euler(t, tmax, y, dx, step=1.):
    ys = []
    while t < tmax:
        y = y + step*dx(t, y)
        ys.append(y)
        t += step
    return ys

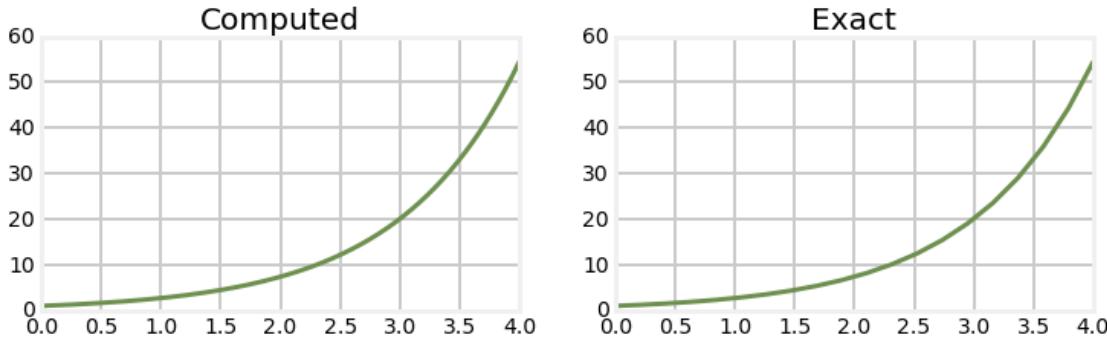
In [15]: def dx(t, y): return y

print(euler(0, 1, 1, dx, step=1.)[-1])
print(euler(0, 2, 1, dx, step=1.)[-1])
```

2.0  
4.0

This looks correct. So now lets plot the result of a much smaller step size.

```
In [16]: ys = euler(0, 4, 1, dx, step=0.00001)
plt.subplot(1,2,1)
plt.title('Computed')
plt.plot(np.linspace(0, 4, len(ys)), ys)
plt.subplot(1,2,2)
t = np.linspace(0, 4, 20)
plt.title('Exact')
plt.plot(t, np.exp(t));
```



```
In [17]: print('exact answer=', np.exp(4))
print('euler answer=', ys[-1])
print('difference =', np.exp(4) - ys[-1])
print('iterations =', len(ys))

exact answer= 54.5981500331
euler answer= 54.59705808834125
difference = 0.00109194480299
iterations = 400000
```

Here we see that the error is reasonably small, but it took a very large number of iterations to get three digits of precision. In practice Euler's method is too slow for most problems, and we use more sophisticated methods.

Before we go on, let's formally derive Euler's method, as it is the basis for the more advanced Runge Kutta methods used in the next section. In fact, Euler's method is the simplest form of Runge Kutta.

Here are the first 3 terms of the Euler expansion of  $y$ . An infinite expansion would give an exact answer, so  $O(h^4)$  denotes the error due to the finite expansion.

$$y(t_0 + h) = y(t_0) + hy'(t_0) + \frac{1}{2!}h^2y''(t_0) + \frac{1}{3!}h^3y'''(t_0) + O(h^4)$$

Here we can see that Euler's method is using the first two terms of the Taylor expansion. Each subsequent term is smaller than the previous terms, so we are assured that the estimate will not be too far off from the correct value.

### 7.4.2 Runge Kutta Methods

Runge Kutta integration is the workhorse of numerical integration. There are a vast number of methods in the literature. In practice, using the Runge Kutta algorithm that I present here will solve most any problem you will face. It offers a very good balance of speed, precision, and stability, and it is the 'go to' numerical integration method unless you have a very good reason to choose something different.

Let's dive in. We start with some differential equation

$$\ddot{y} = \frac{d}{dt}\dot{y}$$

We can substitute the derivative of  $y$  with a function  $f$ , like so

$$\ddot{y} = \frac{d}{dt} f(y, t)$$

Deriving these equations is outside the scope of this book, but the Runge-Kutta RK4 method is defined with these equations.

$$y(t + \Delta t) = y(t) + \frac{1}{6}(k_1 + 2k_2 + 2k_3 + k_4) + O(\Delta t^4)$$

$$\begin{aligned} k_1 &= f(y, t)\Delta t \\ k_2 &= f\left(y + \frac{1}{2}k_1, t + \frac{1}{2}\Delta t\right)\Delta t \\ k_3 &= f\left(y + \frac{1}{2}k_2, t + \frac{1}{2}\Delta t\right)\Delta t \\ k_4 &= f(y + k_3, t + \Delta t)\Delta t \end{aligned}$$

Here is the corresponding code:

```
In [18]: def runge_kutta4(y, x, dx, f):
    """computes 4th order Runge-Kutta for dy/dx.
    y is the initial value for y
    x is the initial value for x
    dx is the difference in x (e.g. the time step)
    f is a callable function (y, x) that you supply
    to compute dy/dx for the specified values.
    """
    k1 = dx * f(y, x)
    k2 = dx * f(y + 0.5*k1, x + 0.5*dx)
    k3 = dx * f(y + 0.5*k2, x + 0.5*dx)
    k4 = dx * f(y + k3, x + dx)

    return y + (k1 + 2*k2 + 2*k3 + k4) / 6.
```

Let's use this for a simple example. Let

$$\dot{y} = t\sqrt{y(t)}$$

with the initial values

$$\begin{aligned} t_0 &= 0 \\ y_0 &= y(t_0) = 1 \end{aligned}$$

```
In [19]: import math
import numpy as np
t = 0.
y = 1.
dt = .1

ys, ts = [], []
```

```

def func(y,t):
    return t*math.sqrt(y)

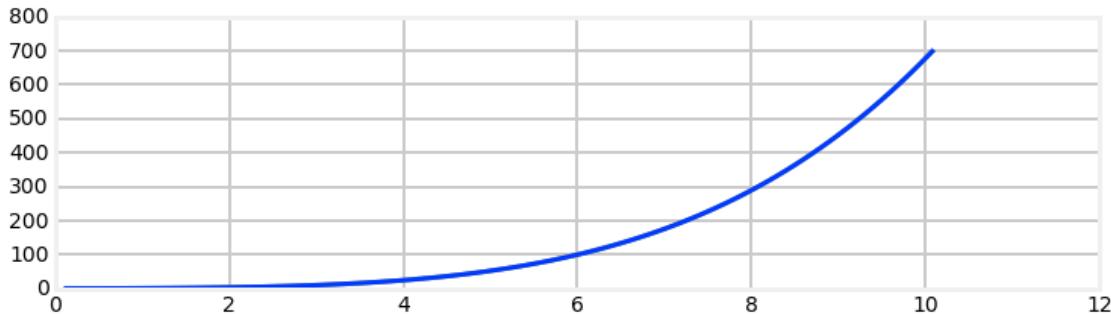
while t <= 10:
    y = runge_kutta4(y, t, dt, func)
    t += dt
    ys.append(y)
    ts.append(t)

exact = [(t**2 + 4)**2 / 16. for t in ts]
plt.plot(ts, ys)
plt.plot(ts, exact)

error = np.array(exact) - np.array(ys)
print("max error {}".format(max(error)))

max error 5.206970035942504e-05

```



## 7.5 Bayesian Filtering

Starting in the Discrete Bayes chapter I used a Bayesian formulation for filtering. Suppose we are tracking an object. We define its state at a specific time as its position, velocity, and so on. For example, we might write the state at time  $t$  as  $\mathbf{x}_t = [x_t \quad \dot{x}_t]^\top$ .

When we take a measurement of the object we are measuring the state. Sensors are noisy, so the measurement is corrupted with noise. Clearly though, the measurement is determined by the state. That is, a change in state may change the measurement, but a change in measurement will not change the state.

In filtering our goal is to compute an optimal estimate for a set of states  $\mathbf{x}_{0:t}$  from time 0 to time  $t$ . If we knew  $\mathbf{x}_{0:t}$  then it would be trivial to compute a set of measurements  $\mathbf{z}_{0:t}$  corresponding to those states. However, we receive a set of measurements  $\mathbf{z}_{0:t}$ , and want to compute the corresponding states  $\mathbf{x}_{0:t}$ . This is called statistical inversion because we are trying to compute the input from the output.

Inversion is a difficult problem because there is typically no unique solution. For a given set of states  $\mathbf{x}_{0:t}$  there is only one possible set of measurements (plus noise), but for a given set of measurements there are many different sets of states that could have led to those measurements.

Recall Bayes Theorem:

$$P(x | z) = \frac{P(z | x)P(x)}{P(z)}$$

where  $P(z | x)$  is the likelihood of the measurement  $z$ ,  $P(x)$  is the prior based on our process model, and  $P(z)$  is a normalization constant.  $P(x | z)$  is the posterior, or the distribution after incorporating the measurement  $z$ , also called the evidence.

This is a statistical inversion as it goes from  $P(z | x)$  to  $P(x | z)$ . The solution to our filtering problem can be expressed as:

$$P(\mathbf{x}_{0:t} | \mathbf{z}_{0:t}) = \frac{P(\mathbf{z}_{0:t} | \mathbf{x}_{0:t})P(\mathbf{x}_{0:t})}{P(\mathbf{z}_{0:t})}$$

That is all well and good until the next measurement  $\mathbf{z}_{t+1}$  comes in, at which point we need to recompute the entire expression for the range  $0 : t + 1$ .

In practice this is intractable because we are trying to compute the posterior distribution  $P(\mathbf{x}_{0:t} | \mathbf{z}_{0:t})$  for the state over the full range of time steps. But do we really care about the probability distribution at the third step (say) when we just received the tenth measurement? Not usually. So we relax our requirements and only compute the distributions for the current time step.

The first simplification is we describe our process (e.g., the motion model for a moving object) as a Markov chain. That is, we say that the current state is solely dependent on the previous state and a transition probability  $P(\mathbf{x}_k | \mathbf{x}_{k-1})$ , which is just the probability of going from the last state to the current one. We write:

$$\mathbf{x}_k \sim P(\mathbf{x}_k | \mathbf{x}_{k-1})$$

The next simplification we make is to define the measurement model as depending on the current state  $\mathbf{x}_k$  with the conditional probability of the measurement given the current state:  $P(\mathbf{z}_t | \mathbf{x}_x)$ . We write:

$$\mathbf{z}_k \sim P(\mathbf{z}_t | \mathbf{x}_x)$$

We have a recurrence now, so we need an initial condition to terminate it. Therefore we say that the initial distribution is the probability of the state  $\mathbf{x}_0$ :

$$\mathbf{x}_0 \sim P(\mathbf{x}_0)$$

These terms are plugged into Bayes equation. If we have the state  $\mathbf{x}_0$  and the first measurement we can estimate  $P(\mathbf{x}_1 | \mathbf{z}_1)$ . The motion model creates the prior  $P(\mathbf{x}_2 | \mathbf{x}_1)$ . We feed this back into Bayes theorem to compute  $P(\mathbf{x}_2 | \mathbf{z}_2)$ . We continue this predictor-corrector algorithm, recursively computing the state and distribution at time  $t$  based solely on the state and distribution at time  $t - 1$  and the measurement at time  $t$ .

The details of the mathematics for this computation varies based on the problem. The **Discrete Bayes** and **Univariate Kalman Filter** chapters gave two different formulations which you should have been able to reason through. The univariate Kalman filter assumes that for a scalar state both the noise and process are linear model are affected by zero-mean, uncorrelated Gaussian noise.

The Multivariate Kalman filter make the same assumption but for states and measurements that are vectors, not scalars. Dr. Kalman was able to prove that if these assumptions hold true then the Kalman filter is optimal. Colloquially this means there is no way to derive more information from the noise. In the remainder of the book I will present filters that relax the constraints on linearity and Gaussian noise.

Before I go on, a few more words about statistical inversion. As Calvetti and Somersalo write in Introduction to Bayesian Scientific Computing, “we adopt the Bayesian point of view: randomness simply means lack of

information.” [4] Our state parametrize physical phenomena that we could in principle measure or compute: velocity, air drag, and so on. We lack enough information to compute or measure their value, so we opt to consider them as random variables. Strictly speaking they are not random, thus this is a subjective position.

They devote a full chapter to this topic. I can spare a paragraph. Bayesian filters are possible because we ascribe statistical properties to unknown parameters. In the case of the Kalman filter we have closed-form solutions to find an optimal estimate. Other filters, such as the discrete Bayes filter or the particle filter which we cover in a later chapter, model the probability in a more ad-hoc, non-optimal manner. The power of our technique comes from treating lack of information as a random variable, describing that random variable as a probability distribution, and then using Bayes Theorem to solve the statistical inference problem.

## 7.6 Converting the Multivariate Equations to the Univariate Case

The multivariate Kalman filter equations do not resemble the equations for the univariate filter. However, if we use one dimensional states and measurements the equations do reduce to the univariate equations. This section will provide you with a strong intuition into what the Kalman filter equations are actually doing. While reading this section is not required to understand the rest of the book, I recommend reading this section carefully as it should make the rest of the material easier to understand.

Here are the multivariate equations for the prediction.

$$\begin{aligned}\bar{\mathbf{x}} &= \mathbf{F}\mathbf{x} + \mathbf{B}\mathbf{u} \\ \bar{\mathbf{P}} &= \mathbf{F}\mathbf{P}\mathbf{F}^T + \mathbf{Q}\end{aligned}$$

For a univariate problem the state  $\mathbf{x}$  only has one variable, so it is a  $1 \times 1$  matrix. Our motion  $\mathbf{u}$  is also a  $1 \times 1$  matrix. Therefore,  $\mathbf{F}$  and  $\mathbf{B}$  must also be  $1 \times 1$  matrices. That means that they are all scalars, and we can write

$$\bar{x} = Fx + Bu$$

Here the variables are not bold, denoting that they are not matrices or vectors.

Our state transition is simple - the next state is the same as this state, so  $F = 1$ . The same holds for the motion transition, so,  $B = 1$ . Thus we have

$$x = x + u$$

which is equivalent to the Gaussian equation from the last chapter

$$\mu = \mu_1 + \mu_2$$

Hopefully the general process is clear, so now I will go a bit faster on the rest. We have

$$\bar{\mathbf{P}} = \mathbf{F}\mathbf{P}\mathbf{F}^T + \mathbf{Q}$$

Again, since our state only has one variable  $\mathbf{P}$  and  $\mathbf{Q}$  must also be  $1 \times 1$  matrix, which we can treat as scalars, yielding

$$\bar{P} = FPF^T + Q$$

We already know  $F = 1$ . The transpose of a scalar is the scalar, so  $F^T = 1$ . This yields

$$\bar{P} = P + Q$$

which is equivalent to the Gaussian equation of

$$\sigma^2 = \sigma_1^2 + \sigma_2^2$$

This proves that the multivariate prediction equations are performing the same math as the univariate equations for the case of the dimension being 1.

These are the equations for the update step:

$$\begin{aligned}\mathbf{K} &= \bar{\mathbf{P}}\mathbf{H}^\top(\mathbf{H}\bar{\mathbf{P}}\mathbf{H}^\top + \mathbf{R})^{-1} \\ \mathbf{y} &= \mathbf{z} - \mathbf{H}\bar{x} \\ \mathbf{x} &= \bar{x} + \mathbf{K}\mathbf{y} \\ \mathbf{P} &= (\mathbf{I} - \mathbf{K}\mathbf{H})\bar{\mathbf{P}}\end{aligned}$$

As above, all of the matrices become scalars.  $H$  defines how we convert from a position to a measurement. Both are positions, so there is no conversion, and thus  $H = 1$ . Let's substitute in our known values and convert to scalar in one step. The inverse of a  $1 \times 1$  matrix is the reciprocal of the value so we will convert the matrix inversion to division.

$$\begin{aligned}K &= \frac{\bar{P}}{\bar{P} + R} \\ y &= z - \bar{x} \\ x &= \bar{x} + Ky \\ P &= (1 - K)\bar{P}\end{aligned}$$

Before we continue with the proof, I want you to look at those equations to recognize what a simple concept these equations implement. The residual  $y$  is nothing more than the measurement minus the prediction. The gain  $K$  is scaled based on how certain we are about the last prediction vs how certain we are about the measurement. We choose a new state  $x$  based on the old value of  $x$  plus the scaled value of the residual. Finally, we update the uncertainty based on how certain we are about the measurement. Algorithmically this should sound exactly like what we did in the last chapter.

Let's finish off the algebra to prove this. Recall that the univariate equations for the update step are:

$$\begin{aligned}\mu &= \frac{\sigma_1^2\mu_2 + \sigma_2^2\mu_1}{\sigma_1^2 + \sigma_2^2}, \\ \sigma^2 &= \frac{1}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}}\end{aligned}$$

Here we will say that  $\mu_1$  is the state  $x$ , and  $\mu_2$  is the measurement  $z$ . Thus it follows that that  $\sigma_1^2$  is the state uncertainty  $P$ , and  $\sigma_2^2$  is the measurement noise  $R$ . Let's substitute those in.

$$\mu = \frac{Px + Rx}{P + R} \sigma^2 = \frac{1}{\frac{1}{P} + \frac{1}{R}}$$

I will handle  $\mu$  first. The corresponding equation in the multivariate case is

$$\begin{aligned}
x &= x + Ky \\
&= x + \frac{P}{P+R}(z - x) \\
&= \frac{P+R}{P+R}x + \frac{Pz - Px}{P+R} \\
&= \frac{Px + Rx + Pz - Px}{P+R} \\
&= \frac{Pz + Rx}{P+R}
\end{aligned}$$

Now let's look at  $\sigma^2$ . The corresponding equation in the multivariate case is

$$\begin{aligned}
P &= (1 - K)P \\
&= (1 - \frac{P}{P+R})P \\
&= (\frac{P+R}{P+R} - \frac{P}{P+R})P \\
&= (\frac{P+R-P}{P+R})P \\
&= \frac{RP}{P+R} \\
&= \frac{1}{\frac{P+R}{RP}} \\
&= \frac{1}{\frac{R}{RP} + \frac{P}{RP}} \\
&= \frac{1}{\frac{1}{P} + \frac{1}{R}} \quad \blacksquare
\end{aligned}$$

We have proven that the multivariate equations are equivalent to the univariate equations when we only have one state variable. I'll close this section by recognizing one quibble - I hand waved my assertion that  $H = 1$  and  $F = 1$ . In general we know this is not true. For example, a digital thermometer may provide measurement in volts, and we need to convert that to temperature, and we use  $H$  to do that conversion. I left that issue out to keep the explanation as simple and streamlined as possible. It is very straightforward to add that generalization to the equations above, redo the algebra, and still have the same results.

## 7.7 Converting Kalman Filter to a g-h Filter

I've stated that the Kalman filter is a form of the g-h filter. It just takes some algebra to prove it. It's more straightforward to do with the one dimensional case, so I will do that. Recall

$$\mu_x = \frac{\sigma_1^2 \mu_2 + \sigma_2^2 \mu_1}{\sigma_1^2 + \sigma_2^2}$$

which I will make more friendly for our eyes as:

$$\mu_x = \frac{ya + xb}{a + b}$$

We can easily put this into the g-h form with the following algebra

$$\begin{aligned}
\mu_x &= (x - x) + \frac{ya + xb}{a + b} \\
\mu_x &= x - \frac{a + b}{a + b}x + \frac{ya + xb}{a + b} \\
\mu_x &= x + \frac{-x(a + b) + xb + ya}{a + b} \\
\mu_x &= x + \frac{-xa + ya}{a + b} \\
\mu_x &= x + \frac{a}{a + b}(y - x)
\end{aligned}$$

We are almost done, but recall that the variance of estimate is given by

$$\sigma_x^2 = \frac{1}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}} = \frac{1}{\frac{1}{a} + \frac{1}{b}}$$

We can incorporate that term into our equation above by observing that

$$\begin{aligned}
\frac{a}{a + b} &= \frac{a/a}{(a + b)/a} = \frac{1}{(a + b)/a} \\
&= \frac{1}{1 + \frac{b}{a}} = \frac{1}{\frac{b}{b} + \frac{b}{a}} \\
&= \frac{1}{b} \frac{1}{\frac{1}{b} + \frac{1}{a}} \\
&= \frac{\sigma_{x'}^2}{b}
\end{aligned}$$

We can tie all of this together with

$$\begin{aligned}
\mu_x &= x + \frac{a}{a + b}(y - x) \\
&= x + \frac{\sigma_{x'}^2}{b}(y - x) \\
&= x + g_n(y - x)
\end{aligned}$$

where

$$g_n = \frac{\sigma_x^2}{\sigma_y^2}$$

The end result is multiplying the residual of the two measurements by a constant and adding to our previous value, which is the  $g$  equation for the g-h filter.  $g$  is the variance of the new estimate divided by the variance of the measurement. Of course in this case  $g$  is not a constant as it varies with each time step as the variance changes. We can also derive the formula for  $h$  in the same way. It is not a particularly illuminating derivation and I will skip it. The end result is

$$h_n = \frac{COV(x, \dot{x})}{\sigma_y^2}$$

The takeaway point is that  $g$  and  $h$  are specified fully by the variance and covariances of the measurement and predictions at time  $n$ . In other words, we are picking a point between the measurement and prediction by a scale factor determined by the quality of each of those two inputs.

## 7.8 References

- [1] Matrix Exponential [http://en.wikipedia.org/wiki/Matrix\\_exponential](http://en.wikipedia.org/wiki/Matrix_exponential)
- [2] LTI System Theory [http://en.wikipedia.org/wiki/LTI\\_system\\_theory](http://en.wikipedia.org/wiki/LTI_system_theory)
- [3] C.F. van Loan, “Computing Integrals Involving the Matrix Exponential,” IEEE Transactions Automatic Control, June 1978.
- [4] Calvetti, D and Somersalo E, “Introduction to Bayesian Scientific Computing: Ten Lectures on Subjective Computing,” Springer, 2007.



# Chapter 8

## Designing Kalman Filters

### 8.1 Introduction

In the last chapter we worked with ‘textbook’ problems. These are problems that are easy to state, program in a few lines of code, and teach. Real world problems are rarely this simple. In this chapter, and the remainder of the book, we will work with more realistic examples.

We will begin by tracking a robot in a 2D space, such as a field or warehouse. We will start with a simple noisy sensor that outputs noisy  $(x, y)$  coordinates which we will need to filter to generate a 2D track. Once we have mastered this concept, we will extend the problem significantly with more sensors and then adding control inputs.

We will then move to a nonlinear problem. The world is nonlinear, but the Kalman filter is linear. Sometimes you can get away with using it for mildly nonlinear problems, sometimes you can’t. I’ll show you examples of both. This will set the stage for the remainder of the book, where we learn techniques for nonlinear problems.

### 8.2 Tracking a Robot

This first attempt at tracking a robot will closely resemble the 1D dog tracking problem of previous chapters. This will allow us to ‘get our feet wet’ with Kalman filtering. Instead of a sensor that outputs position in a hallway, we now have a sensor that supplies a noisy measurement of position in a 2D space. At each time  $t$  it will provide an  $(x, y)$  coordinate pair of the noisy measurement of the sensor’s position in the field.

Implementation of code to interact with real sensors is beyond the scope of this book, so as before we will program simple simulations of the sensors. We will develop several of these sensors as we go, each with more complications, so as I program them I will just append a number to the function name.

So let’s start with a very simple sensor, one that simulates tracking an object traveling in a straight line. It is initialized with the initial position, velocity, and noise standard deviation. Each call to `read()` updates the position by one time step and returns the new measurement.

```
In [2]: from numpy.random import randn
import copy
class PosSensor1(object):
    def __init__(self, pos=(0, 0), vel=(0, 0), noise_std=1.):
        self.vel = vel
        self.noise_std = noise_std
        self.pos = [pos[0], pos[1]]
```

```

def read(self):
    self.pos[0] += self.vel[0]
    self.pos[1] += self.vel[1]

    return [self.pos[0] + randn() * self.noise_std,
            self.pos[1] + randn() * self.noise_std]

```

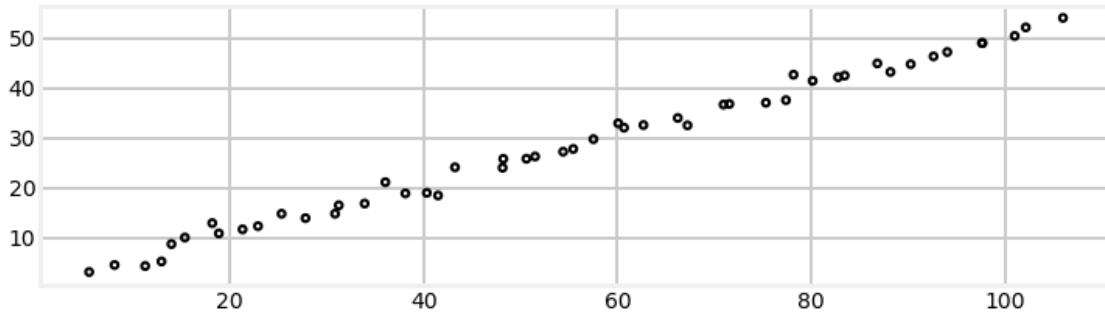
A quick test to verify that it works as we expect.

```

In [3]: import matplotlib.pyplot as plt
        import numpy as np
        import book_plots as bp

        pos, vel = (4, 3), (2, 1)
        sensor = PosSensor1(pos, vel, noise_std=1)
        ps = np.array([sensor.read() for _ in range(50)])
        bp.plot_measurements(ps[:, 0], ps[:, 1]);

```



That looks correct. The slope is  $1/2$ , as we would expect with a velocity of  $(2, 1)$ , and the data seems to start at near  $(6, 4)$ . It doesn't look realistic. This is still a 'textbook' representation. As we continue we will add complications that adds real world behavior.

### 8.2.1 Step 1: Choose the State Variables

As always, the first step is to choose our state variables. We are tracking in two dimensions and have a sensor that gives us a reading in each of those two dimensions, so we know that we have the two observed variables  $x$  and  $y$ . If we created our Kalman filter using only those two variables the performance would not be very good because we would be ignoring the information velocity can provide to us. We will want to incorporate velocity into our equations as well. I will represent this as

$$\mathbf{x} = [x \ \dot{x} \ y \ \dot{y}]^T$$

There is nothing special about this organization. I could have used  $[x \ y \ \dot{x} \ \dot{y}]^T$  or something less logical. I just need to be consistent in the rest of the matrices. I like keeping positions and velocities next to each other because it keeps the covariances between positions and velocities in the same sub block of the covariance matrix. In my formulation  $P[1,0]$  contains the covariance of  $x$  and  $\dot{x}$ . In the alternative formulation that covariance is at  $P[2, 0]$ . This gets worse as the number of dimension increases (e.g. 3D space, accelerations).

Let's pause and address how you identify the hidden variables. This example is somewhat obvious because we've already worked through the 1D case, but other problems won't be obvious. There is no easy answer to this question. The first thing to ask yourself is what is the interpretation of the first and second derivatives of the data from the sensors. We do that because obtaining the first and second derivatives is mathematically trivial if you are reading from the sensors using a fixed time step. The first derivative is just the difference between two successive readings. In our tracking case the first derivative has an obvious physical interpretation: the difference between two successive positions is velocity.

Beyond this you can start looking at how you might combine the data from two or more different sensors to produce more information. This opens up the field of sensor fusion, and we will be covering examples of this in later sections. For now, recognize that choosing the appropriate state variables is paramount to getting the best possible performance from your filter. Once you have chosen hidden variables, you must run many tests to ensure that you are generating real results for them. The Kalman filter runs whatever model you give it; if your model cannot generate good information for the hidden variables the Kalman filter output will be nonsense.

### 8.2.2 Step 2: Design State Transition Function

Our next step is to design the state transition function. Recall that the state transition function is implemented as a matrix  $\mathbf{F}$  that we multiply with the previous state of our system to get the next state, like so.

$$\bar{\mathbf{x}} = \mathbf{F}\mathbf{x}$$

I will not belabor this as it is very similar to the 1-D case we did in the previous chapter. The state equations are

$$\begin{aligned} x &= 1x + \Delta t \dot{x} + 0y + 0\dot{y} \\ v_x &= 0x + 1\dot{x} + 0y + 0\dot{y} \\ y &= 0x + 0\dot{x} + 1y + \Delta t \dot{y} \\ v_y &= 0x + 0\dot{x} + 0y + 1\dot{y} \end{aligned}$$

Laying it out that way shows us both the values and row-column organization required for  $\mathbf{F}$ . In linear algebra, we would write this as:

$$\begin{bmatrix} x \\ \dot{x} \\ y \\ \dot{y} \end{bmatrix} = \begin{bmatrix} 1 & \Delta t & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & \Delta t \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \\ y \\ \dot{y} \end{bmatrix}$$

So, let's do this in Python. It is very simple; the only thing new here is setting `dim_z` to 2. We will see why it is set to 2 in step 4.

```
In [4]: from filterpy.kalman import KalmanFilter
        tracker = KalmanFilter(dim_x=4, dim_z=2)
        dt = 1.    # time step 1 second

        tracker.F = np.array([[1, dt, 0, 0],
                             [0, 1, 0, 0],
                             [0, 0, 1, dt],
                             [0, 0, 0, 1]])
```

### 8.2.3 Step 3: Design the Process Noise Matrix

We will use FilterPy's functions to compute the **Q** matrix for us. For simplicity I will assume the noise is a discrete time Wiener process - that it is constant for each time period. This assumption allows me to use a variance to specify how much I think the model changes between steps. Revisit the Kalman Filter Math chapter if this is not clear.

```
In [5]: from scipy.linalg import block_diag
        from filterpy.common import Q_discrete_white_noise

        q = Q_discrete_white_noise(dim=2, dt=dt, var=0.05)
        tracker.Q = block_diag(q, q)
        print(tracker.Q)

[[ 0.013  0.025  0.     0.     ]
 [ 0.025  0.05   0.     0.     ]
 [ 0.     0.     0.013  0.025]
 [ 0.     0.     0.025  0.05 ]]
```

Here I assume the noise in  $x$  and  $y$  are independent, so the covariances between any  $x$  and  $y$  variable should be zero. This allows me to compute **Q** for one dimension, and then use `block.diag` to copy it for the  $x$  and  $y$  axis. The printed result shows how much noise we add during each evolution (prediction).

### 8.2.4 Step 4: Design the Control Function

We haven't yet added controls to our robot, so there is nothing to be done for this step. The `KalmanFilter` class initializes **B** to zero under the assumption that there is no control input, so there is no code to write. If you like, you can be explicit and set `tracker.B` to 0, but as you can see it already has that value.

```
In [6]: print(tracker.B)
```

```
0
```

### 8.2.5 Step 5: Design the Measurement Function

The measurement function **H** defines how we go from the state variables to the measurements using the equation  $\mathbf{z} = \mathbf{Hx}$ . In this case we have measurements for  $(x,y)$ , so we will design **z** as  $[x \ y]^T$  which is dimension 2x1. Our state variable is size 4x1. We can deduce the required size for **H** by recalling that multiplying a matrix of size MxN by NxP yields a matrix of size MxP. Thus,

$$(2 \times 1) = (a \times b)(4 \times 1) = (2 \times 4)(4 \times 1)$$

So, **H** is 2x4.

Filling in the values for **H** is easy because the measurement is the position of the robot, which is the  $x$  and  $y$  variables of the state **x**. Let's make this slightly more interesting by deciding we want to change units. The measurements are returned in feet, and that we desire to work in meters. **H** changes from state to measurement, so the conversion is  $\text{feet} = \text{meters}/0.3048$ . This yields

$$\mathbf{H} = \begin{bmatrix} \frac{1}{0.3048} & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{0.3048} & 0 \end{bmatrix}$$

which corresponds to these linear equations

$$\begin{aligned} z_x &= \left(\frac{x}{0.3048}\right) + (0 * v_x) + (0 * y) + (0 * v_y) = \frac{x}{0.3048} \\ z_y &= (0 * x) + (0 * v_x) + \left(\frac{y}{0.3048}\right) + (0 * v_y) = \frac{y}{0.3048} \end{aligned}$$

This is a simple problem, and we could have found the equations directly without going through the dimensional analysis that I did above. But it is useful to remember that the equations of the Kalman filter imply a specific dimensionality for all of the matrices, and when I start to get lost as to how to design something it is useful to look at the matrix dimensions.

Here is my implementation:

```
In [7]: tracker.H = np.array([[1/0.3048, 0, 0, 0], [0, 0, 1/0.3048, 0]])
```

### 8.2.6 Step 6: Design the Measurement Noise Matrix

In this step we model the noise in our sensor. For now we will make the simple assumption that the  $x$  and  $y$  variables are independent white Gaussian processes. That is, the noise in  $x$  is not in any way dependent on the noise in  $y$ , and the noise is normally distributed about the mean 0. For now let's set the variance for  $x$  and  $y$  to be 5 meters<sup>2</sup>. They are independent, so there is no covariance, and our off diagonals will be 0. This gives us:

$$\mathbf{R} = \begin{bmatrix} \sigma_x^2 & \sigma_y \sigma_x \\ \sigma_x \sigma_y & \sigma_y^2 \end{bmatrix} = \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$$

It is a  $2 \times 2$  matrix because we have 2 sensor inputs, and covariance matrices are always of size  $n \times n$  for  $n$  variables. In Python we write:

```
In [8]: tracker.R = np.array([[5, 0], [0, 5]])
print(tracker.R)

[[5 0]
 [0 5]]
```

### 8.2.7 Initial Conditions

For our simple problem we will set the initial position at (0,0) with a velocity of (0,0). Since that is a pure guess, we will set the covariance matrix  $\mathbf{P}$  to a large value.

$$\mathbf{x} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \mathbf{P} = \begin{bmatrix} 500 & 0 & 0 & 0 \\ 0 & 500 & 0 & 0 \\ 0 & 0 & 500 & 0 \\ 0 & 0 & 0 & 500 \end{bmatrix}$$

The Python implementation is

```
In [9]: tracker.x = np.array([[0, 0, 0, 0]]).T
tracker.P = np.eye(4) * 500.
```

### 8.2.8 Implement the Filter

Design is complete, now we just have to write the code to run the filter and output the data in the format of our choice. We will run the code for 30 iterations.

```
In [10]: def tracker1():
    tracker = KalmanFilter(dim_x=4, dim_z=2)
    dt = 1.0    # time step

    tracker.F = np.array([[1, dt, 0, 0],
                          [0, 1, 0, 0],
                          [0, 0, 1, dt],
                          [0, 0, 0, 1]])
    tracker.u = 0.
    tracker.H = np.array([[1/0.3048, 0, 0, 0],
                          [0, 0, 1/0.3048, 0]])

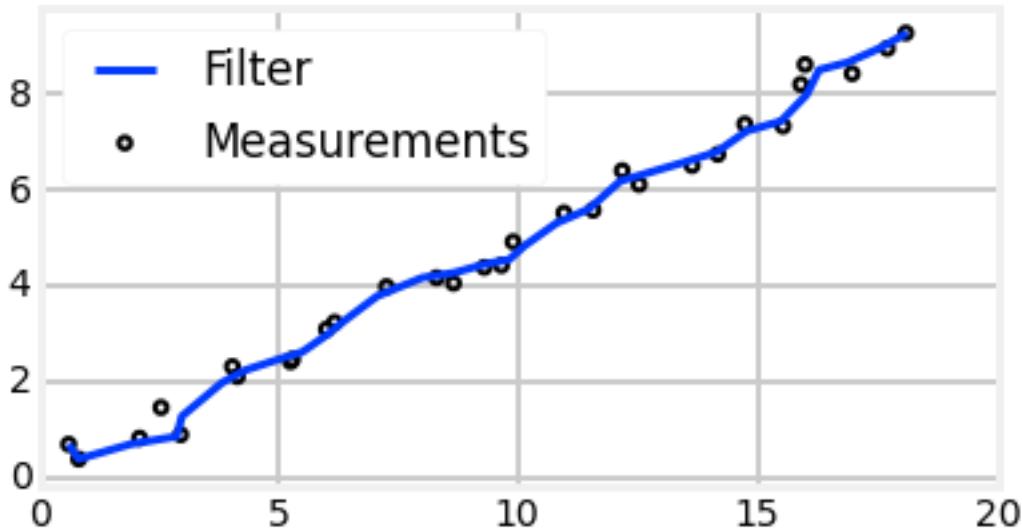
    tracker.R = np.eye(2) * 5
    q = Q_discrete_white_noise(dim=2, dt=dt, var=0.05)
    tracker.Q = block_diag(q, q)
    tracker.x = np.array([[0, 0, 0, 0]]).T
    tracker.P = np.eye(4) * 500.
    return tracker

# simulate robot movement
N = 30
sensor = PosSensor1 ((0, 0), (2, 1), 1.)

zs = np.array([np.array([sensor.read()]).T for _ in range(N)])

# run filter
robot_tracker = tracker1()
mu, cov, _, _ = robot_tracker.batch_filter(zs)

# plot results
zs *= .3048 # convert to meters
bp.plot_filter(mu[:, 0], mu[:, 2])
bp.plot_measurements(zs[:, 0], zs[:, 1])
plt.legend(loc=2)
plt.gca().set_aspect('equal')
plt.xlim((0, 20));
```



I encourage you to play with this, setting  $\mathbf{Q}$  and  $\mathbf{R}$  to various values. However, we did a fair amount of that sort of thing in the last chapters, and we have a lot of material to cover, so I will move on to more complicated cases where we will also have a chance to experience changing these values.

Now I will run the same Kalman filter with the same settings, but also plot the covariance ellipse for  $x$  and  $y$ . First, the code without explanation, so we can see the output. I print the last covariance to see what it looks like. But before you scroll down to look at the results, what do you think it will look like? You have enough information to figure this out but this is still new to you, so don't be discouraged if you get it wrong.

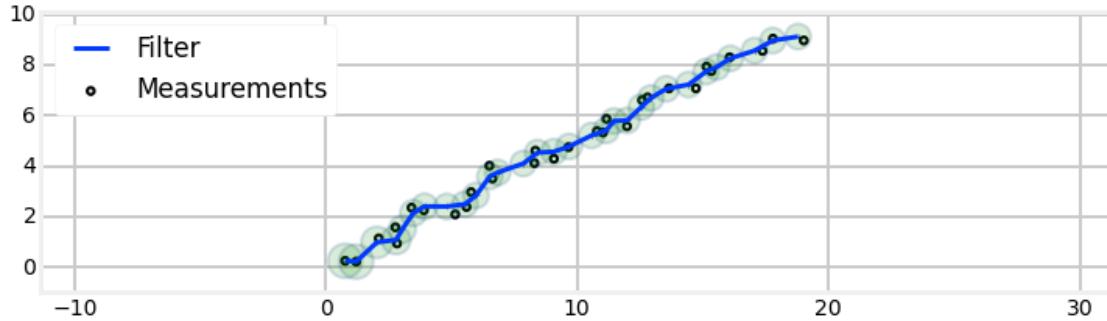
```
In [11]: from filterpy.stats import plot_covariance_ellipse

# simulate robot movement
N = 30
sensor = PosSensor1([0, 0], (2, 1), 1.)
zs = np.array([np.array([sensor.read()]).T for _ in range(N)])

# run filter
robot_tracker = tracker1()
mu, cov, _, _ = robot_tracker.batch_filter(zs)

for x, P in zip(mu, cov):
    # covariance of x and y
    cov = np.array([[P[0, 0], P[2, 0]],
                   [P[0, 2], P[2, 2]]])
    mean = (x[0, 0], x[2, 0])
    plot_covariance_ellipse(mean, cov=cov, fc='g', alpha=0.15)

# plot results
zs *= .3048 # convert to meters
bp.plot_filter(mu[:, 0], mu[:, 2])
bp.plot_measurements(zs[:, 0], zs[:, 1])
plt.legend(loc=2)
plt.gca().set_aspect('equal')
plt.xlim(0, 20);
```



Did you correctly predict what the covariance matrix and plots would look like? Perhaps you were expecting a tilted ellipse, as in the last chapters. If so, recall that in those chapters we were not plotting  $x$  against  $y$ , but  $x$  against  $\dot{x}$ .  $x$  is correlated to  $\dot{x}$ , but  $x$  is not correlated or dependent on  $y$ . Therefore our ellipses are not tilted. Furthermore, the noise for both  $x$  and  $y$  are modeled to have the same value, 5, in  $\mathbf{R}$ . If we were to set  $\mathbf{R}$  to, for example,

$$\mathbf{R} = \begin{bmatrix} 10 & 0 \\ 0 & 5 \end{bmatrix}$$

we would be telling the Kalman filter that there is more noise in  $x$  than  $y$ , and our ellipses would be longer than they are tall.

The final  $\mathbf{P}$  tells us everything we need to know about the correlation between the state variables. If we look at the diagonal alone we see the variance for each variable. In other words  $\mathbf{P}_{0,0}$  is the variance for  $x$ ,  $\mathbf{P}_{1,1}$  is the variance for  $\dot{x}$ ,  $\mathbf{P}_{2,2}$  is the variance for  $y$ , and  $\mathbf{P}_{3,3}$  is the variance for  $\dot{y}$ . We can extract the diagonal of a matrix using `numpy.diag()`.

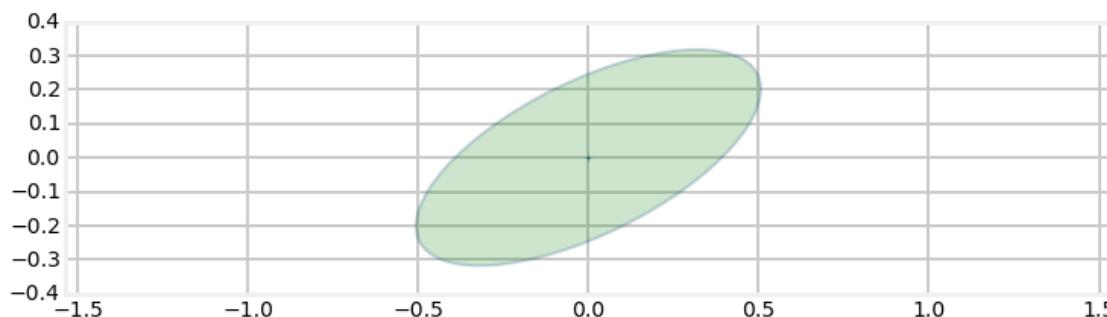
```
In [12]: print(np.diag(robot_tracker.P))

[ 0.257  0.101  0.257  0.101]
```

The covariance matrix contains four  $2\times 2$  matrices that you should be able to easily pick out. This is due to the correlation of  $x$  to  $\dot{x}$ , and of  $y$  to  $\dot{y}$ . The upper left hand side shows the covariance of  $x$  to  $\dot{x}$ .

```
In [13]: c = robot_tracker.P[0:2, 0:2]
print(c)
plot_covariance_ellipse((0, 0), cov=c, fc='g', alpha=0.2)

[[ 0.257  0.102]
 [ 0.102  0.101]]
```



The covariance contains the data for  $x$  and  $\dot{x}$  in the upper left because of how it is organized. Recall that entries  $\mathbf{P}_{i,j}$  and  $\mathbf{P}_{j,i}$  contain  $\sigma_i\sigma_j$ .

Finally, let's look at the lower left side of  $\mathbf{P}$ , which is all 0s. Why 0s? Consider  $\mathbf{P}_{3,0}$ . That stores the term  $\sigma_3\sigma_0$ , which is the covariance between  $\dot{y}$  and  $x$ . These are independent, so the term will be 0. The rest of the terms are for similarly independent variables.

In [14]: `robot_tracker.P[2:4, 0:2]`

Out[14]: `array([[ 0., 0.],  
 [ 0., 0.]])`

## 8.3 The Effect of Filter Order

We have only studied tracking position and velocity. It has worked well, but only because I have been selecting problems for which this is an appropriate choice. You now have enough experience with the Kalman filter to consider this in more general terms.

What do I mean by order? In the context of these system models it is the number of derivatives required to accurately model a system. Consider a system that does not change, such as the height of a building. There is no change, so there is no need for a derivative, and the order of the system is zero. We could express this in an equation as  $x = 312.5$ .

A first order system has a first derivative. For example, change of position is velocity, and we can write this as

$$v = \frac{dx}{dt}$$

which we integrate into the Newtonian equation

$$x = vt + x_0.$$

This is also called a constant velocity model, because of the assumption of a constant velocity.

A second order has a second derivative. The second derivative of position is acceleration, with the equation

$$a = \frac{d^2x}{dt^2}$$

which we integrate into

$$x = \frac{1}{2}at^2 + v_0t + x_0.$$

This is also known as a constant acceleration model.

Another, equivalent way of looking at this is to consider the order of the polynomial. The constant acceleration model has a second derivative, so it is second order. Likewise, the polynomial  $x = \frac{1}{2}at^2 + v_0t + x_0$  is second order.

When we design the state variables and process model we must choose the order of the system we want to model. Let's say we are tracking something with a constant velocity. No real world process is perfect, and so there will be slight variations in the velocity over short time period. You might reason that the best

approach is to use a second order filter, allowing the acceleration term to deal with the slight variations in velocity.

In practice that doesn't work well. To thoroughly understand this issue lets see the effects of using a process model that does not match the order of the system being filtered.

First we need a system to filter. I'll write a class to simulate an object with constant velocity. Essentially no physical system has a truly constant velocity, so on each update we alter the velocity by a small amount. I also write a sensor to simulate Gaussian noise in a sensor. The code is below, and a plot an example run to verify that it is working correctly.

```
In [15]: class ConstantVelocityObject(object):
    def __init__(self, x0=0, vel=1., noise_scale=0.06):
        self.x = x0
        self.vel = vel
        self.noise_scale = noise_scale

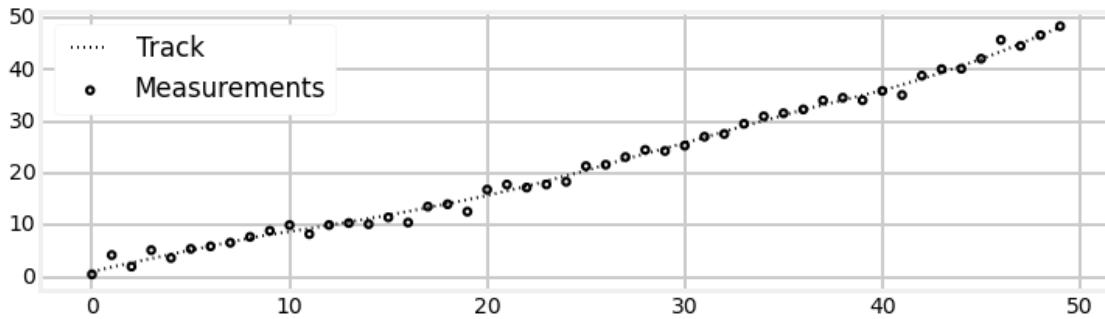
    def update(self):
        self.vel += randn() * self.noise_scale
        self.x += self.vel
        return (self.x, self.vel)

    def sense(x, noise_scale=1.):
        return x[0] + randn()*noise_scale

np.random.seed(124)
obj = ConstantVelocityObject()

xs, zs = [], []
for i in range(50):
    x = obj.update()
    z = sense(x)
    xs.append(x)
    zs.append(z)

xs = np.asarray(xs)
bp.plot_track(xs[:, 0])
bp.plot_measurements(range(len(zs)), zs)
plt.legend(loc='best');
```



I am satisfied with this plot. The track is not perfectly straight due to the noise that we added to the system

- this could be the track of a person walking down the street, or perhaps of an aircraft being buffeted by variable winds. There is no intentional acceleration here, so we call it a constant velocity system. Again, you may be asking yourself that since there is in fact a tiny bit of acceleration going on why would we not use a second order Kalman filter to account for those changes? Let's find out.

How does one design a zero order, first order, or second order Kalman filter? We have been doing it all along, but just not using those terms. It might be slightly tedious, but I will elaborate fully on each - if the concept is clear to you feel free to skim a bit.

### 8.3.1 Zero Order Kalman Filter

A zero order Kalman filter is just a filter that tracks with no derivatives. We are tracking position, so that means we only have a state variable for position (no velocity or acceleration), and the state transition function also only accounts for position. Using the matrix formulation we would say that the state variable is

$$\mathbf{x} = [x]$$

The state transition function is very simple. There is no change in position, so we need to model  $x = x$ ; in other words,  $\underline{x}$  at time  $t+1$  is the same as it was at time  $t$ . In matrix form, our state transition function is

$$\mathbf{F} = [1]$$

The measurement function is very easy. Recall that we need to define how to convert the state variable  $\mathbf{x}$  into a measurement. We will assume that our measurements are positions. The state variable only contains a position, so we get

$$\mathbf{H} = [1]$$

Let's write a function that constructs and returns a zero order Kalman filter.

```
In [16]: def ZeroOrderKF(R, Q, P=20):
    """ Create zero order Kalman filter.
    Specify R and Q as floats."""
    kf = KalmanFilter(dim_x=1, dim_z=1)
    kf.x = np.array([0.])
    kf.R *= R
    kf.Q *= Q
    kf.P *= P
    kf.F = np.eye(1)
    kf.H = np.eye(1)
    return kf
```

### 8.3.2 First Order Kalman Filter

A first order Kalman filter tracks a first order system, such as position and velocity. We already did this for the dog tracking problem above, so this should be very clear. But let's do it again.

A first order system has position and velocity, so the state variable needs both of these. The matrix formulation could be

$$\mathbf{x} = \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

So now we have to design our state transition. The Newtonian equations for a time step are:

$$\begin{aligned}x_t &= x_{t-1} + v\Delta t \\v_t &= v_{t-1}\end{aligned}$$

Recall that we need to convert this into the linear equation

$$\begin{bmatrix} x \\ \dot{x} \end{bmatrix} = \mathbf{F} \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

Setting

$$\mathbf{F} = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix}$$

gives us the equations above.

Finally, we design the measurement function. The measurement function needs to implement

$$\mathbf{z} = \mathbf{H}\mathbf{x}$$

Our sensor still only reads position, so it should take the position from the state, and 0 out the velocity, like so:

$$\mathbf{H} = [1 \ 0]$$

This function constructs and returns a first order Kalman filter.

```
In [17]: def FirstOrderKF(R, Q, dt):
    """ Create first order Kalman filter.
    Specify R and Q as floats."""
    kf = KalmanFilter(dim_x=2, dim_z=1)
    kf.x = np.zeros(2)
    kf.P *= np.array([[100, 0], [0, 1]])
    kf.R *= R
    kf.Q = Q_discrete_white_noise(2, dt, Q)
    kf.F = np.array([[1., dt],
                    [0., 1.]])
    kf.H = np.array([[1., 0.]])
    return kf
```

### 8.3.3 Second Order Kalman Filter

A second order Kalman filter tracks a second order system, such as position, velocity and acceleration. The state variable will be

$$\mathbf{x} = \begin{bmatrix} x \\ \dot{x} \\ \ddot{x} \end{bmatrix}$$

So now we have to design our state transition. The Newtonian equations for a time step are:

$$\begin{aligned}x_t &= x_{t-1} + v_{t-1}\Delta t + 0.5a_{t-1}\Delta t^2 \\v_t &= v_{t-1} + a_{t-1}\Delta t \\a_t &= a_{t-1}\end{aligned}$$

Recall that we need to convert this into the linear equation

$$\begin{bmatrix} x \\ \dot{x} \\ \ddot{x} \end{bmatrix} = \mathbf{F} \begin{bmatrix} x \\ \dot{x} \\ \ddot{x} \end{bmatrix}$$

Setting

$$\mathbf{F} = \begin{bmatrix} 1 & \Delta t & .5\Delta t^2 \\ 0 & 1 & \Delta t \\ 0 & 0 & 1 \end{bmatrix}$$

gives us the equations above.

Finally, we design the measurement function. The measurement function needs to implement

$$z = \mathbf{H}\mathbf{x}$$

Our sensor still only reads position, so it should take the position from the state, and 0 out the velocity, like so:

$$\mathbf{H} = [1 \ 0 \ 0]$$

This function constructs and returns a second order Kalman filter.

```
In [18]: def SecondOrderKF(R_std, Q, dt, P=100):
    """ Create second order Kalman filter.
    Specify R and Q as floats."""
    kf = KalmanFilter(dim_x=3, dim_z=1)
    kf.x = np.zeros(3)
    kf.P[0, 0] = P
    kf.P[1, 1] = 1
    kf.P[2, 2] = 1
    kf.R *= R_std**2
    kf.Q = Q_discrete_white_noise(3, dt, Q)
    kf.F = np.array([[1., dt, .5*dt*dt],
                    [0., 1., dt],
                    [0., 0., 1.]])
    kf.H = np.array([[1., 0., 0.]])
    return kf
```

### 8.3.4 Evaluating the Performance

Now we can run each Kalman filter against the simulation and evaluate the results.

How do we evaluate the results? We can do this qualitatively by plotting the track and the Kalman filter output and eyeballing the results. However, a rigorous approach uses mathematics. Recall that system covariance matrix  $\mathbf{P}$  contains the computed variance and covariances for each of the state variables. The diagonal contains the variance. Remember that roughly 99% of all measurements fall within  $3\sigma$  if the noise

is Gaussian. If this is not clear please review the Gaussian chapter before continuing, as this is an important point.

So we can evaluate the filter by looking at the residuals between the estimated state and actual state and comparing them to the standard deviations which we derive from  $\mathbf{P}$ . If the filter is performing correctly 99% of the residuals will fall within  $3\sigma$ . This is true for all the state variables, not just for the position.

I must mention that this is only true for simulated systems. Real sensors are not perfectly Gaussian, and you may need to expand your criteria to, say,  $5\sigma$  with real sensor data.

So let's run the first order Kalman filter against our first order system and access it's performance. You can probably guess that it will do well, but let's look at it using the standard deviations.

First, let's write a routine to generate the noisy measurements for us.

```
In [19]: def simulate_system(Q, count):
    obj = ConstantVelocityObject(x0=.0, vel=0.5, noise_scale=Q)
    xs, zs = [], []
    for i in range(count):
        x = obj.update()
        z = sense(x)
        xs.append(x)
        zs.append(z)
    return np.asarray(xs), np.asarray(zs)
```

And now a routine to perform the filtering.

```
In [20]: def filter_data(kf, zs):
    xs, ps = [], []
    for z in zs:
        kf.predict()
        kf.update(z)

        xs.append(kf.x)
        ps.append(kf.P.diagonal()) # just save variances

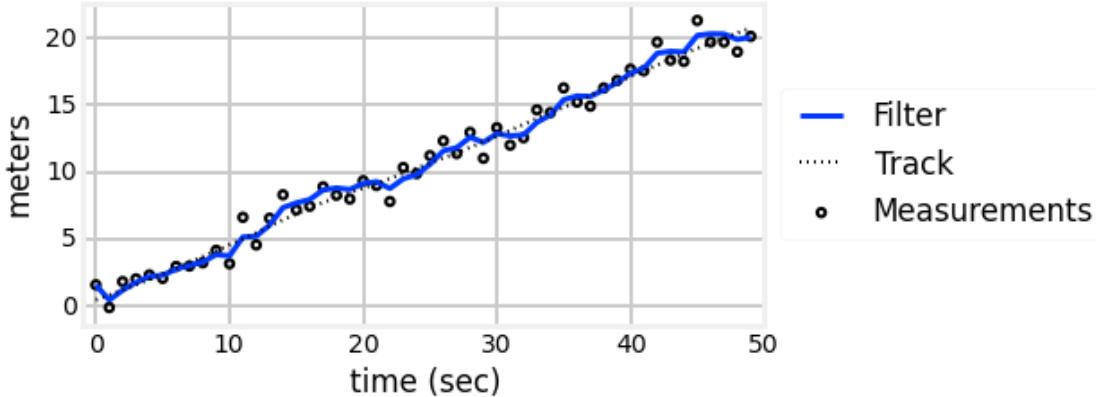
    return np.asarray(xs), np.asarray(ps)
```

Now we are prepared to run the filter and look at the results.

```
In [21]: R, Q = 1, 0.03
xs, zs = simulate_system(Q=Q, count=50)

kf = FirstOrderKF(R, Q, dt=1)
filter_xs1, ps1 = filter_data(kf, zs)

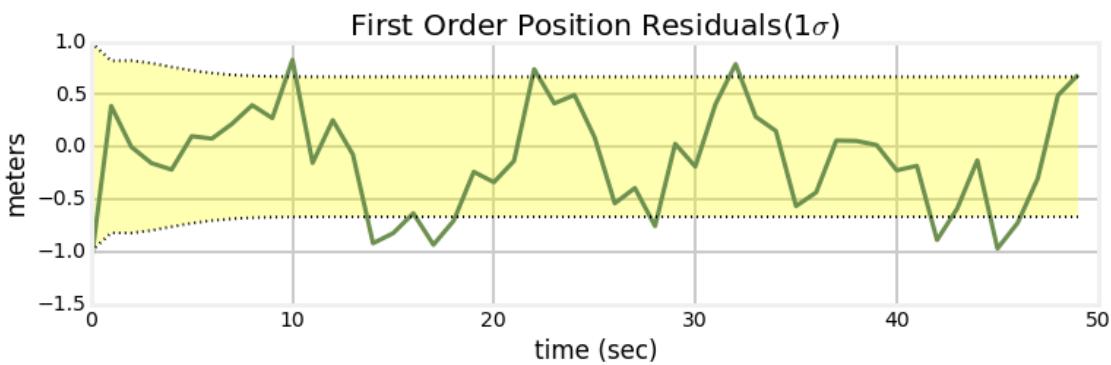
plt.figure()
bp.plot_kf_output(xs, filter_xs1, zs)
```



It looks like the filter is performing well, but it is hard to tell exactly how well. Let's look at the residuals and see if they help. You may have noticed that in the code above I saved the covariance at each step. I did that to use in the following plot. The `ConstantVelocityObject` class returns a tuple of (position, velocity) for the real object, and this is stored in the array `xs`, and the filter's estimates are in `filter_xs1`.

```
In [22]: def plot_residuals(xs, filter_xs, Ps, title, y_label, stds=1):
    res = xs - filter_xs
    plt.plot(res)
    bp.plot_residual_limits(Ps, stds)
    bp.set_labels(title, 'time (sec)', y_label)
    plt.show()

In [23]: plot_residuals(xs[:, 0], filter_xs1[:, 0], ps1[:, 0],
                      title='First Order Position Residuals( $1\sigma$ )',
                      y_label='meters')
```

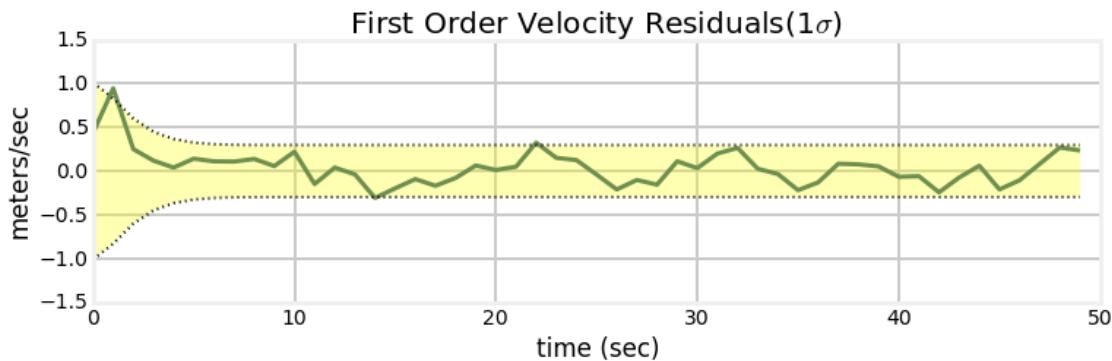


How do we interpret this plot? The residual is drawn as the jagged line - the difference between the measurement and the predicted position. If there was no measurement noise and the Kalman filter prediction was always perfect the residual would always be zero. So the ideal output would be a horizontal line at 0. We can see that the residual is centered around 0, so this gives us confidence that the noise is Gaussian (because the errors fall equally above and below 0). The yellow area between dotted lines show the theoretical

performance of the filter for 1 standard deviations. In other words, approximately 68% of the errors should fall within the dotted lines. The residual falls within this range, so we see that the filter is performing well, and that it is not diverging.

Let's look at the residuals for velocity.

```
In [24]: plot_residuals(xs[:, 1], filter_xs1[:, 1], ps1[:, 1],
                      title='First Order Velocity Residuals( $1\sigma$ )',
                      y_label='meters/sec')
```

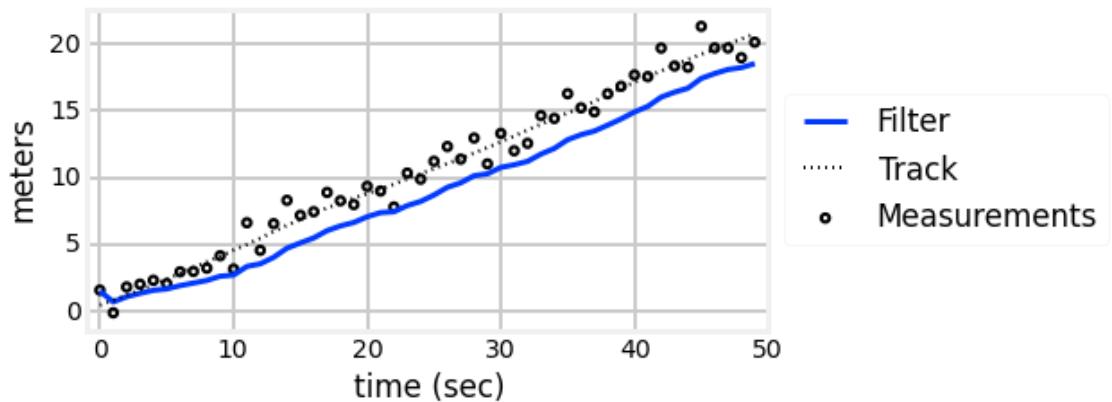


Again, as expected, the residual falls within the theoretical performance of the filter, so we feel confident that the filter is well designed for this system.

Now let's do the same thing using the zero order Kalman filter. All of the code and math is largely the same, so let's just look at the results without discussing the implementation much.

```
In [25]: from book_plots import plot_kf_output
```

```
kf0 = ZeroOrderKF(R, Q)
filter_xs0, ps0 = filter_data(kf0, zs)
plot_kf_output(xs, filter_xs0, zs)
```



As we would expect, the filter has problems. Think back to the g-h filter, where we incorporated acceleration into the system. The g-h filter always lagged the input because there were not enough terms to allow the filter to adjust quickly enough to the changes in velocity. The same thing is happening here, just one order lower. On every `predict()` step the Kalman filter assumes that there is no change in position - if the current position is 4.3 it will predict that the position at the next time period is 4.3. Of course, the actual position is closer to 5.3. The measurement, with noise, might be 5.4, so the filter chooses an estimate part way between 4.3 and 5.4, causing it to lag the actual value of 5.3 by a significant amount. This same thing happens in the next step, the next one, and so on. The filter never catches up.

This raises a very important point. The assumption of ‘constant’ is an assumption of constant-ness between discrete samples only. The filter’s output can still change over time.

Now let’s look at the residuals. We are not tracking velocity, so we can only look at the residual for position.

```
In [26]: plot_residuals(xs[:, 0], filter_xs0[:, 0], ps0[:, 0],
                      title='Zero Order Position Residuals(3$\sigma)'),
                      y_label='meters',
                      stds=3)
```

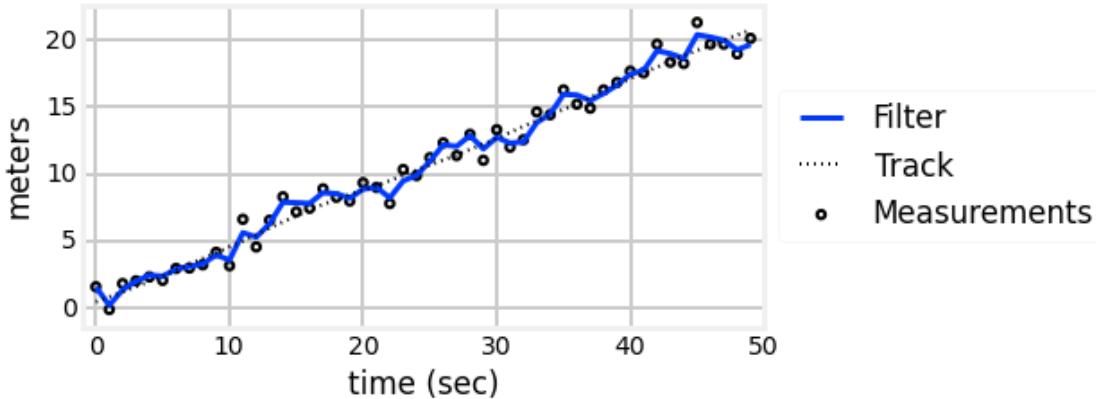


We can see that the filter diverges almost immediately. After a few seconds the residual exceeds the bounds of three standard deviations. It is important to understand that the covariance matrix  $\mathbf{P}$  is only reporting the theoretical performance of the filter assuming all of the inputs are correct. In other words, this Kalman filter is diverging, but  $\mathbf{P}$  implies that the Kalman filter’s estimates are getting better and better with time because the variance is getting smaller. The filter has no way to know that you are lying to it about the system. This is sometimes referred to a smug filter - it is overconfident in its performance.

In this system the divergence is immediate and striking. In many systems it will only be gradual, and/or slight. It is important to look at charts like these for your systems to ensure that the performance of the filter is within the bounds of its theoretical performance.

Now let’s try a second order system. This might strike you as a good thing to do. After all, we know there is a bit of noise in the movement of the simulated object, which implies there is some acceleration. Why not model the acceleration with a second order model? If there is no acceleration, the acceleration should just be estimated to be 0, right?. But is that what happens? Think about it before going on.

```
In [27]: kf2 = SecondOrderKF(R, Q, dt=1)
filter_xs2, ps2 = filter_data(kf2, zs)
plot_kf_output(xs, filter_xs2, zs)
```



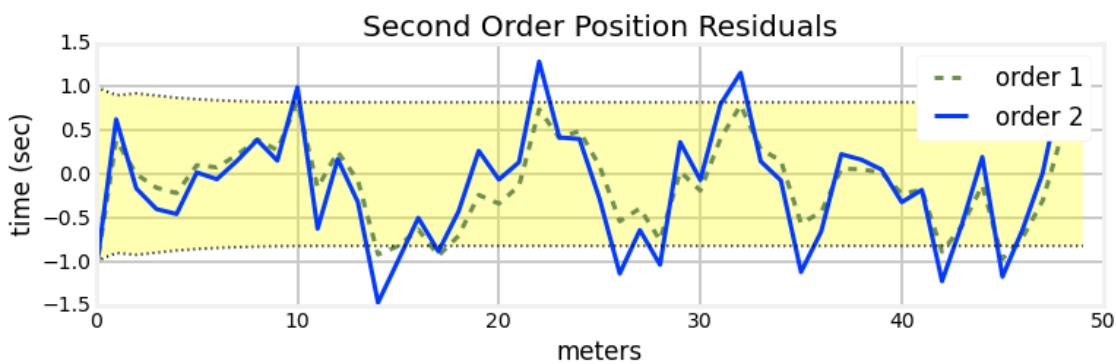
Did this perform as you expected?

We can see that second order filter performs poorly compared to the first order filter. Why? This filter models acceleration, and so the large changes in the measurement gets interpreted as acceleration instead of noise. Thus the filter close tracks the noise. Not only that, but it overshoots the noise in places if the noise is consistently above or below the track because the filter incorrectly assumes an acceleration that does not exist, and so its prediction goes further and further away from the track on each measurement. This is not a good state of affairs.

Still, the track doesn't look horrible. Let's see the story that the residuals tell. I will add a wrinkle here. The residuals for the second order system do not look terrible in that they do not diverge or exceed three standard deviations. However, it is very telling to look at the residuals for the first order vs the second order filter, so I have plotted both on the same graph.

```
In [28]: res = xs[:, 0] - filter_xs2[:, 0]
res1 = xs[:, 0] - filter_xs1[:, 0]

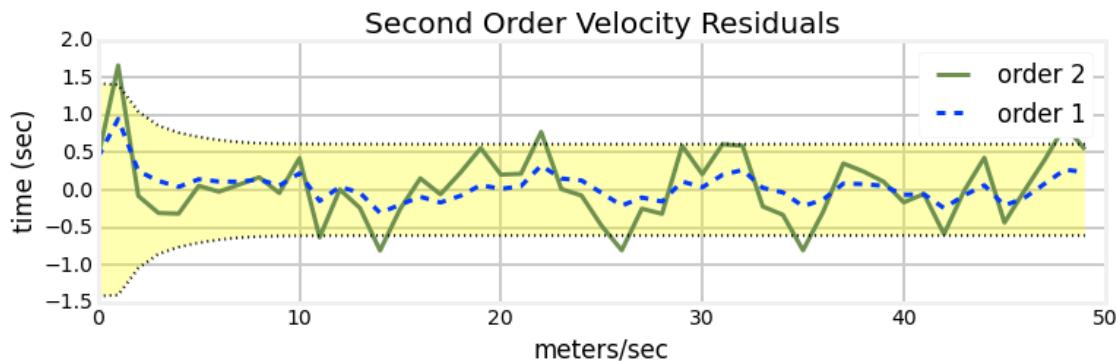
plt.plot(res1, ls="--", label='order 1')
plt.plot(res, label='order 2')
bp.plot_residual_limits(ps2[:, 0])
bp.set_labels('Second Order Position Residuals',
              'meters', 'time (sec)')
plt.legend();
```



The second order position residuals are slightly worse than the residuals of the first order filter, but they still fall within the theoretical limits of the filter. There is nothing very alarming here.

Now let's look at the residuals for the velocity.

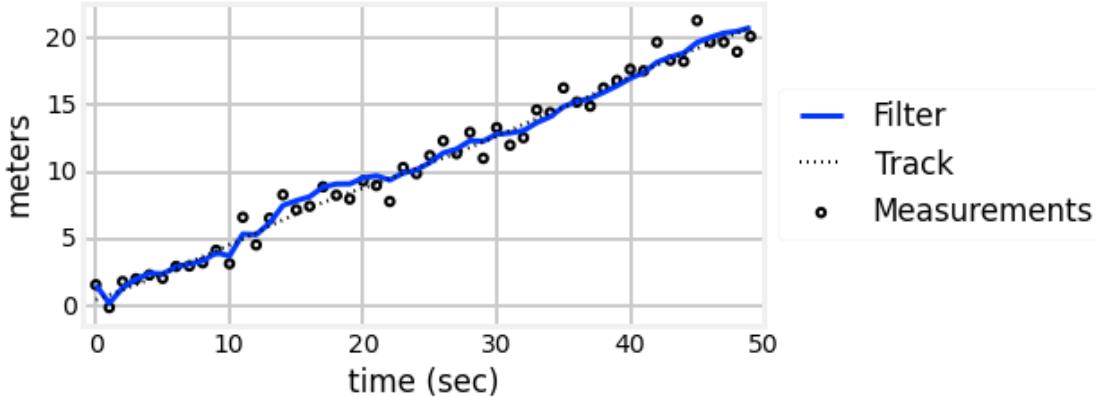
```
In [29]: res = xs[:, 1] - filter_xs2[:, 1]
res1 = xs[:, 1] - filter_xs1[:, 1]
plt.plot(res, label='order 2')
plt.plot(res1, ls='--', label='order 1')
bp.plot_residual_limits(ps2[:, 1])
bp.set_labels('Second Order Velocity Residuals',
              'meters/sec', 'time (sec)')
plt.legend();
```



Here the story is very different. While the residuals of the second order system fall within the theoretical bounds of the filter's performance, we can see that the residuals are far worse than for the first order filter. This is the usual result for this scenario. The filter is assuming that there is acceleration that does not exist. It mistakes noise in the measurement as acceleration and this gets added into the velocity estimate on every predict cycle. Of course the acceleration is not actually there and so the residual for the velocity is much larger than is optimum.

I have one more trick up my sleeve. We have a first order system; i.e. the velocity is more-or-less constant. Real world systems are never perfect, so of course the velocity is never exactly the same between time periods. When we use a first order filter we account for that slight variation in velocity with the process noise. The matrix **Q** is computed to account for this slight variation. If we move to a second order filter we are now accounting for the changes in velocity. Perhaps now we have no process noise, and we can set **Q** to zero!

```
In [30]: kf2 = SecondOrderKF(R, 0, dt=1)
filter_xs2, ps2 = filter_data(kf2, zs)
plot_kf_output(xs, filter_xs2, zs)
```



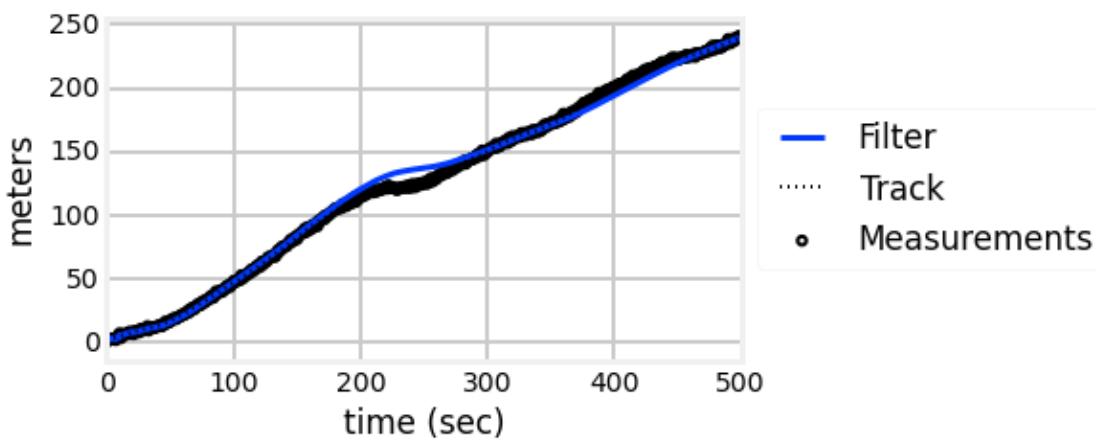
To my eye it looks like the filter quickly converges to the actual track. Success!

Or, maybe not. Setting the process noise to 0 tells the filter that the process model is perfect. Let's look at the performance of the filter over a longer period of time.

```
In [31]: np.random.seed(25944)
xs500, zs500 = simulate_system(Q=Q, count=500)

kf2 = SecondOrderKF(R, 0, dt=1)
filter_xs2, ps2 = filter_data(kf2, zs500)

plot_kf_output(xs500, filter_xs2, zs500)
plot_residuals(xs500[:, 0], filter_xs2[:, 0], ps2[:, 0],
               'Second Order Position Residuals',
               'meters')
```





We can see that the performance of the filter is abysmal. We can see that in the track plot where the filter diverges from the track for an extended period of time. The residual plot makes the problem more apparent. Just before the 100th update the filter diverges sharply from the theoretical performance. It might be converging at the end, but I doubt it. The entire time, the filter is reporting smaller and smaller variances. **Do not trust the filter's covariance matrix to tell you if the filter is performing well!**

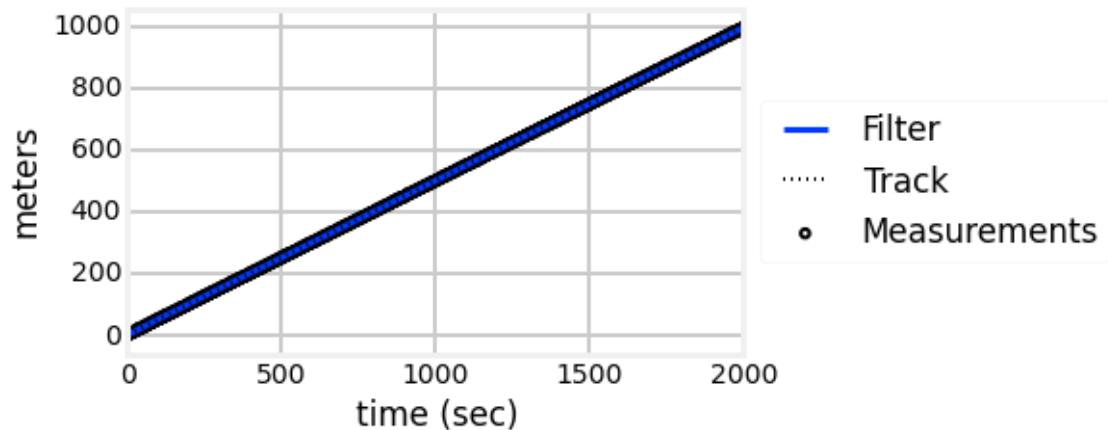
Why is this happening? Recall that if we set the process noise to zero we are telling the filter to use only the process model. The measurements end up getting ignored. The physical system is not perfect, and so the filter is unable to adapt to this imperfect behavior.

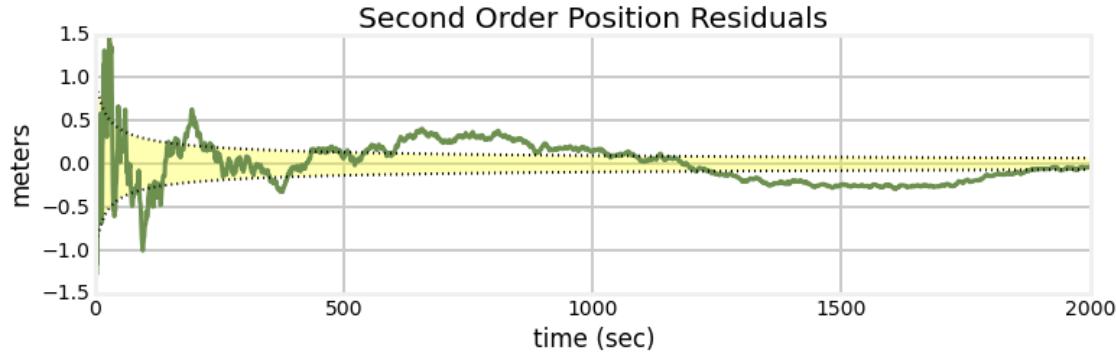
Maybe just a really low process noise? Let's try that.

```
In [32]: np.random.seed(32594)
        xs2000, zs2000 = simulate_system(Q=0.0001, count=2000)

        kf2 = SecondOrderKF(R, 0, dt=1)
        filter_xs2, ps2 = filter_data(kf2, zs2000)

        plot_kf_output(xs2000, filter_xs2, zs2000)
        plot_residuals(xs2000[:, 0], filter_xs2[:, 0], ps2[:, 0],
                       'Second Order Position Residuals',
                       'meters')
```



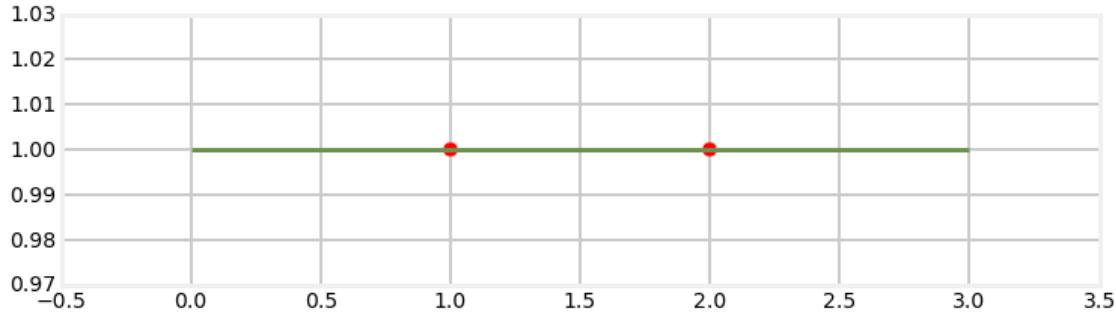


Again, the residual plot tells the story. The track looks very good, but the residual plot shows that the filter is diverging for significant periods of time.

How should you think about all of this? You might argue that the last plot is ‘good enough’ for your application, and perhaps it is. I warn you however that a diverging filter doesn’t always converge. With a different data set, or a physical system that performs differently you can end up with a filter that veers further and further away from the measurements.

Also, let’s think about this in a data fitting sense. Suppose I give you two points, and tell you to fit a straight line to the points.

```
In [33]: plt.scatter([1, 2], [1, 1], s=100, c='r')
plt.plot([0, 3], [1, 1]);
```



A straight line is the only possible answer. Furthermore, the answer is optimal. If I gave you more points you could use a least squares fit to find the best line, and the answer would still be optimal (in a least squares sense).

But suppose I told you to fit a higher order polynomial to those two points. There is now an infinite number of answers to the problem. For example, infinite number of parabolas (which are second order) pass through those points. When the Kalman filter is of higher order than your physical process it also has an infinite number of solutions to choose from. The answer is not just non-optimal, it often diverges and never recovers.

For best performance you need a filter whose order matches the system’s order.. In many cases that will be easy to do - if you are designing a Kalman filter to read the thermometer of a freezer it seems clear that a zero order filter is the right choice. But what order should we use if we are tracking a car? Order one will work well while the car is moving in a straight line at a constant speed, but cars turn, speed up, and

slow down, in which case a second order filter will perform better. That is the problem addressed in the Adaptive Filters chapter. There we will learn how to design a filter that adapts to changing order in the tracked object's behavior.

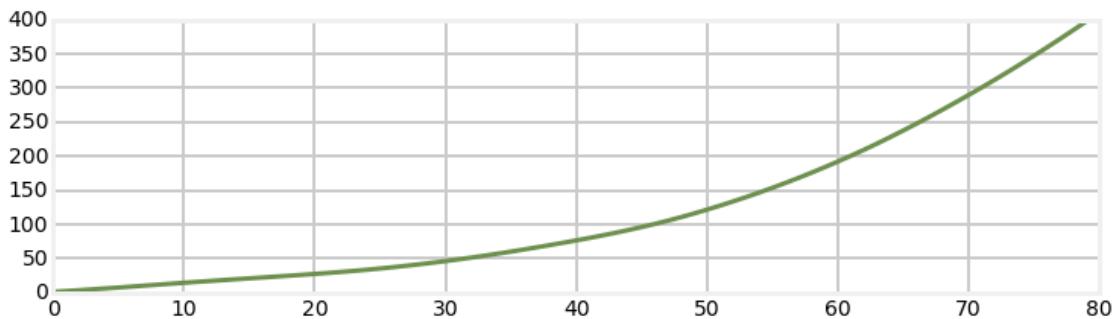
With that said, a lower order filter can track a higher order process so long as you add enough process noise and you keep the discretization period small (100 samples a second are usually locally linear). The results will not be optimal, but they can still be very good, and I always reach for this tool first before trying an adaptive filter. Let's look at an example with acceleration. First, the simulation.

```
In [34]: class ConstantAccelerationObject(object):
    def __init__(self, x0=0, vel=1., acc=0.1, acc_noise=.1):
        self.x = x0
        self.vel = vel
        self.acc = acc
        self.acc_noise_scale = acc_noise

    def update(self):
        self.acc += randn() * self.acc_noise_scale
        self.vel += self.acc
        self.x += self.vel
        return (self.x, self.vel, self.acc)

R, Q = 6., 0.02
def simulate_acc_system(R, Q, count):
    obj = ConstantAccelerationObject(acc_noise=Q)
    zs = []
    xs = []
    for i in range(count):
        x = obj.update()
        z = sense(x, R)
        xs.append(x)
        zs.append(z)
    return np.asarray(xs), zs

np.random.seed(124)
xs, zs = simulate_acc_system(R=R, Q=Q, count=80)
plt.plot(xs[:, 0]);
```



Now we will filter the data using a second order filter.

```
In [35]: np.random.seed(124)
```

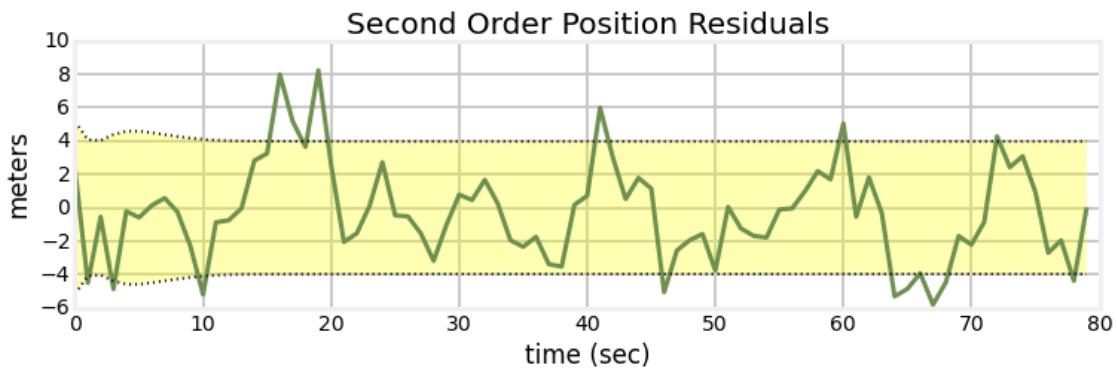
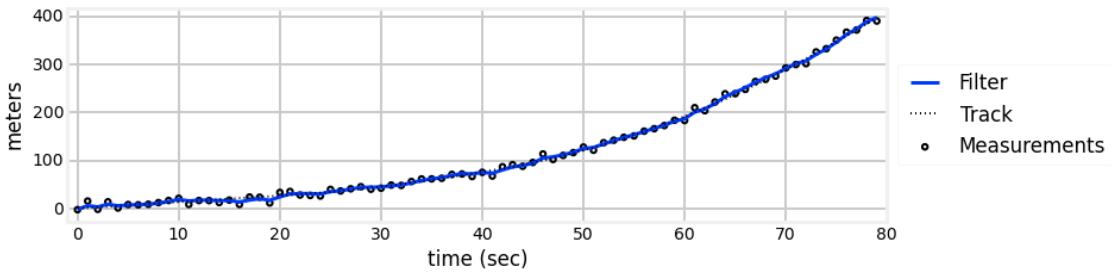
```

xs, zs = simulate_acc_system(R=R, Q=Q, count=80)

kf2 = SecondOrderKF(R, Q, dt=1)
fxs2, ps2 = filter_data(kf2, zs)

plot_kf_output(xs, fxs2, zs, aspect_equal=False)
plot_residuals(xs[:, 0], fxs2[:, 0], ps2[:, 0],
               'Second Order Position Residuals',
               'meters')

```



We can see that the filter is performing within the theoretical limits of the filter.

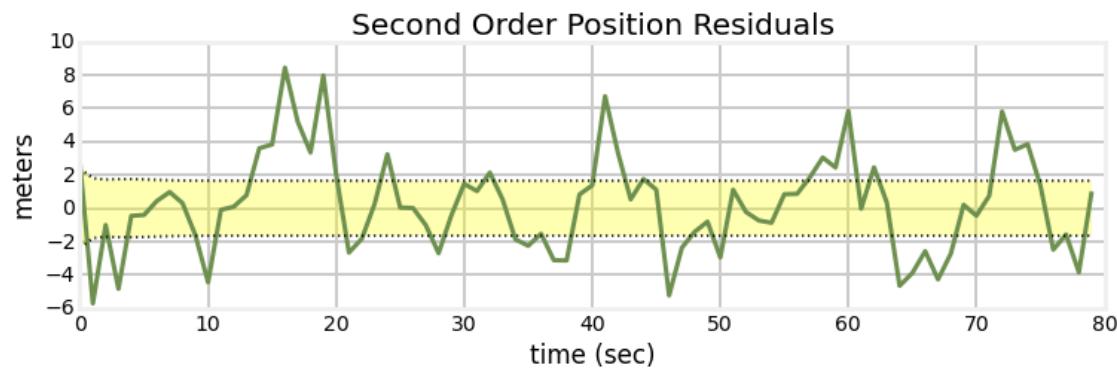
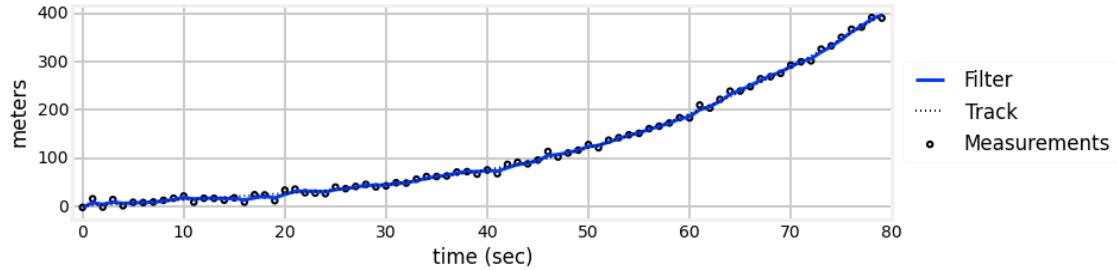
Now let's use a lower order filter. As already demonstrated the lower order filter will lag the signal because it is not modeling the acceleration. However, we can account for that (to an extent) by increasing the size of the process noise. The filter will treat the acceleration as noise in the process model. The result will be suboptimal, but if designed well it will not diverge. Choosing the amount of extra process noise is not an exact science. You will have to experiment with representative data. Here, I've multiplied it by 10, and am getting good results.

```

In [36]: kf2 = FirstOrderKF(R, Q * 10, dt=1)
         fxs2, ps2 = filter_data(kf2, zs)

plot_kf_output(xs, fxs2, zs, aspect_equal=False)
plot_residuals(xs[:, 0], fxs2[:, 0], ps2[:, 0],
               'Second Order Position Residuals',
               'meters')

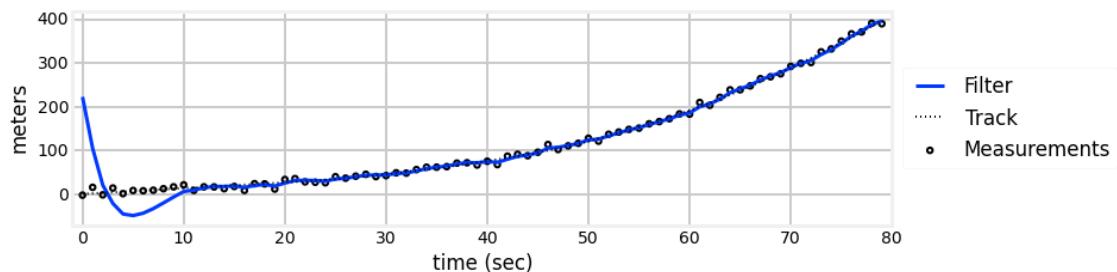
```



Think about what will happen if you make the process noise many times larger than it needs to be. A large process noise tells the filter to favor the measurements, so we would expect the filter to closely mimic the noise in the measurements. Let's find out.

```
In [37]: kf1 = FirstOrderKF(R, Q * 10000, dt=1)
        fxs1, ps1 = filter_data(kf2, zs)

        plot_kf_output(xs, fxs1, zs, aspect_equal=False)
        plot_residuals(xs[:, 0], fxs1[:, 0], ps2[:, 0],
                       'First Order Position Residuals',
                       'meters')
```





## 8.4 Evaluating Filter Performance

It is easy to design a Kalman filter for a simulated situation. You know how much noise you are injecting in your process model, so you specify  $\mathbf{Q}$  to have the same value. You also know how much noise is in the measurement simulation, so the measurement noise matrix  $\mathbf{R}$  is equally trivial to define.

In practice design is more ad hoc. Real sensors rarely perform to spec, and they rarely perform in a Gaussian manner. They are also easily fooled by environmental noise. For example, circuit noise causes voltage fluctuations which can affect the output of a sensor. Creating a process model and noise is even more difficult. Modelling an automobile is very difficult. The steering causes nonlinear behavior, tires slip, people brake and accelerate hard enough to cause tire slips, winds push the car off course. The end result is the Kalman filter is an inexact model of the system. This inexactness causes suboptimal behavior, which in the worst case causes the filter to diverge completely.

Because of the unknowns you will be unable to analytically compute the correct values for the filter's matrices. You will start by making the best estimate possible, and then test your filter against a wide variety of simulated and real data. Your evaluation of the performance will guide you towards what changes you need to make to the matrices. We've done some of this already - I've shown you the effect of  $\mathbf{Q}$  being too large or small.

Now let's consider more analytical ways of assessing performance. If the Kalman filter is performing optimally its estimation errors (the difference between the true state and the estimated state will have these properties:

1. The mean of the estimation error is zero
2. Its covariance is described by the Kalman filter's covariance matrix

### 8.4.1 Normalized Estimated Error Squared (NEES)

The first method is the most powerful, but only possible in simulations. If you are simulating a system you know its true value, or 'ground truth'. It is then trivial to compute the error of the system at any step as the difference between ground truth ( $\mathbf{x}$ ) and the filter's state estimate ( $\hat{\mathbf{x}}$ ):

$$\tilde{\mathbf{x}} = \mathbf{x} - \hat{\mathbf{x}}$$

We can then define the **normalized estimated error squared** (NEES) as

$$\epsilon = \tilde{\mathbf{x}}^T \mathbf{P}^{-1} \tilde{\mathbf{x}}$$

To understand this equation lets look at it if the state's dimension is one. In that case both  $x$  and  $P$  are scalars, so

$$\epsilon = \frac{x^2}{P}$$

If this is not clear, recall that the transpose of a scalar is the same scalar, and that  $a^{-1} = \frac{1}{a}$  for a scalar.

So as the covariance matrix gets smaller NEES gets larger for the same error. The covariance matrix is the filter's estimate of it's error, so if it is small relative to the estimation error then it is performing worse than if it is large relative to the same estimation error.

This computation gives us a scalar result. If  $x$  is dimension  $(n \times 1)$ , then the dimension of the computation is  $(1 \times n)(n \times n)(n \times 1) = (1 \times 1)$ . What do we do with this number?

The math is outside the scope of this book, but a random variable in the form  $\tilde{x}^T P^{-1} \tilde{x}$  is said to be **chi-squared distributed with  $n$  degrees of freedom**, and thus the expected value of the sequence should be  $n$ . Bar-Shalom [1] has an excellent discussion of this topic.

In plain terms, take the average of all the NESS values, and they should be less then the dimension of  $x$ . Let's see that using an example from earlier in the chapter:

```
In [38]: from scipy.linalg import inv

def NEES(xs, est_xs, ps):
    est_err = xs - est_xs
    err = []
    for x, p in zip(est_err, ps):
        err.append(np.dot(x.T, inv(p)).dot(x))
    return err

In [39]: R, Q = 6., 0.02
xs, zs = simulate_acc_system(R=R, Q=Q, count=80)
kf2 = SecondOrderKF(R, Q, dt=1)
est_xs, ps, _, _ = kf2.batch_filter(zs)

nees = NEES(xs, est_xs, ps)
eps = np.mean(nees)

print('mean NEES is: ', eps)
if eps < kf2.dim_x:
    print('passed')
else:
    print('failed')

mean NEES is:  1.04110092605
passed
```

NESS is implemented FilterPy; access it with

```
from filterpy.stats import NESS
```

This is the best measure of the filter's performance, and should be used whenever possible. If your simulation is of limited fidelity then you need to use another approach.

### 8.4.2 Likelihood Function

In statistics the likelihood is very similar to a probability, with a subtle difference that is important to us. A probability is the chance of something happening - as in what is the probability that a fair die rolls 6 three times in five rolls? The likelihood asks the reverse question - given that a die rolled 6 three times in five rolls, what is the likelihood that the die is fair?

This is important to us because we have the filter output and we want to know the likelihood that it is performing optimally given the assumptions of Gaussian noise and linear behavior. If the likelihood is low we know one of our assumptions is wrong. In the **Adaptive Filtering** chapter we will learn how to make use of this information to improve the filter; here we will only learn to make the measurement.

The residual and system uncertainty of the filter is defined as

$$\begin{aligned}\mathbf{y} &= \mathbf{z} - \mathbf{H}\bar{\mathbf{x}} \\ \mathbf{S} &= \mathbf{H}\bar{\mathbf{P}}\mathbf{H}^T + \mathbf{R}\end{aligned}$$

Given these we can compute the likelihood function with

$$\mathcal{L} = \frac{1}{\sqrt{2\pi S}} \exp\left[-\frac{1}{2}\mathbf{y}^T \mathbf{S}^{-1} \mathbf{y}\right]$$

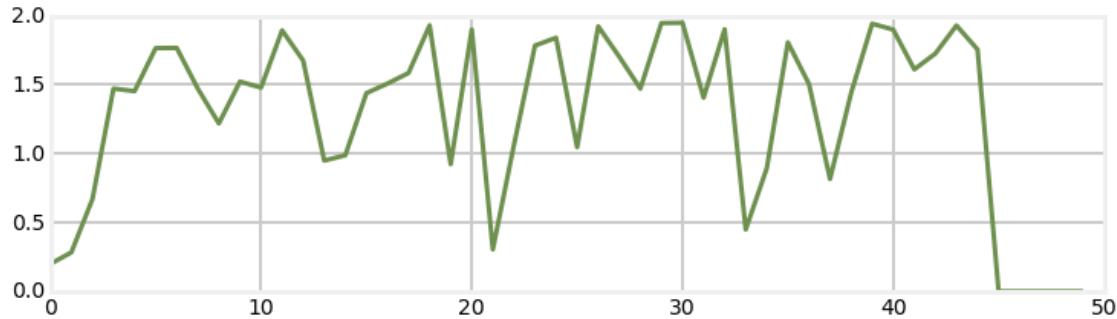
That may look complicated, but notice that the exponential is the equation for a Gaussian. This suggests an implementation of

```
from scipy.stats import multivariate_normal
hx = np.dot(H, x).flatten()
S = np.dot(H, P).dot(H.T) + R
likelihood = multivariate_normal.pdf(z.flatten(), mean=hx, cov=S)
```

This is computed for you when `update()` is called, and is accessible via the `likelihood` data attribute. Let's look at this: I'll run the filter with several measurements within expected range, and then inject measurements far from the expected values:

```
In [40]: R, Q = .05, 0.02
        xs, zs = simulate_acc_system(R=R, Q=Q, count=50)
        zs[-5:-1] = [100, 200, 200, 200] # bad measurements, bad!

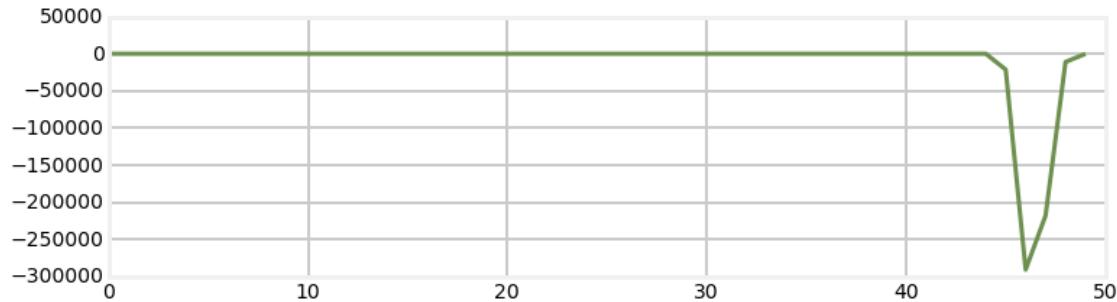
        kf = SecondOrderKF(R, Q, dt=1, P=1)
        likelihoods = []
        for z in zs:
            kf.predict()
            kf.update(z)
            likelihoods.append(kf.likelihood)
plt.plot(likelihoods);
```



The likelihood gets larger as the filter converges during the first few iterations. After that the likelihood bounces around until it reaches the bad measurements, at which time it goes to zero, indicating that if the measurements are valid the filter is very unlikely to be optimal.

The log likelihood is somewhat more convenient to use in some cases. By taking the log we guarantee that the function is monotonically increasing, and its maximum values occur at the same points as the function itself. See how starkly it illustrates where the filter goes ‘bad’.

```
In [41]: kf = SecondOrderKF(R, Q, dt=1, P=1)
log_likelihoods = []
for z in zs:
    kf.predict()
    kf.update(z)
    log_likelihoods.append(kf.log_likelihood)
plt.plot(log_likelihoods);
```



Why does it return to zero at the end? Think about that before reading the answer. The filter begins to adapt to the new measurements by moving the state close to the measurements. The residuals become small, and so the state and residuals agree.

### 8.4.3 NIS

todo

## 8.5 Control Inputs

In the **Discrete Bayes** chapter I introduced the notion of using control signals to improve the filter's performance. Instead of assuming that the object continues to move as it has so far we use our knowledge of the control inputs to predict where the object is. I did not delve into the details, but in the **Univariate Kalman Filter** chapter we made use of the same idea. The predict method of the Kalman filter read

```
def predict(self, u=0.0):
    self.x += u
    self.P += self.Q
```

Here  $u$  is the control which is assumed to be the change in the state  $x$  effected by the control input.

In this chapter we learned that the equation for the state prediction is:

$$\bar{x} = \mathbf{F}x + \mathbf{B}u$$

Our state is a vector, so we need to represent the control input as a vector. Here  $\mathbf{u}$  is the control input, and  $\mathbf{B}$  is a matrix that transforms the control input into a change in  $\mathbf{x}$ . Let's consider a simple example. Suppose the state is  $x = [x \ \dot{x}]$  for a robot we are controlling and the control input is commanded velocity. This gives us a control input of

$$\mathbf{u} = [\dot{x}_{\text{cmd}}]$$

For simplicity we will assume that the robot can respond instantly to changes to this input. That means that the new position and velocity after  $\Delta t$  seconds will be

$$\begin{aligned} x &= x + \dot{x}_{\text{cmd}}\Delta t \\ \dot{x} &= \dot{x}_{\text{cmd}} \end{aligned}$$

We need to represent this set of equation in the form  $\bar{x} = \mathbf{F}x + \mathbf{B}u$ .

I will use the  $\mathbf{F}x$  term to extract the  $x$  for the top equation, and the  $\mathbf{B}u$  term for the rest, like so:

$$\begin{bmatrix} x \\ \dot{x} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \end{bmatrix} + \begin{bmatrix} \Delta t \\ 1 \end{bmatrix} [\dot{x}_{\text{cmd}}]$$

This is a simplification; typical control inputs are changes to steering angle and changes in acceleration. This introduces nonlinearities which we will learn to deal with in a later chapter.

The rest of the Kalman filter will be designed as normal. You've seen this several times by now, so without further comment here is an example.

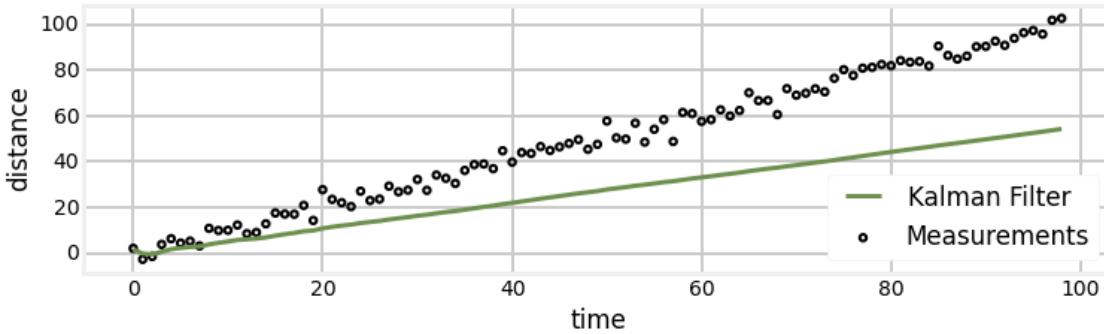
```
In [42]: dt = 0.1
R = 3.
kf = KalmanFilter(dim_x=2, dim_z=1, dim_u = 1)
kf.P *= 10
kf.R *= R
kf.Q = Q_discrete_white_noise(2, dt, 0.1)
kf.F = np.array([[1., 0],
                 [0., 1]])
kf.B = np.array([[dt],
                 [1]])
kf.H = np.array([[1., 0]])
```

```

zs = [i + randn()*R for i in range(1, 100)]
xs = []
cmd_velocity = 1.
for z in zs:
    kf.predict(u=cmd_velocity)
    kf.update(z)
    xs.append(kf.x[0])

plt.plot(xs, label='Kalman Filter')
bp.plot_measurements(zs)
plt.xlabel('time')
plt.legend(loc=4)
plt.ylabel('distance');

```



## 8.6 Sensor Fusion

Early in the g-h filter chapter we discussed designing a filter for two scales, one accurate and one inaccurate. We determined that we should always include the information from the inaccurate filter - we should never discard any information. If you do not recall this discussion, consider this edge case. You have two scales, one accurate to 1 lb, the second to 10 lbs. You weigh yourself, and the first reads 170 lbs, and the second reads 181 lbs. If we only used the data from the first scale we would only know that our weight is in the range of 169 to 171 lbs. However, when we incorporate the measurement from the second scale we can decide that the correct weight can only be 171 lbs, since 171 lbs is exactly at the outside range of the second measurement, which would be 171 to 191 lbs. Of course most measurements don't offer such an exact answer, but the principle applies to any measurement.

So consider a situation where we have two sensors measuring the system. How shall we incorporate that into our Kalman filter?

Suppose we have a train or cart on a railway. It has a sensor attached to the wheels counting revolutions, which can be converted to a distance along the track. Then, suppose we have a GPS-like sensor which I'll call a 'position sensor' mounted to the train which reports position. I'll explain why I don't just use a GPS in the next section. Thus, we have two measurements, both reporting position along the track. Suppose further that the accuracy of the wheel sensor is 1m, and the accuracy of the position sensor is 10m. How do we combine these two measurements into one filter? This may seem quite contrived, but aircraft use sensor fusion to fuse the measurements from sensors such as a GPS, INS, Doppler radar, VOR, the airspeed indicator, and more.

Kalman filters for inertial filters are very difficult, but fusing data from two or more sensors providing measurements of the same state variable (such as position) is quite easy. The relevant matrix is the measurement matrix  $\mathbf{H}$ . Recall that this matrix tells us how to convert from the Kalman filter's state  $\mathbf{x}$  to a measurement  $\mathbf{z}$ . Suppose that we decide that our Kalman filter state should contain the position and velocity of the train, so that

$$\mathbf{x} = \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

We have two measurements for position, so we will define the measurement vector to be a vector of the measurements from the wheel and the position sensor.

$$\mathbf{z} = \begin{bmatrix} x_{wheel} \\ x_{ps} \end{bmatrix}$$

So we have to design the matrix  $\mathbf{H}$  to convert  $\mathbf{x}$  to  $\mathbf{z}$ . They are both positions, so the conversion is nothing more than multiplying by one:

$$\begin{bmatrix} x_{wheel} \\ x_{ps} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

To make it clearer, suppose that the wheel reports not position but the number of rotations of the wheels, where 1 revolution yields 2 meters of travel. In that case we would write

$$\begin{bmatrix} x_{wheel} \\ x_{ps} \end{bmatrix} = \begin{bmatrix} 0.5 & 0 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

Now we have to design the measurement noise matrix  $\mathbf{R}$ . Suppose that the measurement variance for the position is twice the variance of the wheel, and the standard deviation of the wheel is 1.5 meters. That gives us

$$\begin{aligned} \sigma_{wheel} &= 1.5 \\ \sigma_{wheel}^2 &= 2.25 \\ \sigma_{ps} &= 1.5 * 2 = 3 \\ \sigma_{ps}^2 &= 9. \end{aligned}$$

That is pretty much our Kalman filter design. We need to design for  $\mathbf{Q}$ , but that is invariant to whether we are doing sensor fusion or not, so I will just choose some arbitrary value.

So let's run a simulation of this design. I will assume a velocity of 10 m/s with an update rate of 0.1 seconds.

```
In [43]: from numpy import array, asarray
import numpy.random as random

def fusion_test(wheel_sigma, ps_sigma, do_plot=True):
    dt = 0.1
    kf = KalmanFilter(dim_x=2, dim_z=2)

    kf.F = array([[1., dt], [0., 1.]])
    kf.H = array([[1., 0.], [1., 0.]])
    kf.x = array([0., 1.])
    kf.Q **= array([[[(dt**3)/3, (dt**2)/2],
                    [(dt**2)/2, dt]]]) * 0.02
```

```

kf.P *= 100
kf.R[0, 0] = wheel_sigma**2
kf.R[1, 1] = ps_sigma**2

random.seed(1123)
xs, zs, nom = [], [], []
for i in range(1, 100):
    m0 = i + randn()*wheel_sigma
    m1 = i + randn()*ps_sigma
    z = array([[m0], [m1]])

    kf.predict()
    kf.update(z)

    xs.append(kf.x.T[0])
    zs.append(z.T[0])
    nom.append(i)

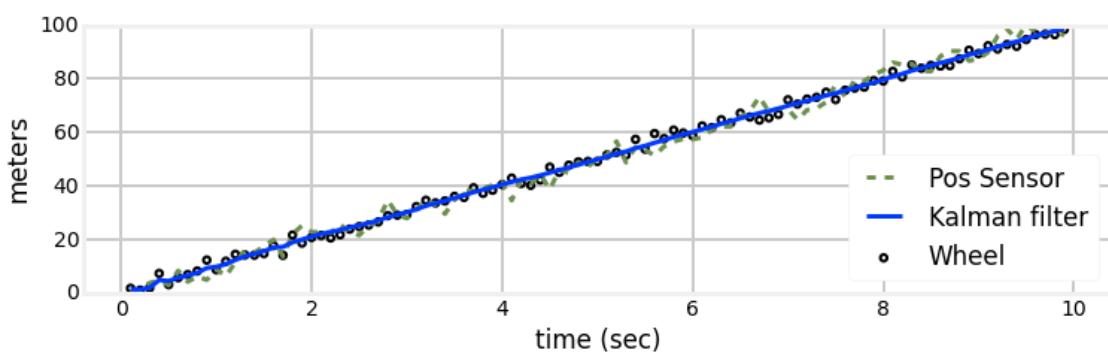
xs = array(xs)
zs = array(zs)
nom = array(nom)

res = nom - xs[:, 0]
print('fusion std: {:.3f}'.format(np.std(res)))
if do_plot:
    ts = np.arange(0.1, 10, .1)
    bp.plot_measurements(ts, zs[:, 0], label='Wheel')
    plt.plot(ts, zs[:, 1], ls='--', label='Pos Sensor')
    bp.plot_filter(ts, xs[:, 0], label='Kalman filter')
    plt.legend(loc=4)
    plt.ylim(0, 100)
    bp.set_labels(x='time (sec)', y='meters')

fusion_test(1.5, 3.0)

```

fusion std: 0.391



We can see the result of the Kalman filter in blue.

It may be somewhat difficult to understand the previous example at an intuitive level. Let's look at a different problem. Suppose we are tracking an object in 2D space, and have two radar systems at different positions. Each radar system gives us a range and bearing to the target. How do the readings from each data affect the results?

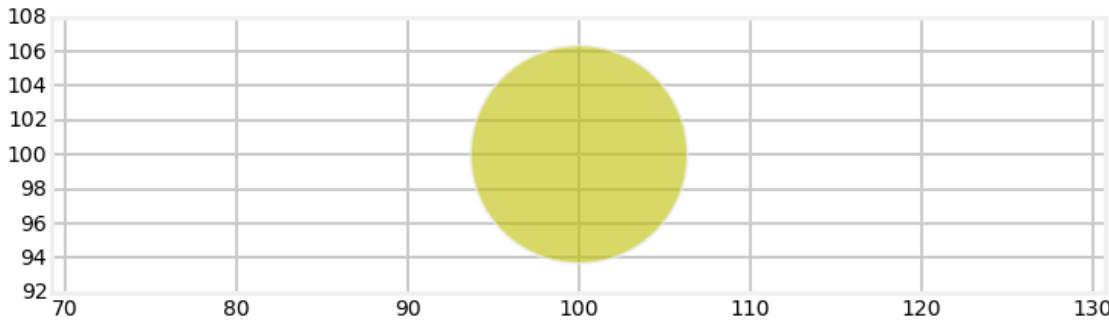
This is a nonlinear problem because we need to use a trigonometry to compute coordinates from a range and bearing, and we have not yet learned how to solve nonlinear problems with Kalman filters. So for this problem ignore the code that I use and just concentrate on the charts that the code outputs. We will revisit this problem in subsequent chapters and learn how to write this code.

I will position the target at (100, 100). The first radar will be at (50, 50), and the second radar at (150, 50). This will cause the first radar to measure a bearing of 45 degrees, and the second will report 135 degrees.

I will create the Kalman filter first, and then plot its initial covariance matrix. I am using a **unscented Kalman filter**.

```
In [44]: from kf_design_internal import sensor_fusion_kf, set_radar_pos
```

```
kf = sensor_fusion_kf()
x0, p0 = kf.x.copy(), kf.P.copy()
plot_covariance_ellipse(x0, p0, fc='y', ec=None, alpha=0.6)
```



We are equally uncertain about the position in x and y, so the covariance is circular.

Now we will update the Kalman filter with a reading from the first radar. I will set the standard deviation of the bearing error at  $0.5^\circ$ , and the standard deviation of the distance error at 3.

```
In [45]: from math import radians
```

```
# set the error of the radar's bearing and distance
kf.R[0, 0] = radians (.5)**2
kf.R[1, 1] = 3.**2

# compute position and covariance from first radar station
set_radar_pos((50, 50))
dist = (50**2 + 50**2) ** 0.5
kf.predict()
kf.update([radians(45), dist])

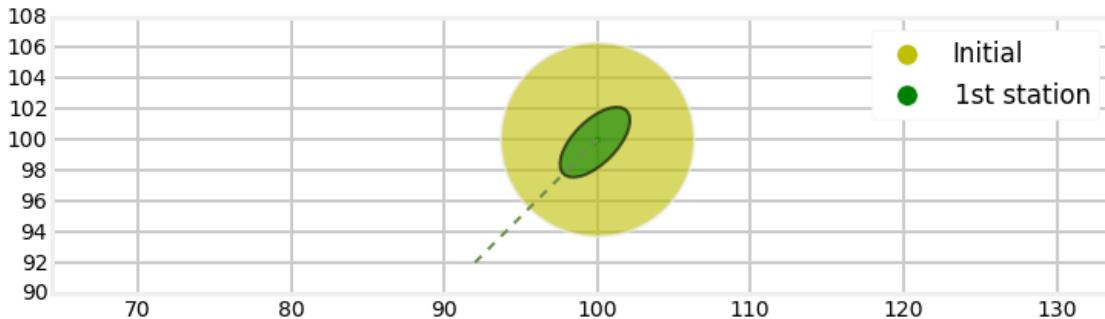
# plot the results
x1, p1 = kf.x.copy(), kf.P.copy()
```

```

plot_covariance_ellipse(x0, p0, fc='y', ec=None, alpha=0.6)
plot_covariance_ellipse(x1, p1, fc='g', ec='k', alpha=0.6)

plt.scatter([100], [100], c='y', label='Initial')
plt.scatter([100], [100], c='g', label='1st station')
plt.legend(scatterpoints=1, markerscale=3)
plt.plot([92, 100], [92, 100], c='g', lw=2, ls='--');

```



We can see the effect of the errors on the geometry of the problem. The radar station is to the lower left of the target. The bearing measurement is extremely accurate at  $\sigma = 0.5^\circ$ , but the distance error is inaccurate at  $\sigma = 3$ . I've shown the radar reading with the dotted green line. We can easily see the effect of the accurate bearing and inaccurate distance in the shape of the covariance ellipse.

Now we can incorporate the second radar station's measurement. The second radar is at (150,50), which is below and to the right of the target. Before you go on, think about how you think the covariance will change when we incorporate this new reading.

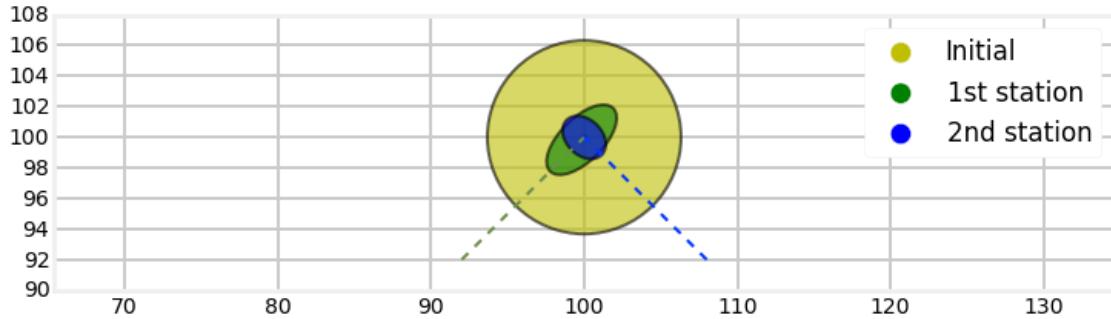
```

In [46]: # compute position and covariance from first radar station
    set_radar_pos((150, 50))
    kf.predict()
    kf.update([radians(135), dist])

    plot_covariance_ellipse(x0, p0, fc='y', ec='k', alpha=0.6)
    plot_covariance_ellipse(x1, p1, fc='g', ec='k', alpha=0.6)
    plot_covariance_ellipse(kf.x, kf.P, fc='b', ec='k', alpha=0.6)

    plt.scatter([100], [100], c='y', label='Initial')
    plt.scatter([100], [100], c='g', label='1st station')
    plt.scatter([100], [100], c='b', label='2nd station')
    plt.legend(scatterpoints=1, markerscale=3)
    plt.plot([92, 100], [92, 100], c='g', lw=2, ls='--')
    plt.plot([108, 100], [92, 100], c='b', lw=2, ls='--');

```



We can see how the second radar measurement altered the covariance. The angle to the target is orthogonal to the first radar station, so the effects of the error in the bearing and range are swapped. So the angle of the covariance matrix switches to match the direction to the second station. It is important to note that the direction did not merely change; the size of the covariance matrix became much smaller as well.

The covariance will always incorporate all of the information available, including the effects of the geometry of the problem. This formulation makes it particularly easy to see what is happening, but the same thing occurs if one sensor gives you position and a second sensor gives you velocity, or if two sensors provide measurements of position.

One final thing before we move on: sensor fusion is a vast topic, and my coverage is simplistic to the point of being misleading. For example, the Kalman Filter Math chapter talks about using iterated least squares to determine the position from a set of pseudorange readings from the satellites without using a Kalman filter. That is the usual but not exclusive way this computation is done in GPS receivers. If you are a hobbyist my coverage may get you started. A commercial grade filter requires very careful design of the fusion process. That is the topic of several books, and you will have to further your education by finding one that covers your domain.

### 8.6.1 Exercise: Can you Kalman Filter GPS outputs?

In the section above I have you apply a Kalman filter to ‘GPS-like’ sensor. Can you apply a Kalman filter to the output of a commercial Kalman filter? In other words, will the output of your filter be better than, worse than, or equal to the GPS’s output?

#### Solution

Commercial GPS’s have a Kalman filter built into them, and their output is the filtered estimate created by that filter. So, suppose you have a steady stream of output from the GPS consisting of a position and position error. Can you not pass those two pieces of data into your own filter?

Well, what are the characteristics of that data stream, and more importantly, what are the fundamental requirements of the input to the Kalman filter?

Inputs to the Kalman filter must be Gaussian and time independent. The output of the GPS is time dependent because the filter bases its current estimate on the recursive estimates of all previous measurements. Hence, the signal is not white, it is not time independent, and if you pass that data into a Kalman filter you have violated the mathematical requirements of the filter. So, the answer is no, you cannot get better estimates by running a KF on the output of a commercial GPS.

Another way to think of it is that Kalman filters are optimal in a least squares sense. There is no way to take an optimal solution, pass it through a filter, any filter, and get a ‘more optimal’ answer because it is a

logical impossibility. At best the signal will be unchanged, in which case it will still be optimal, or it will be changed, and hence no longer optimal.

This is a difficult problem that hobbyists face when trying to integrate GPS, IMU's and other off the shelf sensors.

Let's look at the effect. A commercial GPS reports position, and an estimated error range. The estimated error just comes from the Kalman filter's  $\mathbf{P}$  matrix. So let's filter some noisy data, take the filtered output as the new noisy input to the filter, and see what the result is. In other words,  $\mathbf{x}$  will supply the  $\mathbf{z}$  input, and  $\mathbf{P}$  will supply the measurement covariance  $\mathbf{R}$ . To exaggerate the effects somewhat to make them more obvious I will plot the effects of doing this one, and then a second time. The second iteration doesn't make any 'sense' (no one would try that), it just helps me illustrate a point. First, the code and plots.

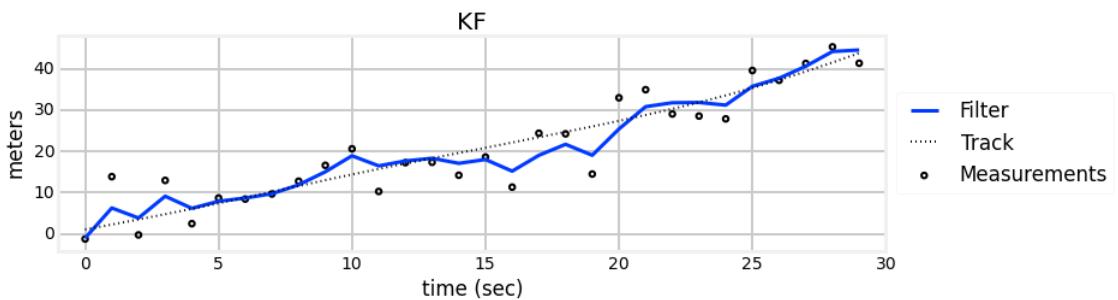
```
In [47]: np.random.seed(124)
R=5.
xs, zs = simulate_acc_system(R=R, Q=Q, count=30)

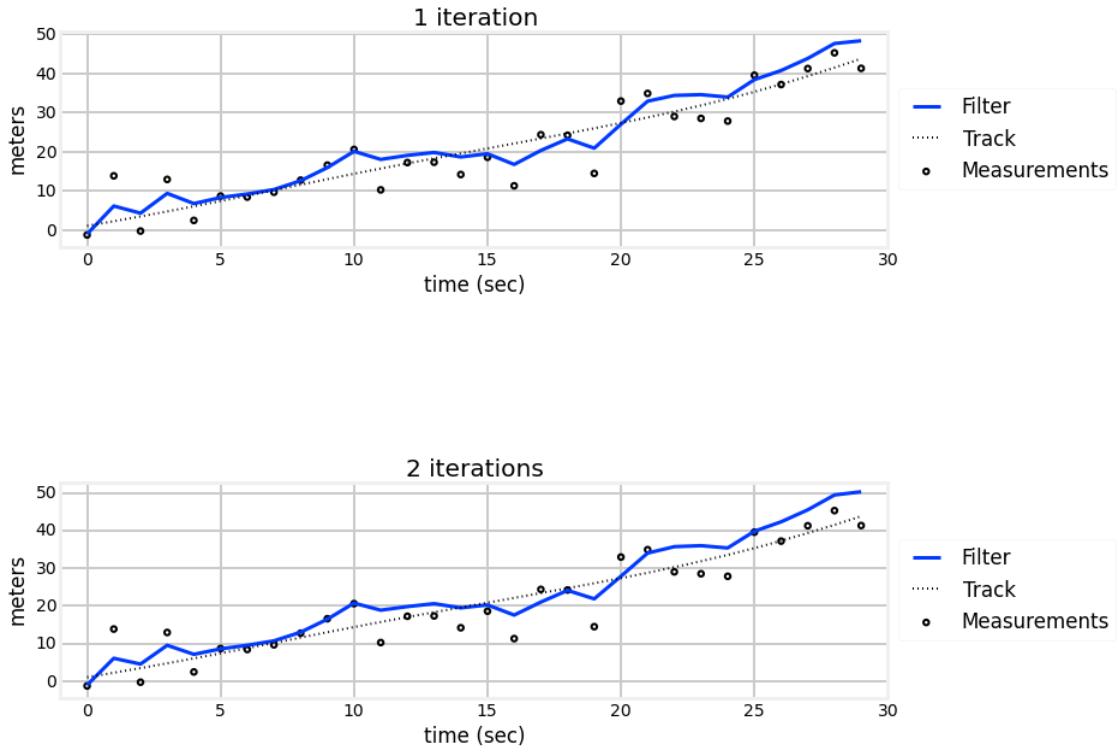
kf0 = SecondOrderKF(R, Q, dt=1)
kf1 = SecondOrderKF(R, Q, dt=1)
kf2 = SecondOrderKF(R, Q, dt=1)

# Filter measurements
fxs0, ps0 = filter_data(kf0, zs)

# filter twice more, using P as error
fxs1, ps1, _, _ = kf1.batch_filter(fxs0[:, 0], ps0[:, 0])
fxs2, _, _, _ = kf2.batch_filter(fxs1[:, 0], ps1[:, 0, 0])

plot_kf_output(xs, fxs0, zs, 'KF', False)
plot_kf_output(xs, fxs1, zs, '1 iteration', False)
plot_kf_output(xs, fxs2, zs, '2 iterations', False)
R,Q
```



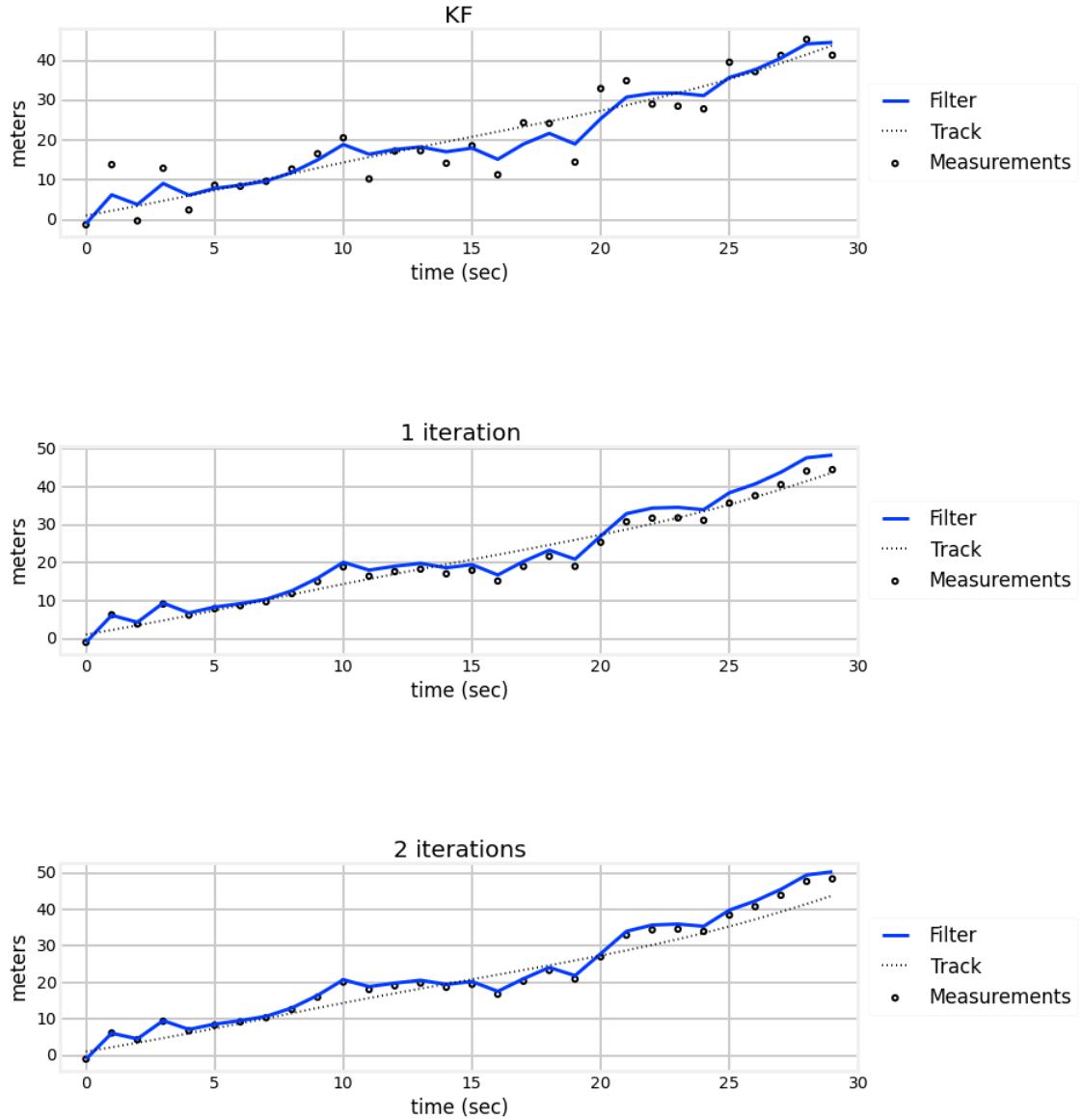


Out [47]: (5.0, 0.02)

We see that the filtered output of the reprocessed signal is smoother, but it also diverges from the track. What is happening? Recall that the Kalman filter requires that the signal not be time correlated. However the output of the Kalman filter is time correlated because it incorporates all previous measurements into its estimate for this time period. So look at the last graph, for 2 iterations. The measurements start with several peaks that are larger than the track. This is ‘remembered’ (that is vague terminology, but I am trying to avoid the math) by the filter, and it has started to compute that the object is above the track. Later, at around 13 seconds we have a period where the measurements all happen to be below the track. This also gets incorporated into the memory of the filter, and the iterated output diverges far below the track.

Now let’s look at this in a different way. The iterated output is not using  $\mathbf{z}$  as the measurement, but the output of the previous Kalman filter estimate. So I will plot the output of the filter against the previous filter’s output.

```
In [48]: plot_kf_output(xs, fxs0, zs, title='KF', aspect_equal=False)
plot_kf_output(xs, fxs1, fxs0[:, 0], '1 iteration', False)
plot_kf_output(xs, fxs2, fxs1[:, 0], '2 iterations', False)
```



I hope the problem with this approach is now apparent. In the bottom graph we can see that the KF is tracking the imperfect estimates of the previous filter, and incorporating delay into the signal as well due to the memory of the previous measurements being incorporated into the signal.

### 8.6.2 Exercise: Prove that the position sensor improves the filter

Devise a way to prove that fusing the position sensor and wheel measurements yields a better result than using the wheel alone.

#### Solution 1

Force the Kalman filter to disregard the position sensor measurement by setting the measurement noise for the position sensor to a near infinite value. Re-run the filter and observe the standard deviation of the

residual.

```
In [49]: fusion_test(1.5, 3.0, do_plot=False)
fusion_test(1.5, 1.e40, do_plot=False)

fusion std: 0.391
fusion std: 0.438
```

Here we can see the error in the filter where the position sensor measurement is almost entirely ignored is greater than that where it is used.

## Solution 2

This is more work, but we can write a Kalman filter that only takes one measurement.

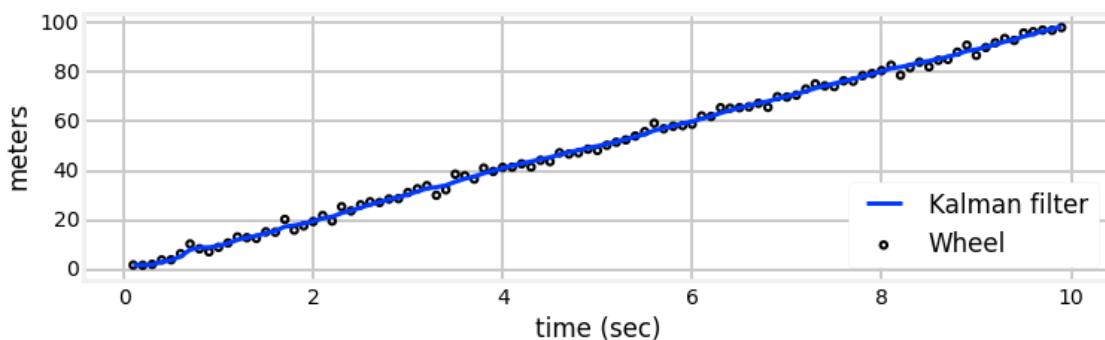
```
In [50]: dt = 0.1
wheel_sigma = 1.5
kf = KalmanFilter(dim_x=2, dim_z=1)
kf.F = array([[1., dt], [0., 1.]])
kf.H = array([[1., 0.]])
kf.x = array([[0.], [1.]])
kf.Q *= 0.01
kf.P *= 100
kf.R[0, 0] = wheel_sigma**2

random.seed(1123)
nom = range(1, 100)
zs = np.array([i + randn()*wheel_sigma for i in nom])
xs, _, _, _ = kf.batch_filter(zs)
ts = np.arange(0.1, 10, .1)

res = nom - xs[:, 0]
print('std: {:.3f}'.format(np.std(res)))

bp.plot_filter(ts, xs[:, 0], label='Kalman filter')
bp.plot_measurements(ts, zs, label='Wheel')
bp.set_labels(x='time (sec)', y='meters')
plt.legend(loc=4);
```

std: 40.386



On this run I got a standard deviation of 0.523 vs the value of 0.391 for the fused measurements.

## 8.7 Nonstationary Processes

So far we have assumed that the various matrices in the Kalman filter are stationary - nonchanging over time. For example, in the robot tracking section we assumed that  $\Delta t = 1.0$  seconds, and designed the state transition matrix to be

$$\mathbf{F} = \begin{bmatrix} 1 & \Delta t & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & \Delta t \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

But what if our data rate changes in some unpredictable manner? Or what if we have two sensors, each running at a different rate? What if the error of the measurement changes?

Handling this is easy; you just alter the Kalman filter matrices to reflect the current situation. Let's go back to our dog tracking problem and assume that the data input is somewhat sporadic. For this problem we designed

$$\bar{\mathbf{x}} = \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

$$\mathbf{F} = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix}$$

and set the Kalman filter variable  $\mathbf{F}$  during initialization like so:

```
dt = 0.1
kf.F = np.array([[1, dt],
                 [0, 1]])
```

How would we handle  $\Delta t$  changing for each measurement? It's easy - just modify the relevant matrices. In this case  $\mathbf{F}$  is variant, so we will need to update this inside the update/predict loop.  $\mathbf{Q}$  is also dependent on time, so it must be assigned during each loop as well. Here is an example of how we might code this:

```
In [51]: kf = KalmanFilter(dim_x=2, dim_z=1)
        kf.x = array([0., 1.])
        kf.H = array([[1, 0]])
        kf.P = np.eye(2) * 50
        kf.R = 1.
        Q = 0.02

        # measurement tuple: (value, time)
        zs = [(1., 1.), (2., 1.1), (3., 0.9), (4.1, 1.23), (5.01, 0.97)]
        for z, dt in zs:
            kf.F = array([[1, dt],
                         [0, 1]])
            kf.Q = Q_discrete_white_noise(dim=2, dt=dt, var=Q)
            kf.predict()
            kf.update(z)
            print(kf.x)
```

```
[ 1.  1.]
[ 2.003  0.916]
[ 2.963  0.997]
[ 4.124  0.97 ]
[ 5.032  0.959]
```

### 8.7.1 Sensor fusion: Different Data Rates

It is rare that two different sensor classes output data at the same rate. Assume that the position sensor produces an update at 1 Hz, and the wheel updates at 4 Hz.

We can do this by setting the data rate to 0.25 seconds, which is 4 Hz. As we loop, on every iteration we call `update()` with only the wheel measurement. Then, every fourth time we will call `update()` with both the wheel and PS measurements.

This means that we vary the amount of data in the `z` parameter. The matrices associated with the measurement are **H** and **R**. In the code above we designed **H** to be

$$\mathbf{H} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix}$$

to account for the two measurements of position. When only the wheel reading is available, we must set

$$\mathbf{H} = [1 \ 0].$$

The matrix **R** specifies the noise in each measurement. In the code above we set

$$\mathbf{R} = \begin{bmatrix} \sigma_{wheel}^2 & 0 \\ 0 & \sigma_{ps}^2 \end{bmatrix}$$

When only the wheel measurement is available, we must set

$$\mathbf{R} = [\sigma_{wheel}^2]$$

The two matrices that incorporate time are **F** and **Q**. For example,

$$\mathbf{F} = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix}.$$

Since the wheel and position sensor reading coincide once every 4 readings we can just set  $\Delta t = 0.25$  and not modify it while filtering. If the readings did not coincide in each iteration you would have to calculate how much time has passed since the last predict, compute a new **F** and **Q**, and then call `predict()` so the filter could make a correct prediction based on the time step required.

```
In [52]: def fusion_test(wheel_sigma, ps_sigma, do_plot=True):
    dt = 0.25
    kf = KalmanFilter(dim_x=2, dim_z=2)

    kf.F = array([[1., dt], [0., 1.]])
    kf.H = array([[1., 0.], [1., 0.]])
    kf.x = array([0., 1.])
    kf.Q *= array([[[(dt**3)/3, (dt**2)/2],
                   [(dt**2)/2, dt]]]) * 0.02
    kf.P *= 100
```

```

kf.R[0, 0] = wheel_sigma**2
kf.R[1, 1] = ps_sigma**2

random.seed(1123)
xs, zs_wheel, zs_ps, nom = [], [], [], []
for i in range(1, 101):
    if i % 4 == 0:
        m0 = i + randn()*wheel_sigma
        m1 = i + randn()*ps_sigma
        z = array([[m0], [m1]])
        kf.H = array([[1., 0.], [1., 0.]])
        R = np.eye(2)
        R[0, 0] = wheel_sigma**2
        R[1, 1] = ps_sigma**2

        zs_wheel.append(m0)
        zs_ps.append((i*dt, m1))
    else:
        m0 = i + randn()*wheel_sigma
        z = array([m0])
        kf.H = array([[1., 0.]])
        R = np.eye(1) * wheel_sigma**2
        zs_wheel.append(m0)
    kf.predict()
    kf.update(z, R)

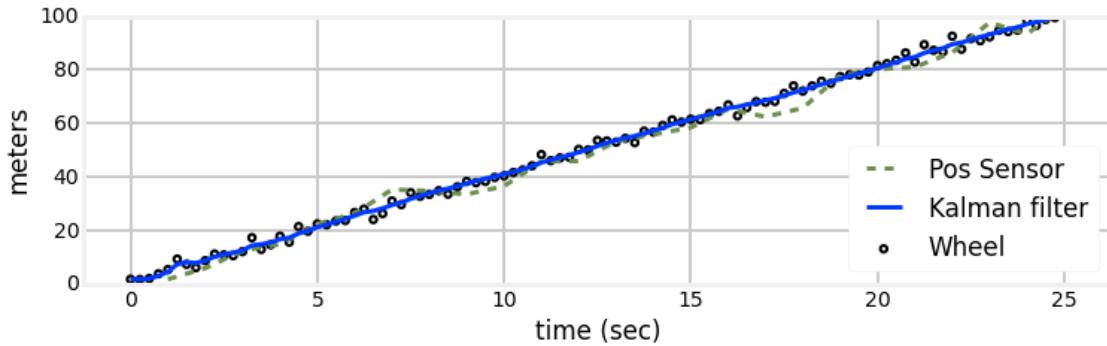
    xs.append(kf.x.T[0])
    nom.append(i)

xs = asarray(xs)
nom = asarray(nom)
res = nom - xs[:, 0]
print('fusion std: {:.3f}'.format(np.std(res)))
if do_plot:
    ts = np.arange(0, 25, .25)
    bp.plot_measurements(ts, zs_wheel, label='Wheel')
    plt.plot(*zip(*zs_ps), ls='--', label='Pos Sensor')
    bp.plot_filter(ts, xs[:, 0], label='Kalman filter')
    plt.legend(loc=4)
    plt.ylim(0, 100)
    bp.set_labels(x='time (sec)', y='meters')

fusion_test(1.5, 3.0);

fusion std: 0.452

```



We can see from the standard deviation that the performance is a bit worse than when the PS and wheel were measured in every update, but better than the wheel alone.

The code is fairly straightforward. The `update()` method optionally takes  $R$  as an argument, and I chose to do that rather than alter `KalmanFilter.R`, mostly to show that it is possible. Either way is fine. I modified `KalmanFilter.H` on each update depending on whether there are 1 or 2 measurements available. The only other difficulty was storing the wheel and PS measurements in two different arrays because there are a different number of measurements for each.

## 8.8 Tracking a Ball

Now let's turn our attention to a situation where the physics of the object that we are tracking is constrained. A ball thrown in a vacuum must obey Newtonian laws. In a constant gravitational field it will travel in a parabola. I will assume you are familiar with the derivation of the formula:

$$y = \frac{g}{2}t^2 + v_{y0}t + y_0$$

$$x = v_{x0}t + x_0$$

where  $g$  is the gravitational constant,  $t$  is time,  $v_{x0}$  and  $v_{y0}$  are the initial velocities in the  $x$  and  $y$  plane. If the ball is thrown with an initial velocity of  $v$  at angle  $\theta$  above the horizon, we can compute  $v_{x0}$  and  $v_{y0}$  as

$$v_{x0} = v \cos \theta$$

$$v_{y0} = v \sin \theta$$

Because we don't have real data we will start by writing a simulator for a ball. As always, we add a noise term independent of time so we can simulate noise sensors.

```
In [53]: from math import radians, sin, cos
import math

def rk4(y, x, dx, f):
    """computes 4th order Runge-Kutta for dy/dx.
    y is the initial value for y
    x is the initial value for x
    dx is the difference in x (e.g. the time step)
    f is a callable function (y, x) that you supply to
    compute dy/dx for the specified values.
```

```

"""
k1 = dx * f(y, x)
k2 = dx * f(y + 0.5*k1, x + 0.5*dx)
k3 = dx * f(y + 0.5*k2, x + 0.5*dx)
k4 = dx * f(y + k3, x + dx)

return y + (k1 + 2*k2 + 2*k3 + k4) / 6.

def fx(x,t):
    return fx.vel

def fy(y,t):
    return fy.vel - 9.8*t

class BallTrajectory2D(object):
    def __init__(self, x0, y0, velocity,
                 theta_deg=0.,
                 g=9.8,
                 noise=[0.0, 0.0]):
        self.x = x0
        self.y = y0
        self.t = 0
        theta = math.radians(theta_deg)
        fx.vel = math.cos(theta) * velocity
        fy.vel = math.sin(theta) * velocity
        self.g = g
        self.noise = noise

    def step(self, dt):
        self.x = rk4(self.x, self.t, dt, fx)
        self.y = rk4(self.y, self.t, dt, fy)
        self.t += dt
        return (self.x + randn()*self.noise[0],
                self.y + randn()*self.noise[1])

```

So to create a trajectory starting at  $(0, 15)$  with a velocity of  $100 \frac{m}{s}$  and an angle of  $60^\circ$  we would write:

```
traj = BallTrajectory2D(x0=0, y0=15, velocity=100, theta_deg=60)
```

and then call `traj.step(t)` for each time step. Let's test this

```
In [54]: def test_ball_vacuum(noise):
    y = 15
    x = 0
    ball = BallTrajectory2D(x0=x, y0=y,
                           theta_deg=60., velocity=100.,
                           noise=noise)
    t = 0
    dt = 0.25
    while y >= 0:
        x, y = ball.step(dt)
```

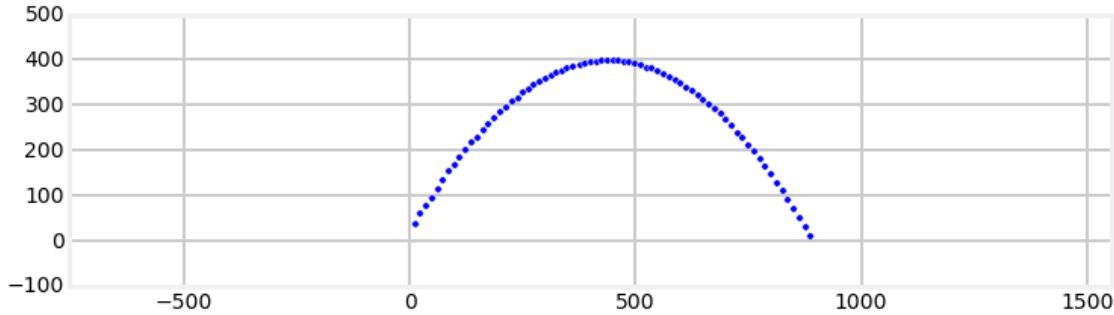
```

t += dt
if y >= 0:
    plt.scatter(x, y)

plt.axis('equal');

#test_ball_vacuum([0, 0]) # plot ideal ball position
test_ball_vacuum([1, 1]) # plot with noise

```



This looks reasonable, so let's continue (exercise for the reader: validate this simulation more robustly).

### 8.8.1 Step 1: Choose the State Variables

We might think to use the same state variables as used for tracking the dog. However, this will not work. Recall that the Kalman filter state transition must be written as  $\bar{\mathbf{x}} = \mathbf{F}\mathbf{x} + \mathbf{B}\mathbf{u}$ , which means we must calculate the current state from the previous state. Our assumption is that the ball is traveling in a vacuum, so the velocity in  $x$  is a constant, and the acceleration in  $y$  is solely due to the gravitational constant  $g$ . We can discretize the Newtonian equations using the well known Euler method in terms of  $\Delta t$  are:

$$\begin{aligned}x_t &= x_{t-1} + v_{x(t-1)}\Delta t \\v_{xt} &= v_{x(t-1)} \\y_t &= y_{t-1} + v_{y(t-1)}\Delta t \\v_{yt} &= -g\Delta t + v_{y(t-1)}\end{aligned}$$

sidebar: Euler's method integrates a differential equation stepwise by assuming the slope (derivative) is constant at time  $t$ . In this case the derivative of the position is velocity. At each time step  $\Delta t$  we assume a constant velocity, compute the new position, and then update the velocity for the next time step. There are more accurate methods, such as Runge-Kutta available to us, but because we are updating the state with a measurement in each step Euler's method is very accurate. If you need to use Runge-Kutta you will have to write your own `predict()` function which computes the state transition for  $\mathbf{x}$ , and then uses the normal Kalman filter equation  $\bar{\mathbf{P}} = \mathbf{F}\mathbf{P}\mathbf{F}^T + \mathbf{Q}$  to update the covariance matrix.

This implies that we need to incorporate acceleration for  $y$  into the Kalman filter, but not for  $x$ . This suggests the following state variables.

$$\mathbf{x} = \begin{bmatrix} x \\ \dot{x} \\ y \\ \dot{y} \\ \ddot{y} \end{bmatrix}$$

However, the acceleration is due to gravity, which is a constant. Instead of asking the Kalman filter to track a constant we can treat gravity as what it really is - a control input. In other words, gravity is a force that alters the behavior of the system in a known way, and it is applied throughout the flight of the ball.

The equation for the state prediction is  $\bar{\mathbf{x}} = \mathbf{F}\mathbf{x} + \mathbf{B}\mathbf{u}$ .  $\mathbf{F}\mathbf{x}$  is the familiar state transition function which we will use to model the position and velocity of the ball. The vector  $\mathbf{u}$  lets you specify a control input into the filter. For a car the control input will be things such as the amount the accelerator and brake are pressed, the position of the steering wheel, and so on. For our ball the control input will be gravity. The matrix  $\mathbf{B}$  models how the control inputs affect the behavior of the system. Again, for a car  $\mathbf{B}$  will convert the inputs of the brake and accelerator into changes of velocity, and the input of the steering wheel into a different position and heading. For our ball tracking problem it will compute the velocity change due to gravity. We will go into the details of that in step 3. For now, we design the state variable to be

$$\mathbf{x} = \begin{bmatrix} x \\ \dot{x} \\ y \\ \dot{y} \end{bmatrix}$$

### 8.8.2 Step 2: Design State Transition Function

Our next step is to design the state transition function. Recall that the state transition function is implemented as a matrix  $\mathbf{F}$  that we multiply with the previous state of our system to get the next state  $\mathbf{x}' = \mathbf{F}\mathbf{x}$ .

I will not belabor this as it is very similar to the 1-D case we did in the previous chapter. Our state equations for position and velocity would be:

$$\begin{aligned} x' &= (1 * x) + (\Delta t * v_x) + (0 * y) + (0 * v_y) \\ v_x &= (0 * x) + (1 * v_x) + (0 * y) + (0 * v_y) \\ y' &= (0 * x) + (0 * v_x) + (1 * y) + (\Delta t * v_y) \\ v_y &= (0 * x) + (0 * v_x) + (0 * y) + (1 * v_y) \end{aligned}$$

Note that none of the terms include  $g$ , the gravitational constant. As I explained in the previous function we will account for gravity using the control input of the Kalman filter. In matrix form we write this as:

$$\mathbf{F} = \begin{bmatrix} 1 & \Delta t & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & \Delta t \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

### 8.8.3 Step 3: Design the Control Input Function

We will use the control input to account for the force of gravity. The term  $\mathbf{B}\mathbf{u}$  is added to  $\bar{\mathbf{x}}$  to account for how much  $\bar{\mathbf{x}}$  changes due to gravity. We can say that  $\mathbf{B}\mathbf{u}$  contains  $[\Delta x_g \quad \Delta \dot{x}_g \quad \Delta y_g \quad \Delta \dot{y}_g]^\top$ .

If we look at the discretized equations we see that gravity only affect the velocity for  $y$ .

$$\begin{aligned}x_t &= x_{t-1} + v_{x(t-1)}\Delta t \\v_{xt} &= vx_{t-1} \\y_t &= y_{t-1} + v_{y(t-1)}\Delta t \\v_{yt} &= -g\Delta t + v_{y(t-1)}\end{aligned}$$

Therefore we want the product  $\mathbf{B}\mathbf{u}$  to equal  $[0 \ 0 \ 0 \ -g\Delta t]^\top$ . In some sense it is arbitrary how we define  $\mathbf{B}$  and  $\mathbf{u}$  so long as multiplying them together yields this result. For example, we could define  $\mathbf{B} = 1$  and  $\mathbf{u} = [0 \ 0 \ 0 \ -g\Delta t]^\top$ . But this doesn't really fit with our definitions for  $\mathbf{B}$  and  $\mathbf{u}$ , where  $\mathbf{u}$  is the control input, and  $\mathbf{B}$  is the control function. The control input is  $-g$  for the velocity of  $y$ . So this is one possible definition.

$$\mathbf{B} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \Delta t \end{bmatrix}, \mathbf{u} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ -g \end{bmatrix}$$

To me this seems a bit excessive; I would suggest we might want  $\mathbf{u}$  to contain the control input for the two dimensions  $x$  and  $y$ , which suggests

$$\mathbf{B} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & \Delta t \end{bmatrix}, \mathbf{u} = \begin{bmatrix} 0 \\ -g \end{bmatrix}$$

You might prefer to only provide control inputs that actually exist, and there is no control input for  $x$ , so we arrive at

$$\mathbf{B} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \Delta t \end{bmatrix}, \mathbf{u} = [-g]$$

I've seen people use

$$\mathbf{B} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \mathbf{u} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ -g\Delta t \end{bmatrix}$$

While this does produce the correct result, I am resistant to putting time into  $\mathbf{u}$  as time is not a control input, it is what we use to convert the control input into a change in state, which is the job of  $\mathbf{B}$ .

#### 8.8.4 Step 4: Design the Measurement Function

The measurement function defines how we go from the state variables to the measurements using the equation  $\mathbf{z} = \mathbf{H}\mathbf{x}$ . We will assume that we have a sensor that provides us with the position of the ball in (x,y), but cannot measure velocities or accelerations. Therefore our function must be:

$$\begin{bmatrix} z_x \\ z_y \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} * \begin{bmatrix} x \\ \dot{x} \\ y \\ \dot{y} \end{bmatrix}$$

where

$$\mathbf{H} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

### 8.8.5 Step 5: Design the Measurement Noise Matrix

As with the robot, we will assume that the error is independent in  $x$  and  $y$ . In this case we will start by assuming that the measurement error in  $x$  and  $y$  are 0.5 meters. Hence,

$$\mathbf{R} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

### 8.8.6 Step 6: Design the Process Noise Matrix

We are assuming a ball moving in a vacuum, so there should be no process noise. We have 4 state variables, so we need a  $4 \times 4$  covariance matrix:

$$\mathbf{Q} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

### 8.8.7 Step 7: Design the Initial Conditions

We already performed this step when we tested the state transition function. Recall that we computed the initial velocity for  $x$  and  $y$  using trigonometry, and set the value of  $\mathbf{x}$  with:

```
omega = radians(omega)
vx = cos(omega) * v0
vy = sin(omega) * v0

f1.x = np.array([[x, vx, y, vy]]).T
```

With all the steps done we are ready to implement our filter and test it. First, the implementation:

In [55]: `from math import sin,cos,radians`

```
def ball_kf(x, y, omega, v0, dt, r=0.5, q=0.):
    kf = KalmanFilter(dim_x=4, dim_z=2, dim_u=1)

    ay = .5*dt**2
    kf.F = np.array([[1, dt, 0, 0],      # x     = x0 + dx*dt
                    [0, 1, 0, 0],      # dx   = dx0
                    [0, 0, 1, dt],     # y     = y0 + dy*dt
                    [0, 0, 0, 1]])    # dy   = dy0

    kf.H = np.array([[1, 0, 0, 0],
                    [0, 0, 1, 0]])
```

```

kf.B = np.array([[0, 0, 0, dt]]).T
kf.R *= r
kf.Q *= q

omega = radians(omega)
vx = cos(omega) * v0
vy = sin(omega) * v0
kf.x = np.array([[x, vx, y, vy]]).T
return kf

```

Now we will test the filter by generating measurements for the ball using the ball simulation class.

```

In [56]: def track_ball_vacuum(dt):
    x, y = 0., 1.
    theta = 35. # launch angle
    v0 = 80.
    g = 9.8 # gravitational constant
    ball = BallTrajectory2D(x0=x, y0=y, theta_deg=theta, velocity=v0,
                           noise=[.2, .2])
    kf = ball_kf(x, y, theta, v0, dt)

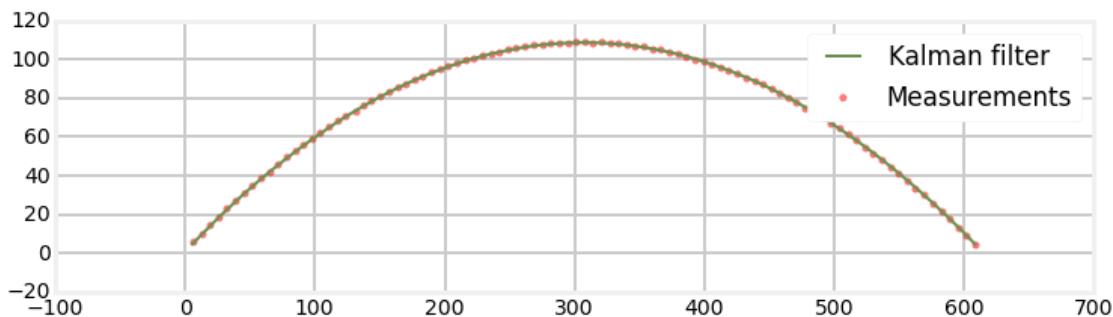
    t = 0
    xs, ys = [], []
    while kf.x[2] > 0:
        t += dt
        x, y = ball.step(dt)
        z = np.array([[x, y]]).T

        kf.update(z)
        xs.append(kf.x[0])
        ys.append(kf.x[2])
        kf.predict(u=-g)
        p1 = plt.scatter(x, y, color='r', marker='.', s=75, alpha=0.5)

    p2, = plt.plot(xs, ys, lw=2)
    plt.legend([p2, p1], ['Kalman filter', 'Measurements'],
               scatterpoints=1)

track_ball_vacuum(dt=1./10)

```



We see that the Kalman filter reasonably tracks the ball. However, as already explained, this is a trivial example because we have no process noise. We can predict trajectories in a vacuum with arbitrary precision; using a Kalman filter in this example is a needless complication. A least squares curve fit would give identical results.

## 8.9 Tracking a Ball in Air

For this problem we assume that we are tracking a ball traveling through the Earth's atmosphere. The path of the ball is influenced by wind, drag, and the rotation of the ball. We will assume that our sensor is a camera; code that we will not implement will perform some type of image processing to detect the position of the ball. This is typically called blob detection in computer vision. However, image processing code is not perfect; in any given frame it is possible to either detect no blob or to detect spurious blobs that do not correspond to the ball. Finally, we will not assume that we know the starting position, angle, or rotation of the ball; the tracking code will have to initiate tracking based on the measurements that are provided. The main simplification that we are making here is a 2D world; we assume that the ball is always traveling orthogonal to the plane of the camera's sensor. We have to make that simplification at this point because we have not yet discussed how we might extract 3D information from a camera, which necessarily provides only 2D data.

### 8.9.1 Implementing Air Drag

Our first step is to implement the math for a ball moving through air. There are several treatments available. A robust solution takes into account issues such as ball roughness (which affects drag non-linearly depending on velocity), the Magnus effect (spin causes one side of the ball to have higher velocity relative to the air vs the opposite side, so the coefficient of drag differs on opposite sides), the effect of lift, humidity, air density, and so on. I assume the reader is not interested in the details of ball physics, and so will restrict this treatment to the effect of air drag on a non-spinning baseball. I will use the math developed by Nicholas Giordano and Hisao Nakanishi in Computational Physics [1997]. This treatment does not take all the factors into account. The most detailed treatment is by Alan Nathan on his website at <http://baseball.physics.illinois.edu/index.html>. I use his math in my own work in computer vision, but I do not want to get distracted by a more complicated model.

**Important:** Before I continue, let me point out that you will not have to understand this next piece of physics to proceed with the Kalman filter. My goal is to create a reasonably accurate behavior of a baseball in the real world, so that we can test how our Kalman filter performs with real-world behavior. In real world applications it is usually impossible to completely model the physics of a real world system, and we make do with a process model that incorporates the large scale behaviors. We then tune the measurement noise and process noise until the filter works well with our data. There is a real risk to this; it is easy to finely tune a Kalman filter so it works perfectly with your test data, but performs badly when presented with slightly different data. This is perhaps the hardest part of designing a Kalman filter, and why it gets referred to with terms such as 'black art'.

I dislike books that implement things without explanation, so I will now develop the physics for a ball moving through air. Move on past the implementation of the simulation if you are not interested.

A ball moving through air encounters wind resistance. This imparts a force on the ball, called drag, which alters the flight of the ball. In Giordano this is denoted as

$$F_{drag} = -B_2 v^2$$

where  $B_2$  is a coefficient derived experimentally, and  $v$  is the velocity of the object.  $F_{drag}$  can be factored into  $x$  and  $y$  components with

$$F_{drag,x} = -B_2vv_x F_{drag,y} = -B_2vv_y$$

If  $m$  is the mass of the ball, we can use  $F = ma$  to compute the acceleration as

$$a_x = -\frac{B_2}{m}vv_x a_y = -\frac{B_2}{m}vv_y$$

Giordano provides the following function for  $\frac{B_2}{m}$ , which takes air density, the cross section of a baseball, and its roughness into account. Understand that this is an approximation based on wind tunnel tests and several simplifying assumptions. It is in SI units: velocity is in meters/sec and time is in seconds.

$$\frac{B_2}{m} = 0.0039 + \frac{0.0058}{1 + \exp[(v - 35)/5]}$$

Starting with this Euler discretization of the ball path in a vacuum:

$$\begin{aligned} x &= v_x \Delta t \\ y &= v_y \Delta t \\ v_x &= v_x \\ v_y &= v_y - 9.8 \Delta t \end{aligned}$$

We can incorporate this force (acceleration) into our equations by incorporating  $accel * \Delta t$  into the velocity update equations. We should subtract this component because drag will reduce the velocity. The code to do this is quite straightforward, we just need to break out the Force into  $x$  and  $y$  components.

I will not belabor this issue further because computational physics is beyond the scope of this book. Recognize that a higher fidelity simulation would require incorporating things like altitude, temperature, ball spin, and several other factors. The aforementioned work by Alan Nathan covers this if you are interested. My intent here is to impart some real-world behavior into our simulation to test how our simpler prediction model used by the Kalman filter reacts to this behavior. Your process model will never exactly model what happens in the world, and a large factor in designing a good Kalman filter is carefully testing how it performs against real world data.

The code below computes the behavior of a baseball in air, at sea level, in the presence of wind. I plot the same initial hit with no wind, and then with a tail wind at 10 mph. Baseball statistics are universally done in US units, and we will follow suit here ([http://en.wikipedia.org/wiki/United\\_States\\_customary\\_units](http://en.wikipedia.org/wiki/United_States_customary_units)). Note that the velocity of 110 mph is a typical exit speed for a baseball for a home run hit.

```
In [57]: from math import sqrt, exp, cos, sin, radians

def mph_to_mps(x):
    return x * .447

def drag_force(velocity):
    """ Returns the force on a baseball due to air drag at
    the specified velocity. Units are SI"""
    return velocity * (0.0039 + 0.0058 /
        (1. + exp((velocity-35.)/5.)))

v = mph_to_mps(110.)
x, y = 0., 1.
dt = .1
theta = radians(35)
```

```

def solve(x, y, vel, v_wind, launch_angle):
    xs = []
    ys = []
    v_x = vel*cos(launch_angle)
    v_y = vel*sin(launch_angle)
    while y >= 0:
        # Euler equations for x and y
        x += v_x*dt
        y += v_y*dt

        # force due to air drag
        velocity = sqrt ((v_x-v_wind)**2 + v_y**2)
        F = drag_force(velocity)

        # euler's equations for vx and vy
        v_x = v_x - F*(v_x-v_wind)*dt
        v_y = v_y - 9.8*dt - F*v_y*dt

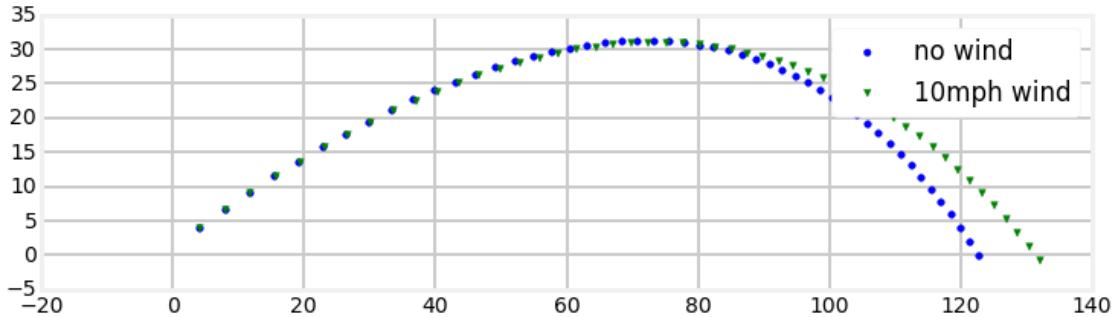
    xs.append(x)
    ys.append(y)

return xs, ys

x, y = solve(x=0, y=1, vel=v, v_wind=0, launch_angle=theta)
p1 = plt.scatter(x, y, color='blue', label='no wind')

wind = mph_to_mps(10)
x, y = solve(x=0, y=1, vel=v, v_wind=wind, launch_angle=theta)
p2 = plt.scatter(x, y, color='green', marker="v",
                  label='10mph wind')
plt.legend(scatterpoints=1);

```



We can easily see the difference between the trajectory in a vacuum and in the air. I used the same initial velocity and launch angle in the ball in a vacuum section above. We computed that the ball in a vacuum would travel over 240 meters (nearly 800 ft). In the air, the distance is just over 120 meters, or roughly 400 ft. 400ft is a realistic distance for a well hit home run ball, so we can be confident that our simulation is reasonably accurate.

Without further ado we will create a ball simulation that uses the math above to create a more realistic ball trajectory. I will note that the nonlinear behavior of drag means that there is no analytic solution to

the ball position at any point in time, so we need to compute the position step-wise. I use Euler's method to propagate the solution; use of a more accurate technique such as Runge-Kutta is left as an exercise for the reader. That modest complication is unnecessary for what we are doing because the accuracy difference between the techniques will be small for the time steps we will be using.

```
In [58]: from math import radians, sin, cos, sqrt, exp

class BaseballPath(object):
    def __init__(self, x0, y0, launch_angle_deg, velocity_ms,
                 noise=(1.0, 1.0)):
        """
        Create 2D baseball path object
        (x = distance from start point in ground plane,
         y=height above ground)

        x0,y0           initial position
        launch_angle_deg angle ball is travelling respective to
                           ground plane
        velocity_ms     speeed of ball in meters/second
        noise           amount of noise to add to each position
                       in (x, y)
        """

        omega = radians(launch_angle_deg)
        self.v_x = velocity_ms * cos(omega)
        self.v_y = velocity_ms * sin(omega)

        self.x = x0
        self.y = y0
        self.noise = noise

    def drag_force(self, velocity):
        """
        Returns the force on a baseball due to air drag at
        the specified velocity. Units are SI
        """
        B_m = 0.0039 + 0.0058 / (1. + exp((velocity-35.)/5.))
        return B_m * velocity

    def update(self, dt, vel_wind=0.):
        """
        compute the ball position based on the specified time
        step and wind velocity. Returns (x, y) position tuple.
        """

        # Euler equations for x and y
        self.x += self.v_x*dt
        self.y += self.v_y*dt

        # force due to air drag
        v_x_wind = self.v_x - vel_wind
        v = sqrt(v_x_wind**2 + self.v_y**2)
        F = self.drag_force(v)

        # Euler's equations for velocity
```

```

    self.v_x = self.v_x - F*v_x_wind*dt
    self.v_y = self.v_y - 9.81*dt - F*self.v_y*dt

    return (self.x + randn()*self.noise[0],
            self.y + randn()*self.noise[1])

```

Now we can test the Kalman filter against measurements created by this model.

In [59]: `x, y = 0, 1.`

```

theta = 35. # launch angle
v0 = 50.
dt = 1/10. # time step
g = 9.8

ball = BaseballPath(x0=x, y0=y, launch_angle_deg=theta,
                     velocity_ms=v0, noise=[.3,.3])
f1 = ball_kf(x, y, theta, v0, dt, r=1.)
f2 = ball_kf(x, y, theta, v0, dt, r=10.)
t = 0
xs, ys = [], []
xs2, ys2 = [], []

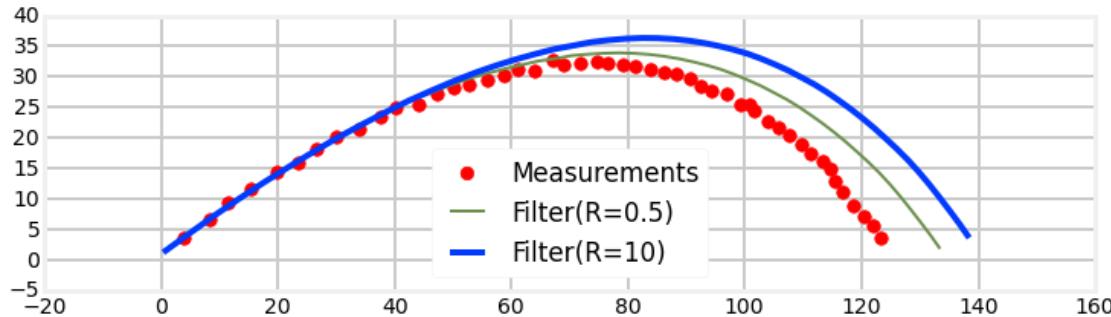
while f1.x[2] > 0:
    t += dt
    x, y = ball.update(dt)
    z = np.array([[x, y]]).T

    f1.update(z)
    f2.update(z)
    xs.append(f1.x[0])
    ys.append(f1.x[2])
    xs2.append(f2.x[0])
    ys2.append(f2.x[2])
    f1.predict(u=-g)
    f2.predict(u=-g)

p1 = plt.scatter(x, y, color='r', marker='o', s=75)

p2, = plt.plot(xs, ys, lw=2)
p3, = plt.plot(xs2, ys2, lw=4)
plt.legend([p1, p2, p3],
           ['Measurements', 'Filter(R=0.5)', 'Filter(R=10)'],
           loc='best', scatterpoints=1);

```



I have plotted the output of two different Kalman filter settings. The measurements are depicted as green circles, a Kalman filter with  $R=0.5$  as a thin green line, and a Kalman filter with  $R=10$  as a thick blue line. These  $R$  values are chosen merely to show the effect of measurement noise on the output, they are not intended to imply a correct design.

We can see that neither filter does very well. At first both track the measurements well, but as time continues they both diverge. This happens because the state model for air drag is nonlinear and the Kalman filter assumes that it is linear. If you recall our discussion about nonlinearity in the g-h filter chapter we showed why a g-h filter will always lag behind the acceleration of the system. We see the same thing here - the acceleration is negative, so the Kalman filter consistently overshoots the ball position. There is no way for the filter to catch up so long as the acceleration continues, so the filter will continue to diverge.

What can we do to improve this? The best approach is to perform the filtering with a nonlinear Kalman filter, and we will do this in subsequent chapters. However, there is also what I will call an ‘engineering’ solution to this problem as well. Our Kalman filter assumes that the ball is in a vacuum, and thus that there is no process noise. However, since the ball is in air the atmosphere imparts a force on the ball. We can think of this force as process noise. This is not a particularly rigorous thought; for one thing, this force is anything but Gaussian. Secondly, we can compute this force, so throwing our hands up and saying ‘it’s random’ will not lead to an optimal solution. But let’s see what happens if we follow this line of thought.

The following code implements the same Kalman filter as before, but with a non-zero process noise. I plot two examples, one with  $Q=.1$ , and one with  $Q=0.01$ .

```
In [60]: def plot_ball_with_q(q, r=1., noise=0.3):
    x, y = 0., 1.
    theta = 35. # launch angle
    v0 = 50.
    dt = 1/10. # time step
    g = 9.8

    ball = BaseballPath(x0=x,
                        y0=y,
                        launch_angle_deg=theta,
                        velocity_ms=v0,
                        noise=[noise,noise])
    f1 = ball_kf(x, y, theta, v0, dt, r=r, q=q)
    t = 0
    xs, ys = [], []

    while f1.x[2] > 0:
        t += dt
        x, y = ball.update(dt)
```

```

z = np.array([[x, y]]).T

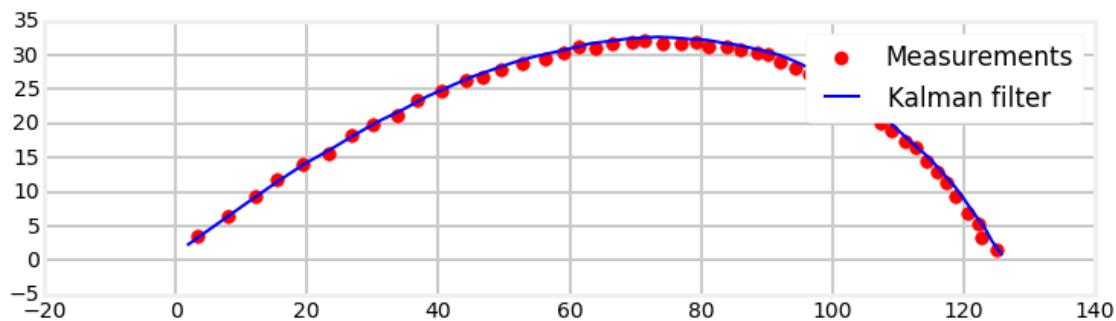
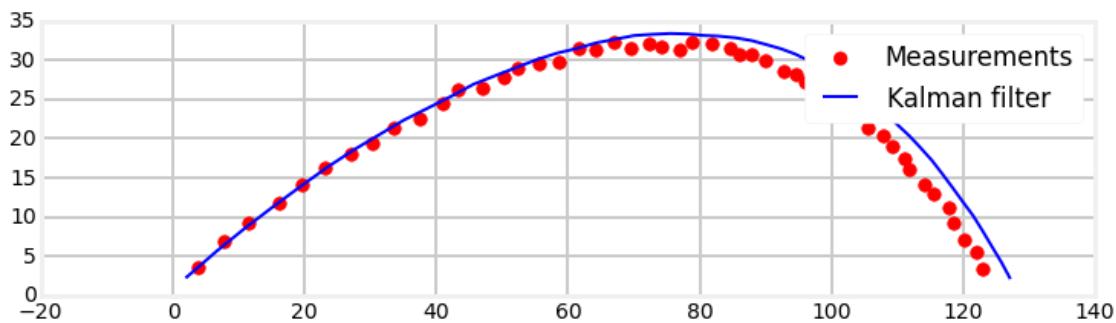
f1.update(z)
xs.append(f1.x[0])
ys.append(f1.x[2])
f1.predict(u=-g)

p1 = plt.scatter(x, y, color='r', marker='o', s=75)

p2, = plt.plot(xs, ys, lw=2, color='b')
plt.legend([p1, p2], ['Measurements', 'Kalman filter'])
plt.show()

plot_ball_with_q(0.01)
plot_ball_with_q(0.1)

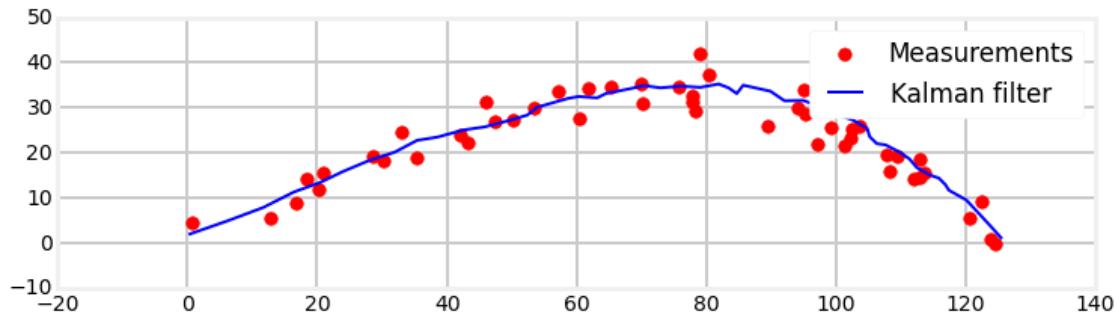
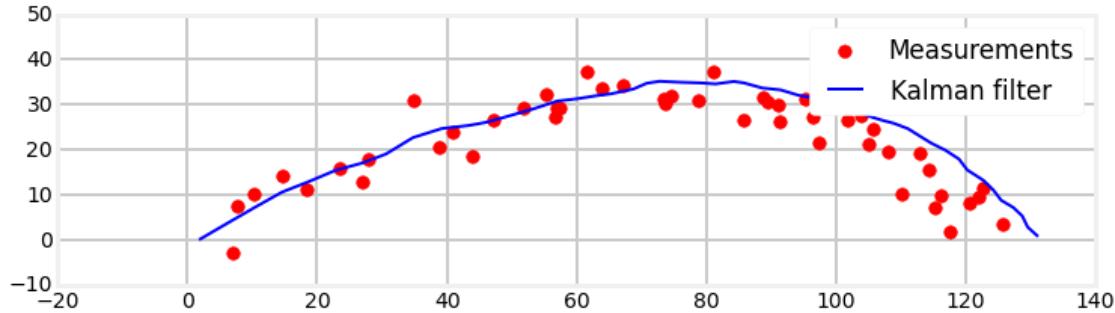
```



The second filter tracks the measurements fairly well. There appears to be a bit of lag, but very little.

Is this a good technique? Usually not, but it depends. Here the nonlinearity of the force on the ball is fairly constant and regular. Assume we are trying to track an automobile - the accelerations will vary as the car changes speeds and turns. When we make the process noise higher than the actual noise in the system the filter will opt to weigh the measurements higher. If you don't have a lot of noise in your measurements this might work for you. However, consider this next plot where I have increased the noise in the measurements.

In [61]: `plot_ball_with_q(0.01, r=3, noise=3.)`  
`plot_ball_with_q(0.1, r=3, noise=3.)`



This output is terrible. The filter has no choice but to give more weight to the measurements than the process (prediction step), but when the measurements are noisy the filter output will just track the noise. This inherent limitation of the linear Kalman filter is what lead to the development of nonlinear versions of the filter.

With that said, it is certainly possible to use the process noise to deal with small nonlinearities in your system. This is part of the ‘black art’ of Kalman filters. Our model of the sensors and of the system are never perfect. Sensors are non-Gaussian and our process model is never perfect. You can mask some of this by setting the measurement errors and process errors higher than their theoretically correct values, but the trade off is a non-optimal solution. Certainly it is better to be non-optimal than to have your Kalman filter diverge. However, as we can see in the graphs above, it is easy for the output of the filter to be very bad. It is also very common to run many simulations and tests and to end up with a filter that performs very well under those conditions. Then, when you use the filter on real data the conditions are slightly different and the filter ends up performing terribly.

For now we will set this problem aside, as we are clearly misapplying the Kalman filter in this example. We will revisit this problem in subsequent chapters to see the effect of using various nonlinear techniques. In some domains you will be able to get away with using a linear Kalman filter for a nonlinear problem, but usually you will have to use one or more of the techniques you will learn in the rest of this book.

## 8.10 References

- [1] Bar-Shalom, Yaakov, et al. Estimation with Applications to Tracking and Navigation. John Wiley & Sons, 2001.

# Chapter 9

## Nonlinear Filtering

### 9.1 Introduction

The Kalman filter that we have developed uses linear equations, and so the filter can only handle linear problems. But the world is nonlinear, and so the classic filter that we have been studying to this point can have very limited utility.

There can be nonlinearity in the process model. Suppose we wanted to track the motion of a weight on a spring, such as an automobile's suspension. The equation for the motion with  $m$  being the mass,  $k$  the spring constant, and  $c$  the damping force is

$$m \frac{d^2x}{dt^2} + c \frac{dx}{dt} + kx = 0$$

It is not possible to write a linear equation for this second order differential equation, and therefore we cannot design a Kalman filter using the theory that we have learned.

A second source of nonlinearity comes from the measurements. For example, radars measure the slant range to an object, and we are typically interested in the aircraft's position over the ground. We invoke Pythagoras and get the nonlinear equation:

$$x = \sqrt{\text{slant}^2 - \text{altitude}^2}$$

These facts were not lost on the early adopters of the Kalman filter. Soon after Dr. Kalman published his paper people began working on how to extend the Kalman filter for nonlinear problems.

It is almost true to state that the only equation anyone knows how to solve is  $\mathbf{Ax} = \mathbf{b}$ . We only really know how to do linear algebra. I can give you any linear set of equations and you can either solve it or prove that it has no solution.

Anyone with formal education in math or physics has spent years learning various analytic ways to solve integrals, differential equations and so on. Yet even trivial physical systems produce equations that cannot be solved analytically. I can take an equation that you are able to integrate, insert a log term, and render it insolvable. It is stuff like this that leads to jokes about physicists stating "assume a spherical cow on a frictionless surface in a vacuum...". Without making extreme simplifications most physical problems do not have analytic solutions.

How do we do things like model airflow over an aircraft in a computer, or predict weather, or track missiles with a Kalman filter? We retreat to what we know:  $\mathbf{Ax} = \mathbf{b}$ . We find some way to linearize the problem, turning it into a set of linear equations, and then use linear algebra software packages to compute an approximate solution. Linearizing a nonlinear problem gives us inexact answers, and in a recursive algorithm

like a Kalman filter or weather tracking system these small errors can sometimes reinforce each other at each step, quickly causing the algorithm to spit out nonsense.

What we are about to embark upon is a difficult problem. There is not one obvious, correct, mathematically optimal solution anymore. We will be using approximations, we will be introducing errors into our computations, and we will forever be battling filters that diverge, that is, filters whose numerical errors overwhelm the solution.

In the remainder of this short chapter I will illustrate the specific problems the nonlinear Kalman filter faces. You can only design a filter after understanding the particular problems the nonlinearity in your problem causes. Subsequent chapters will then teach you how to design and implement different kinds of nonlinear filters.

## 9.2 The Problem with Nonlinearity

As asserted in the introduction the only math you really know how to do is linear math, that is, equations of the form

$$A\mathbf{x} = \mathbf{b}$$

That may strike you as hyperbole. After all, in this book we have integrated a polynomial to get distance from velocity and time. We know how to integrate polynomials, and so we are able to find the closed form equation for distance given velocity and time:

$$\int (vt + v_0) dt = \frac{a}{2}t^2 + v_0 t + d_0$$

That's a nonlinear equation that we know how to solve. But it is also a very special form. You spent a lot of time, probably at least a year, learning how to integrate various terms, and you still can not integrate some arbitrary equation - no one can. We don't know how. If you took freshman Physics you perhaps remember homework involving sliding frictionless blocks on a perfect plane and other toy problems. At the end of the course you were almost entirely unequipped to solve real world problems because the real world is nonlinear, and you were taught linear, closed forms of equations. It made the math tractable, but mostly useless.

The mathematics of the Kalman filter is beautiful in part due to the Gaussian equation being so special. It is nonlinear, but when we add and multiply them using linear algebra we get another Gaussian equation as a result. That is very rare.  $\sin x * \sin y$  does not yield a sin as an output.

What I mean by linearity may be obvious, but there are some subtleties. The mathematical requirements are twofold:

- additivity:  $f(x + y) = f(x) + f(y)$
- homogeneity:  $f(ax) = af(x)$

This leads us to say that a linear system is defined as a system whose output is linearly proportional to the sum of all its inputs. A consequence of this is that to be linear if the input is zero than the output must also be zero. Consider an audio amp - if I sing into a microphone, and you start talking, the output should be the sum of our voices (input) scaled by the amplifier gain. But if amplifier outputs a nonzero signal such as a hum for a zero input the additive relationship no longer holds. This is because you linearity requires that  $\text{amp}(\text{voice}) = \text{amp}(\text{voice} + 0)$  This clearly should give the same output, but if  $\text{amp}(0)$  is nonzero, then

$$\begin{aligned} \text{amp}(\text{voice}) &= \text{amp}(\text{voice} + 0) \\ &= \text{amp}(\text{voice}) + \text{amp}(0) \\ &= \text{amp}(\text{voice}) + \text{non\_zero\_value} \end{aligned}$$

which is clearly nonsense. Hence, an apparently linear equation such as

$$L(f(t)) = f(t) + 1$$

is not linear because  $L(0) = 1!$  Be careful!

### 9.3 An Intuitive Look at the Problem

I particularly like the following way of looking at the problem, which I am borrowing from Dan Simon's Optimal State Estimation [1]. Consider a tracking problem where we get the range and bearing to a target, and we want to track its position. The reported distance is 50 km, and the reported angle is  $90^\circ$ . Assume that the errors in both range and angle are distributed in a Gaussian manner. Given an infinite number of measurements what is the expected value of the position?

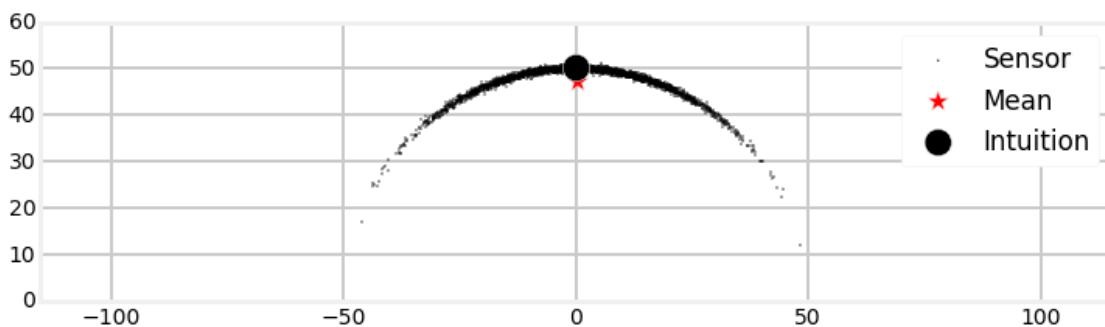
I have been recommending using intuition to gain insight, so let's see how it fares for this problem (hint: nonlinear problems are not intuitive). We might reason that since the mean of the range will be 50 km, and the mean of the angle will be  $90^\circ$ , that clearly the answer will be  $x=0$  km,  $y=50$  km.

Well, let's plot that and find out. Here are 3000 points plotted with a normal distribution of the distance of 0.4 km, and the angle having a normal distribution of 0.35 radians. We compute the average of the all of the positions, and display it as a star. Our intuition is displayed with a large circle.

```
In [15]: import numpy as np
        from numpy.random import randn
        import matplotlib.pyplot as plt

N = 3000
a = np.pi/2. + (randn(N) * 0.35)
r = 50.0      + (randn(N) * 0.4)
xs = r * np.cos(a)
ys = r * np.sin(a)

plt.scatter(xs, ys, label='Sensor', color='k', marker='.', s=2)
xs, ys = sum(xs)/N, sum(ys)/N
plt.scatter(xs, ys, c='r', marker='*', s=200, label='Mean')
plt.scatter(0, 50, c='k', marker='o', s=300, label='Intuition')
plt.axis('equal')
plt.legend();
```

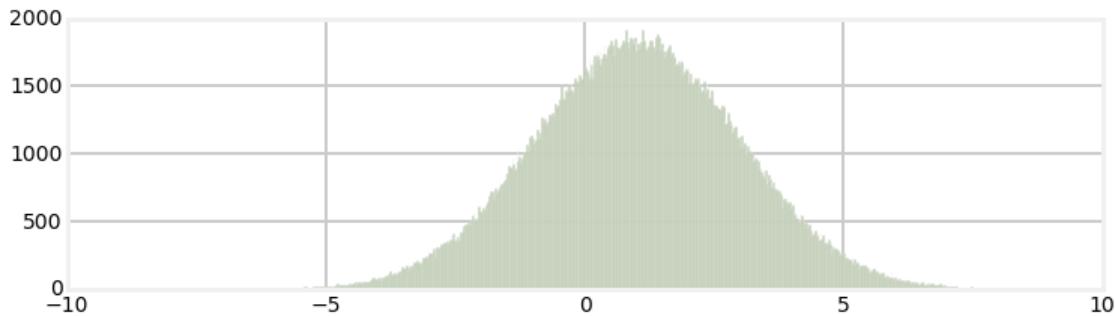


We can see that our intuition failed us because the nonlinearity of the problem forced all of the errors to be biased in one direction. This bias, over many iterations, can cause the Kalman filter to diverge. Even if it doesn't diverge the solution will not be optimal. Linear approximations applied to nonlinear problems yields inaccurate results.

## 9.4 The Effect of Nonlinear Functions on Gaussians

Gaussians are not closed under an arbitrary nonlinear function. Recall the equations of the Kalman filter - at each evolution (prediction) we pass the Gaussian representing the state through the process function to get the Gaussian at time  $k$ . Our process function was always linear, so the output was always another Gaussian. Let's look at that on a graph. I will take an arbitrary Gaussian and pass it through the function  $f(x) = 2x + 1$  and plot the result. We know how to do this analytically, but let's use sampling. I will generate 500,000 points with a normal distribution, pass them through  $f(x)$ , and plot the results. I do it this way because the next example will be nonlinear, and we will have no way to compute this analytically.

```
In [16]: import numpy as np
        from numpy.random import normal
        gaussian = (0., 1.)
        data = normal(loc=gaussian[0], scale=gaussian[1], size=500000)
        plt.hist(2*data + 1, 1000);
```

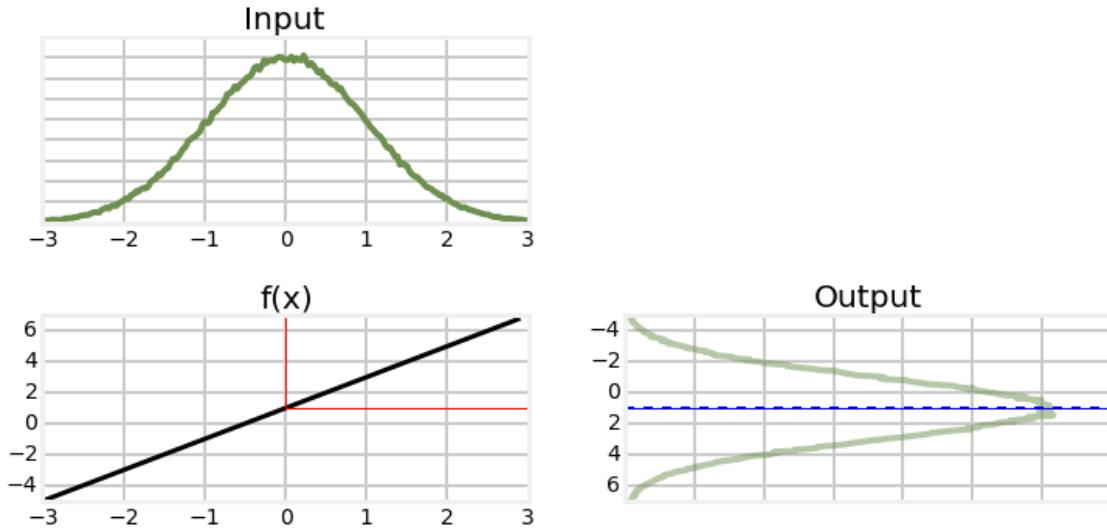


This is an unsurprising result. The result of passing the Gaussian through  $f(x) = 2x + 1$  is another Gaussian centered around 1. Let's look at the input, nonlinear function, and output at once.

```
In [17]: from book_format import set_figsize, figsize
        from nonlinear_plots import plot_nonlinear_func

        def g1(x):
            return 2*x+1

        with figsize(y=5):
            plot_nonlinear_func(data, g1, gaussian)
```



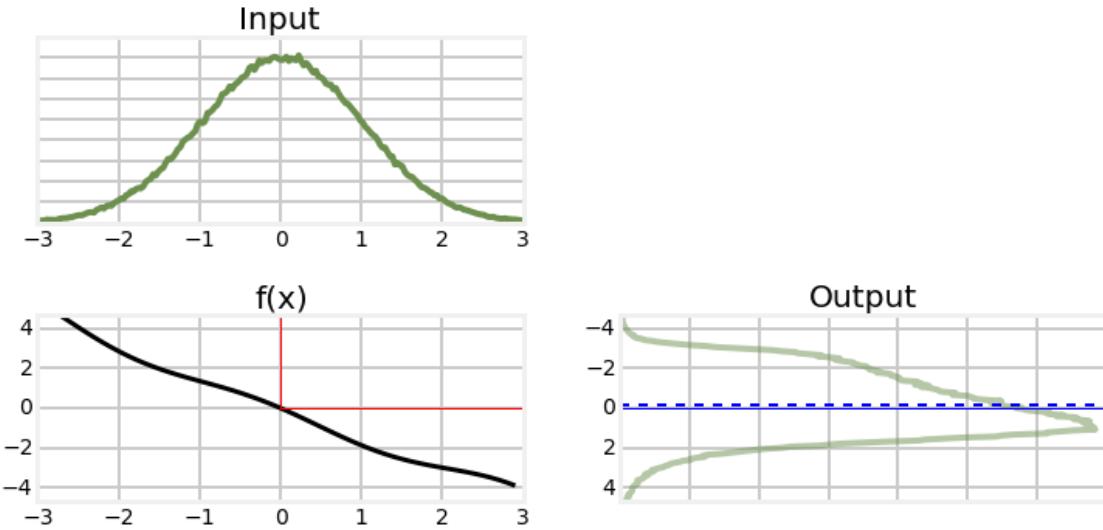
I explain how to plot Gaussians, and much more, in the Notebook [Computing\\_and\\_Plotting\\_PDFs](#) in the Supporting\_Notebooks folder. You can also read it online [here](#)[1]

The plot labeled ‘Input’ is the histogram of the original data. This is passed through the function  $f(x) = 2x + 1$  which is displayed in the chart on the bottom left. The red lines shows how one value,  $x = 0$  is passed through the function. Each value from input is passed through in the same way to the output function on the right. For the output I computed the mean by taking the average of all the points, and drew the results with the dotted blue line. A solid blue line shows the actual mean for the point  $x = 0$ . The output looks like a Gaussian, and is in fact a Gaussian. We can see that the variance in the output is larger than the variance in the input, and the mean has been shifted from 0 to 1, which is what we would expect given the transfer function  $f(x) = 2x + 1$ . The  $2x$  affects the variance, and the  $+1$  shifts the mean. The computed mean, represented by the dotted blue line, is nearly equal to the actual mean. If we used more points in our computation we could get arbitrarily close to the actual value.

Now let’s look at a nonlinear function and see how it affects the probability distribution.

```
In [18]: def g2(x):
    return (np.cos(3*(x/2 + 0.7))) * np.sin(0.3*x) - 1.6*x

    with figsize(y=5):
        plot_nonlinear_func(data, g2, gaussian)
```



This result may be somewhat surprising to you. The function looks “fairly” linear - it is pretty close to a straight line, but the probability distribution of the output is completely different from a Gaussian. Recall the equations for multiplying two univariate Gaussians:

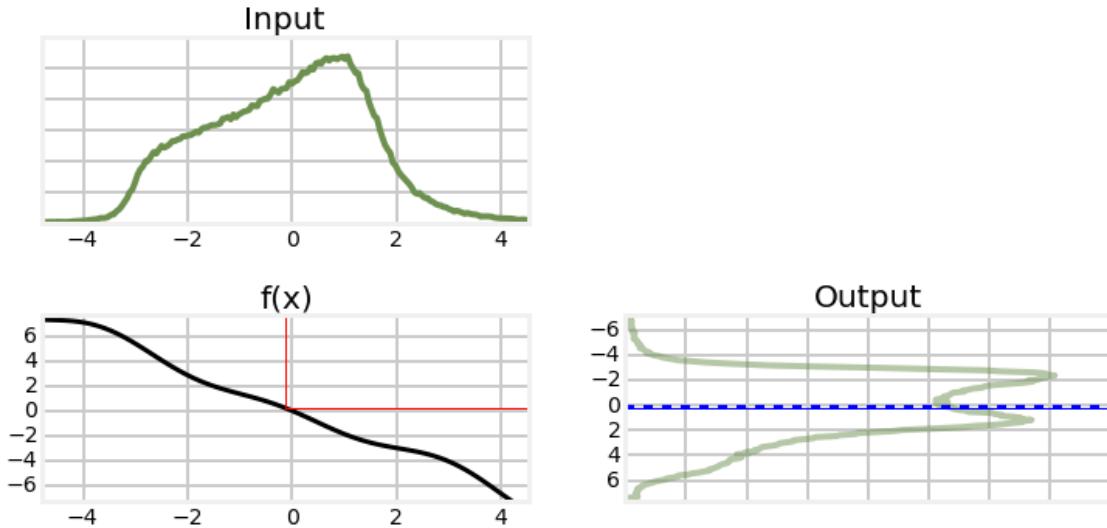
$$\mu = \frac{\sigma_1^2 \mu_2 + \sigma_2^2 \mu_1}{\sigma_1^2 + \sigma_2^2}$$

$$\sigma = \sqrt{\frac{1}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}}}$$

These equations do not hold for non-Gaussians, and certainly do not hold for the probability distribution shown in the ‘Output’ chart above.

Think of what this implies for the Kalman filter algorithm of the previous chapter. All of the equations assume that a Gaussian passed through the process function results in another Gaussian. If this is not true then all of the assumptions and guarantees of the Kalman filter do not hold. Let’s look at what happens when we pass the output back through the function again, simulating the next step time step of the Kalman filter.

```
In [19]: y = g2(data)
gaussian2 = (np.mean(y), np.var(y))
with figsize(y=5):
    plot_nonlinear_func(y, g2, gaussian2)
```



As you can see the probability function is further distorted from the original Gaussian. However, the graph is still somewhat symmetric around  $x=0$ , let's see what the mean is.

```
In [20]: print('input mean, variance: %.4f, %.4f' %
           (np.mean(data), np.var(data)))
print('output mean, variance: %.4f, %.4f' %
      (np.mean(y), np.var(y)))

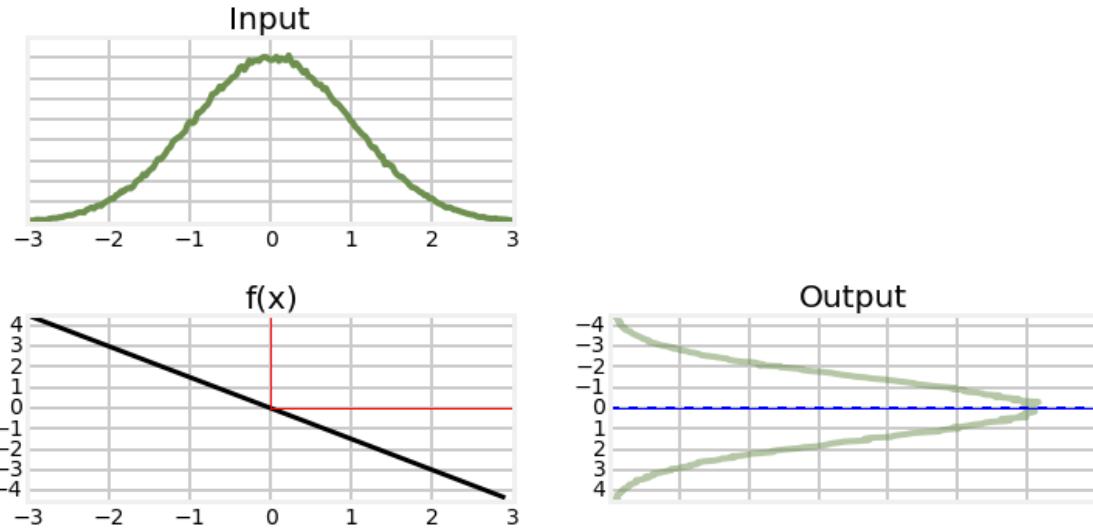
input mean, variance: -0.0005, 0.9999
output mean, variance: -0.1240, 2.4029
```

Let's compare that to the linear function that passes through  $(-2,3)$  and  $(2,-3)$ , which is very close to the nonlinear function we have plotted. Using the equation of a line we have

$$m = \frac{-3 - 3}{2 - (-2)} = -1.5$$

```
In [21]: def g3(x):
    return -1.5 * x

with figsize(y=5):
    plot_nonlinear_func(data, g3, gaussian)
out = g3(data)
print('output mean, variance: %.4f, %.4f' %
      (np.mean(out), np.var(out)))
```



```
output mean, variance: 0.0007, 2.2498
```

Although the shapes of the output are very different, the mean and variance of each are almost the same. This may lead us to reasoning that perhaps we can ignore this problem if the nonlinear equation is ‘close to’ linear. To test that, we can iterate several times and then compare the results.

```
In [22]: out = g3(data)
out2 = g2(data)

for i in range(10):
    out = g3(out)
    out2 = g2(out2)
print('linear      output mean, variance: %.4f, %.4f' %
      (np.average(out), np.std(out)**2))
print('nonlinear output mean, variance: %.4f, %.4f' %
      (np.average(out2), np.std(out2)**2))

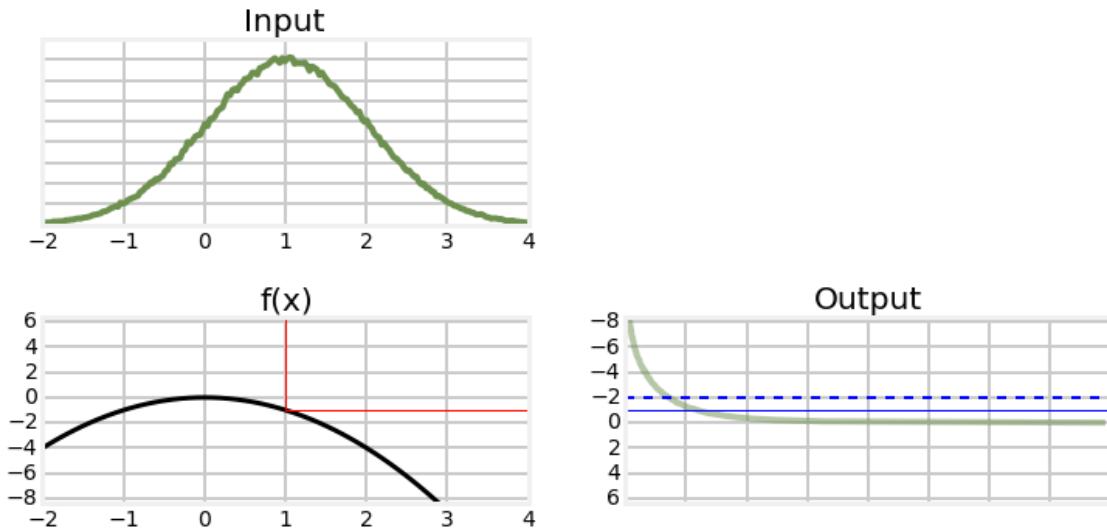
linear      output mean, variance: 0.0404, 7481.2295
nonlinear output mean, variance: -9.0832, 30461.4199
```

Unfortunately the nonlinear version is not stable. It drifted significantly from the mean of 0, and the variance is half an order of magnitude larger.

I minimized the issue by using a function that is quite close to a straight line. What happens if the function is  $y(x) = x^2$ ?

```
In [23]: def g3(x):
    return -x*x
x0 = (1, 1)
data = normal(loc=x0[0], scale=x0[1], size=500000)

with figsize(y=5):
    plot_nonlinear_func(data, g3, gaussian=x0)
```

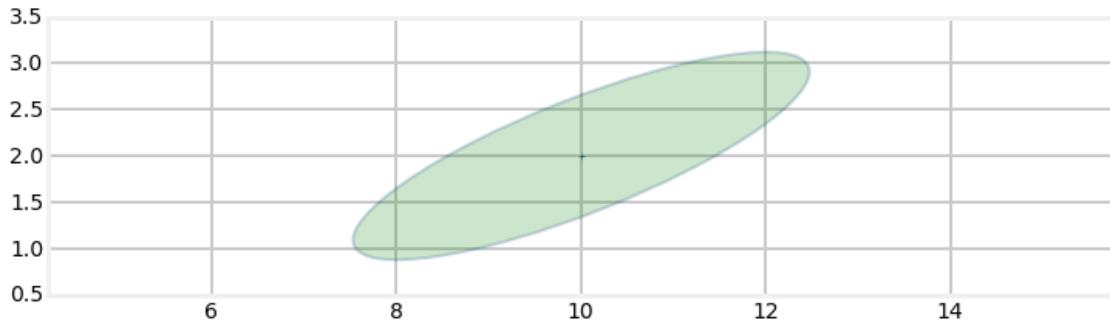


Despite the curve being smooth and reasonably straight at  $x = 1$  the probability distribution of the output doesn't look anything like a Gaussian and the computed mean of the output is quite different than the value computed directly. This is not an unusual function - a ballistic object moves in a parabola, and this is the sort of nonlinearity your filter will need to handle. If you recall we've tried to track a ball and failed miserably. This graph should give you insight into why the filter performed so poorly.

## 9.5 A 2D Example

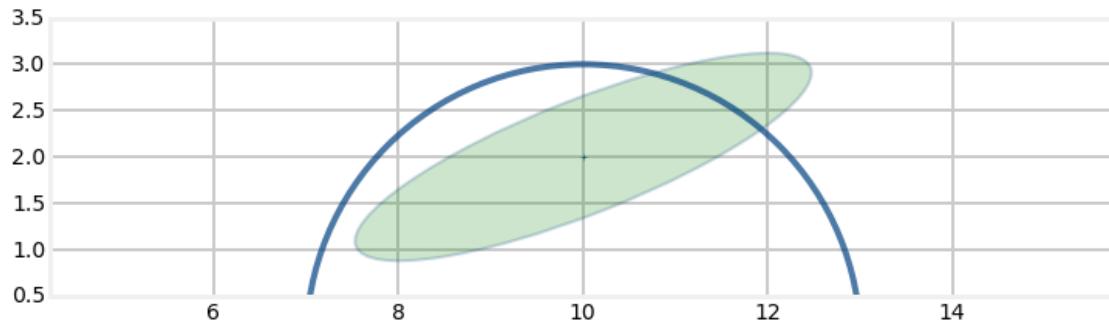
It is hard to look at probability distributions and reason about what will happen in a filter. So let's think about tracking an aircraft with radar. The estimate may have a covariance that looks like this:

```
In [24]: import nonlinear_internal
nonlinear_internal.plot1()
```



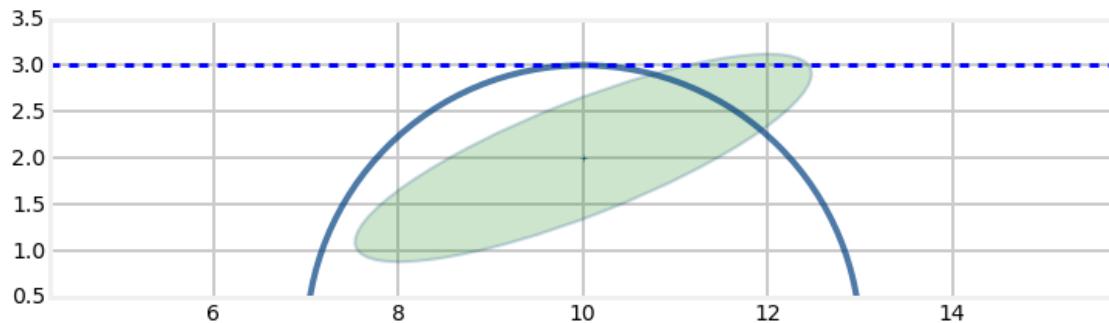
What happens when we try to linearize this problem? The radar gives us a range to the aircraft. Suppose the radar is directly under the aircraft ( $x=10$ ) and the next measurement states that the aircraft is 3 miles away ( $y=3$ ). The positions that could match that measurement form a circle with radius 3 miles, like so.

In [25]: `nonlinear_internal.plot2()`



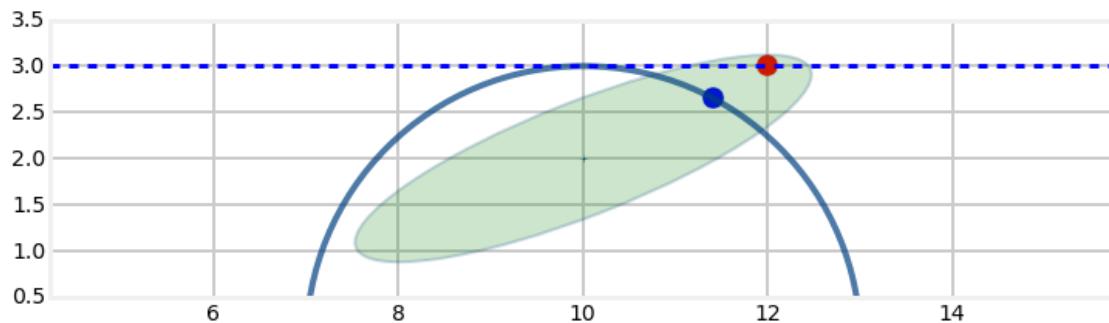
We can see by inspection that the probable position of the aircraft is somewhere near  $x=11.4$ ,  $y=2.7$  because that is where the covariance ellipse and range measurement overlap. But the range measurement is nonlinear so we have to linearize it. We haven't covered this material yet, but the Extended Kalman filter will linearize at the last position of the aircraft -  $(10, 2)$ . At  $x=10$  the range measurement has  $y=3$ , and so we linearize at that point.

In [26]: `nonlinear_internal.plot3()`



Now we have a linear representation of the problem (literally a straight line) which we can solve. Unfortunately you can see that the intersection of the line and the covariance ellipse is a long way from the actual aircraft position.

In [27]: `nonlinear_internal.plot4()`



That sort of error often leads to disastrous results. The error in this estimate is large. But in the next innovation of the filter that very bad estimate will be used to linearize the next radar measurement, so the next estimate is likely to be markedly worse than this one. After only a few iterations the Kalman filter will diverge, and start producing results that have no correspondence to reality.

This covariance ellipse spans miles. I exaggerated the size to illustrate the difficulties of highly nonlinear systems. In real radar tracking problems the nonlinearity is usually not that bad, but the errors will still accumulate. Other systems you may work on could have this amount of nonlinearity - this was not an exaggeration only to make a point. You will always be battling divergence when working with nonlinear systems.

## 9.6 The Algorithms

You may be impatient to solve a specific problem, and wondering which filter to use. I will quickly survey the options. The subsequent chapters are somewhat independent of each other, and you can fruitfully skip around, though I recommend reading linearly if you truly want to master all of the material.

The workhorses of nonlinear filters are the linearized Kalman filter and extended Kalman filter (EKF). These two techniques were invented shortly after Kalman published his paper and they have been the main techniques used since then. The flight software in airplanes, the GPS in your car or phone almost certainly use one of these techniques.

However, these techniques are extremely demanding. The EKF linearizes the differential equations at one point, which requires you to find a solution to a matrix of partial derivatives (a Jacobian). This can be difficult or impossible to do analytically. If impossible, you have to use numerical techniques to find the Jacobian, but this is expensive computationally and introduces more error into the system. Finally, if the problem is quite nonlinear the linearization leads to a lot of error being introduced in each step, and the filters frequently diverge. You can not throw some equations into some arbitrary solver and expect to get good results. It's a difficult field for professionals. I note that most Kalman filtering textbooks merely gloss over the EKF despite it being the most frequently used technique in real world applications.

Recently the field has been changing in exciting ways. First, computing power has grown to the point that we can use techniques that were once beyond the ability of a supercomputer. These use Monte Carlo techniques - the computer generates thousands to tens of thousands of random points and tests all of them against the measurements. It then probabilistically kills or duplicates points based on how well they match the measurements. A point far away from the measurement is unlikely to be retained, whereas a point very close is quite likely to be retained. After a few iterations there is a clump of particles closely tracking your object, and a sparse cloud of points where there is no object.

This has two benefits. First, the algorithm is robust even for extremely nonlinear problems. Second, the algorithm can track arbitrarily many objects at once - some particles will match the behavior of one object, and other particles will match other objects. So this technique is often used to track automobile traffic, people in crowds, and so on.

The costs should be clear. It is computationally expensive to test tens of thousands of points for every step in the filter. But modern CPUs are very fast, and this is a good problem for GPUs because the part of the algorithm is parallelizable. Another cost is that the answer is not mathematical. With a Kalman filter my covariance matrix gives me important information about the amount of error in the estimate. The particle filter does not give me a rigorous way to compute this. Finally, the output of the filter is a cloud of points; I then have to figure out how to interpret it. Usually you will be doing something like taking the mean and standard deviations of the points, but this is a difficult problem. There are still many points that do not 'belong' to a tracked object, so you first have to run some sort of clustering algorithm to first find the points that seem to be tracking an object, and then you need another algorithm to produce an state estimate from those points. None of this is intractable, but it is all quite computationally expensive.

Finally, we have a new algorithm called the unscented Kalman filter (UKF) that does not require you to find analytic solutions to nonlinear equations, and yet almost always performs better than the EKF. It does especially well with highly nonlinear problems - problems where the EKF has significant difficulties. Designing the filter is extremely easy. Some will say the jury is still out on the UKF, but to my mind the UKF is superior in almost every way to the EKF. I suggest that the UKF should be the starting point for any implementation, especially if you are not a Kalman filter professional with a graduate degree in control theory. The main downside is that the UKF can be a few times slower than the EKF, but this really depends on whether the EKF solves the Jacobian analytically or numerically. If numerically the UKF is almost certainly faster. It has not been proven (and probably it cannot be proven) that the UKF always yields more accurate results than the EKF. In practice it almost always does, often significantly so. It is very easy to understand and implement, and I strongly suggest this filter as your starting point.

## 9.7 Summary

The world is nonlinear, but we only really know how to solve linear problems. This introduces significant difficulties for Kalman filters. We've looked at how nonlinearity affects filtering in 3 different but equivalent ways, and I've given you a brief summary of the major approaches: the linearized Kalman filter, the extended Kalman filter, the Unscented Kalman filter, and the particle filter.

Until recently the linearized Kalman filter and EKF have been the standard way to solve these problems. They are very difficult to understand and use, and they are also potentially very unstable.

Recent developments have offered what are to my mind superior approaches. The UKF dispenses with the need to find solutions to partial differential equations, yet it is also usually more accurate than the EKF. It is easy to use and understand. I can get a basic UKF going in a few minutes by using FilterPy. The particle filter dispenses with mathematical modeling completely in favor of a Monte Carlo technique of generating a random cloud of thousands of points. It runs slowly, but it can solve otherwise intractable problems with relative ease.

I get more email about the EKF than anything else; I suspect that this is because most treatments in books, papers, and on the internet use the EKF. If your interest is in mastering the field of course you will want to learn about the EKF. But if you are just trying to get good results I point you to the UKF and particle filter first. They are much easier to implement, understand, and use, and they are typically far more stable than the EKF.

Some will quibble with that advice. A lot of recent publications are devoted to a comparison of the EKF, UKF, and perhaps a few other choices for a given problem. Do you not need to perform a similar comparison for your problem? If you are sending a rocket to Mars then of course you do. You will be balancing issues such as accuracy, round off errors, divergence, mathematical proof of correctness, and the computational effort required. I can't imagine not knowing the EKF intimately.

On the other hand the UKF works spectacularly! I use it at work for real world applications. I mostly haven't even tried to implement an EKF for these applications because I can verify that the UKF is working fine. Is it possible that I might eke out another 0.2% of performance from the EKF in certain situations? Sure! Do I care? No! I completely understand the UKF implementation, it is easy to test and verify, I can pass the code to others and be confident that they can understand and modify it, and I am not a masochist that wants to battle difficult equations when I already have a working solution. If the UKF or particle filters start to perform poorly for some problem then I will turn other techniques, but not before then. And realistically, the UKF usually provides substantially better performance than the EKF over a wide range of problems and conditions. If "really good" is good enough I'm going to spend my time working on other problems.

I'm belaboring this point because in most textbooks the EKF is given center stage, and the UKF is either not mentioned at all or just given a 2 page gloss that leaves you completely unprepared to use the filter. The UKF is still relatively new, and it takes time to write new editions of books. At the time many books were written the UKF was either not discovered yet, or it was just an unproven but promising curiosity. But I

am writing this now, the UKF has had enormous success, and it needs to be in your toolkit. That is what I will spend most of my effort trying to teach you.

## 9.8 References

- [1] [https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python/blob/master/Supporting\\_Notebooks/Computing\\_and](https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python/blob/master/Supporting_Notebooks/Computing_and)



## Chapter 10

# The Unscented Kalman Filter

In the last chapter we discussed the difficulties that nonlinear systems pose. This nonlinearity can appear in two places. It can be in our measurements, such as a radar that is measuring the slant range to an object. Slant range requires you to take a square root to compute the x,y coordinates:

$$x = \sqrt{\text{slant}^2 - \text{altitude}^2}$$

The nonlinearity can also occur in the process model - we may be tracking a ball traveling through the air, where the effects of gravity and air drag lead to nonlinear behavior. The standard Kalman filter performs poorly or not at all with these sorts of problems.

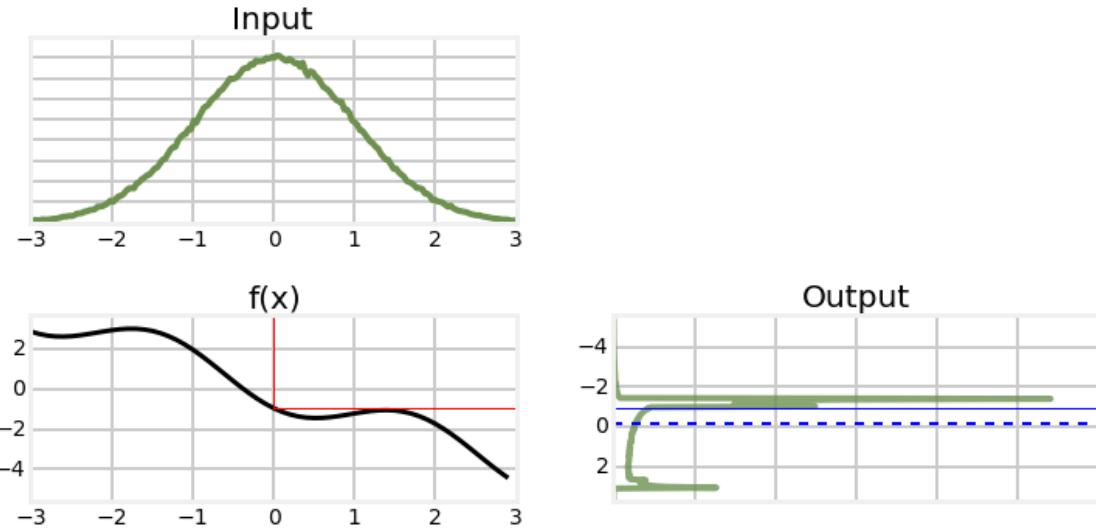
In the last chapter I showed you a plot like this. I have altered the equation somewhat to emphasize the effects of nonlinearity.

```
In [2]: from book_format import set_figsize, figsize
        import matplotlib.pyplot as plt
        from nonlinear_plots import plot_nonlinear_func
        from numpy.random import normal
        import numpy as np

        gaussian=(0., 1.)
        data = normal(loc=gaussian[0], scale=gaussian[1], size=500000)

        def g(x):
            return (np.cos(4*(x/2 + 0.7))) - 1.3*x

        with figsize(y=5):
            plot_nonlinear_func(data, g, gaussian=gaussian)
```



I generated this by taking 500,000 samples from the input, passing it through the nonlinear transform, and building a histogram of the result. We call these points sigma points. From the output histogram we can compute a mean and standard deviation which would give us an updated, albeit approximated Gaussian.

It has perhaps occurred to you that this sampling process constitutes a solution to our problem. Suppose for every update we generated 500,000 points, passed them through the function, and then computed the mean and variance of the result. This is called a Monte Carlo approach, and it used by some Kalman filter designs, such as the Ensemble filter and particle filter. Sampling requires no specialized knowledge, and does not require a closed form solution. No matter how nonlinear or poorly behaved the function is, as long as we sample with enough sigma points we will build an accurate output distribution.

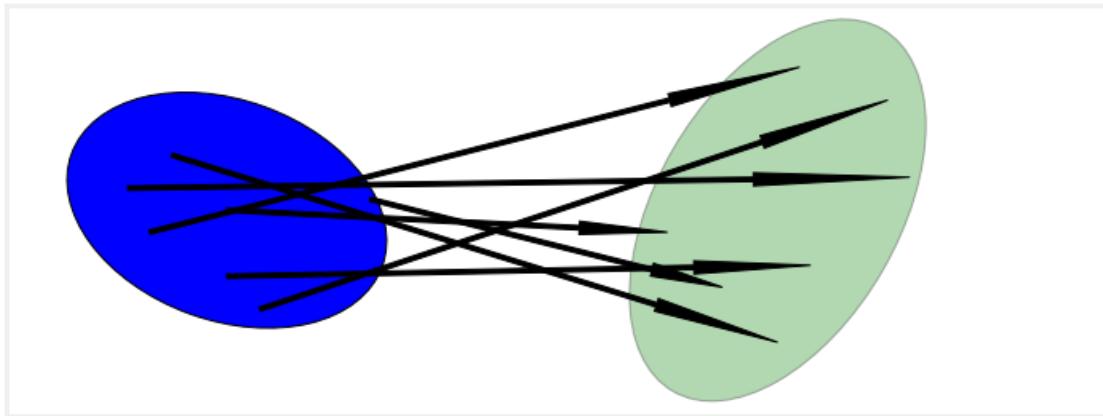
“Enough points” is the rub. The graph above was created with 500,000 sigma points, and the output is still not smooth. What’s worse, this is only for 1 dimension. In general, the number of points required increases by the power of the number of dimensions. If you only needed 500 points for 1 dimension, you’d need 500 squared, or 250,000 points for two dimensions, 500 cubed, or 125,000,000 points for three dimensions, and so on. So while this approach does work, it is very computationally expensive. Ensemble filters and particle filters use clever techniques to significantly reduce this dimensionality, but the computational burdens are still very large. The unscented Kalman filter uses sigma points but drastically reduces the amount of computation by using a deterministic method to choose the points.

## 10.1 Sigma Points - Sampling from a Distribution

Let’s look at the problem in terms of a 2D covariance ellipse. I choose 2D merely because it is easy to plot; this extends to any number of dimensions. Assuming some arbitrary nonlinear function, we will take random points from the first covariance ellipse, pass them through the nonlinear function, and plot their new position. Then we can compute the the mean and covariance of the transformed points, and use that as our estimate of the mean and probability distribution.

```
In [3]: import ukf_internal

with figsize(y=4.3):
    ukf_internal.show_2d_transform()
```



On the left we show an ellipse depicting the  $1\sigma$  distribution of two state variables. The arrows show how several randomly sampled points might be transformed by some arbitrary nonlinear function to a new distribution. The ellipse on the right is drawn semi-transparently to indicate that it is an estimate of the mean and variance of this collection of points.

Let's write a function which passes 10,000 points randomly drawn from the Gaussian

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 32 & 15 \\ 15 & 40 \end{bmatrix}$$

through the nonlinear system:

$$\begin{cases} x = x + y \\ y = 0.1x^2 + y^2 \end{cases}$$

```
In [4]: import numpy as np
from numpy.random import multivariate_normal
from nonlinear_plots import plot_monte_carlo_mean

def f(x,y):
    return x+y, .1*x**2 + y*y

mean = (0, 0)
p = np.array([[32, 15], [15., 40.]])

# Compute linearized mean
mean_fx = f(*mean)

#generate random points
xs, ys = multivariate_normal(mean=mean, cov=p, size=10000).T
plot_monte_carlo_mean(xs, ys, f, mean_fx, 'Linearized Mean');
```

Difference in mean x=0.087, y=42.362



This plot shows the strong nonlinearity that occurs with this function, and the large error that would result if we linearized in the way of the Extended Kalman filter (we will be learning this in the next chapter).

## 10.2 Choosing Sigma Points

I used 10,000 randomly generated sigma points to generate this solution. While the computed mean is quite accurate, computing 10,000 points for every update would cause our filter to be very slow. So, what would be fewest number of sampled points that we can use, and what kinds of constraints does this problem formulation put on the points? We will assume that we have no special knowledge about the nonlinear function as we want to find a generalized algorithm that works for any function.

Let's consider the simplest possible case and see if it offers any insight. The simplest possible system is the **identity function**:  $f(x) = x$ . If our algorithm does not work for the identity function then the filter will never converge. In other words, if the input is 1 (for a one dimensional system), the output must also be 1. If the output was different, such as 1.1, then when we fed 1.1 into the transform at the next time step, we'd get out yet another number, maybe 1.23. This filter diverges.

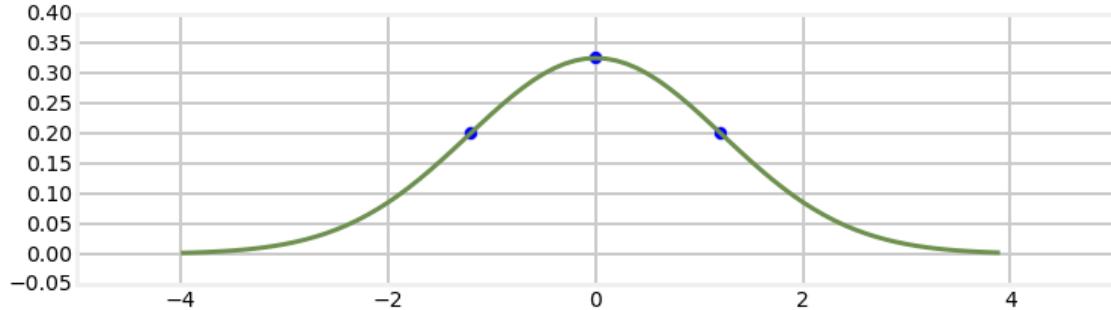
The fewest number of points that we can use is one per dimension. This is the number that the linear Kalman filter uses. The input to a Kalman filter for the distribution  $\mathcal{N}(\mu, \sigma^2)$  is  $\mu$  itself. So while this works for the linear case, it is not a good answer for the nonlinear case.

Perhaps we can use one point per dimension, but altered somehow. However, if we were to pass some value  $\mu + \Delta$  into the identity function  $f(x) = x$  it would not converge, so this will not work. If we didn't alter  $\mu$  then this would be the standard Kalman filter. We must conclude that one sample will not work.

What is the next lowest number we can choose? Two. Consider the fact that Gaussians are symmetric, and that we probably want to always have one of our sample points be the mean of the input for the identity function to work. Two points would require us to select the mean, and then one other point. That one other point would introduce an asymmetry in our input that we probably don't want. It would be very difficult to make this work for the identity function  $f(x) = x$ .

The next lowest number is 3 points. 3 points allows us to select the mean, and then one point on each side of the mean, as depicted on the chart below.

In [5]: `ukf_internal.show_3_sigma_points()`



So we can take those three points, pass them through a nonlinear function  $f(x)$ , and compute a new mean and variance. We can compute the mean as the average of the 3 points, but that is not very general. For example, for a very nonlinear problem we might want to weight the center point much higher than the outside points, or we might want to weight the outside points higher if the distribution is not Gaussian. A more general approach is to compute the mean as  $\mu = \sum_i w_i f(\mathcal{X}_i)$ .

For this to work for the identity function we want the sums of the weights for the mean to equal one. We can always come up with counterexamples, but in general if the sum is greater or less than one the sampling will not yield the correct output. Given that, we then have to select sigma points  $\mathcal{X}$  and their corresponding weights so that they compute to the mean and variance of the input Gaussian.

It is possible to use different weights for the mean ( $w^m$ ) and for the variance ( $w^c$ ). So we can write

**Constraints :**

$$\begin{aligned} 1 &= \sum_i w_i^m \\ 1 &= \sum_i w_i^c \\ \mu &= \sum_i w_i^m f(\mathcal{X}_i) \\ \Sigma &= \sum_i w_i^c (f(\mathcal{X})_i - \mu)(f(\mathcal{X})_i - \mu)^T \end{aligned}$$

The first two equations are the constraint that the weights must sum to one. The third equation is how you compute a weight mean. The forth equation may be less familiar, but recall that the equation for a covariance is:

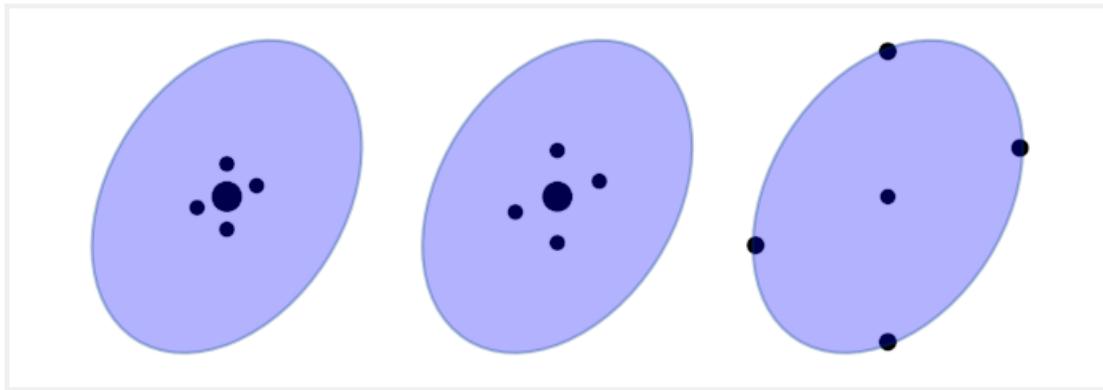
$$COV(x, y) = \frac{\sum(x - \bar{x})(y - \bar{y})}{n}$$

and you should see where it came from.

These constraints do not limit us to a unique answer. For example, if you choose a smaller weight for the point at the mean for the input, you could compensate by choosing larger weights for the rest of the  $\mathcal{X}$ , and vice versa. We can use different weights for the mean and variances, or the same weights. Indeed, these equations do not require that any of the points be the mean of the input at all, though it seems ‘nice’ to do so, so to speak.

We want an algorithm that satisfies the constraints, preferably with only 3 points per dimension. Before we go on I want to make sure the idea is clear. Below are three different examples for the same covariance ellipse with different sigma points. The size of the sigma points is proportional to the weight given to each.

```
In [6]: with figsize(y=4):
    ukf_internal.show_sigma_selections()
```



The points do not lie along the major and minor axis of the ellipse; nothing in the constraints above require me to do that. Furthermore, in each case I show the points evenly spaced; again, the constraints above do not require that.

We can see that the arrangement and spacing of the sigma points will affect how we sample our distribution. Points that are close together will sample local effects, and thus probably work better for very nonlinear problems. Points that are far apart, or far off the axis of the ellipse will sample non-local effects and non Gaussian behavior. However, by varying the weights used for each point we can mitigate this. If the points are far from the mean but weighted very slightly we will incorporate some of the knowledge about the distribution without allowing the nonlinearity of the problem to create a bad estimate.

### 10.3 The Unscented Transform

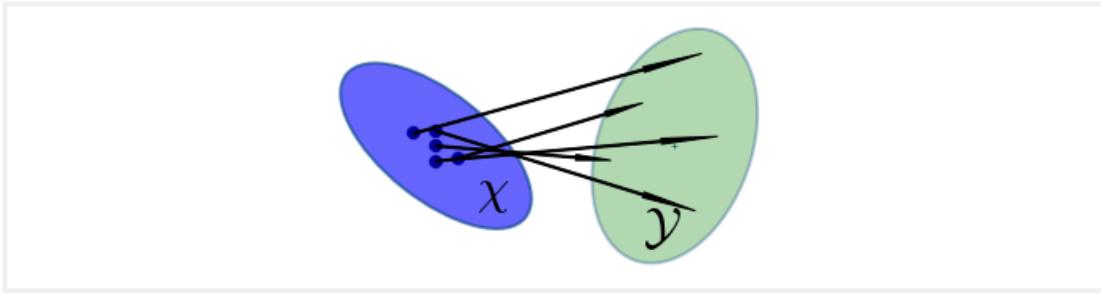
For the moment assume an algorithm exists for selecting the sigma points and weights. How are the sigma points used to implement a filter?

The unscented transform is the core of the algorithm yet it is remarkably simple. We pass the sigma points through a (usually nonlinear) function yielding a transformed set of points.

$$\mathbf{y} = f(\boldsymbol{\chi})$$

We then compute a new mean and covariance of the transformed sigma points. That mean and covariance becomes the new estimate. The figure below depicts the operation of the unscented transform. The green ellipse on the right represents the computed mean and covariance to the transformed sigma points.

```
In [7]: ukf_internal.show_sigma_transform(with_text=True)
```



The mean and covariances are computed using these equations.

$$\mu = \sum_i w_i^m \mathbf{y}_i$$

$$\Sigma = \sum_i w_i^c (\mathbf{y}_i - \mu)(\mathbf{y}_i - \mu)^T$$

These equations should be familiar - they are the constraint equations we developed above.

### 10.3.1 Unscented Transform Example

Earlier we wrote a function that found the mean of a distribution by passing 50,000 points through a nonlinear function. Let's now pass 5 sigma points through the same function, and compute their mean with the unscented transform.

```
In [8]: from filterpy.kalman import unscented_transform
        from filterpy.kalman import MerweScaledSigmaPoints as SigmaPoints
        import scipy.stats as stats

def f_nonlinear_xy(x, y):
    return np.array([x + y, .1*x**2 + y*y])

#initial mean and covariance
mean = (0, 0)
p = np.array([[32., 15],
              [15., 40.]])

### create sigma points - we will learn about this later
### in the chapter
points = SigmaPoints(n=2, alpha=.1, beta=2., kappa=1.)
Wm, Wc = points.weights()
sigmas = points.sigma_points(mean, p)

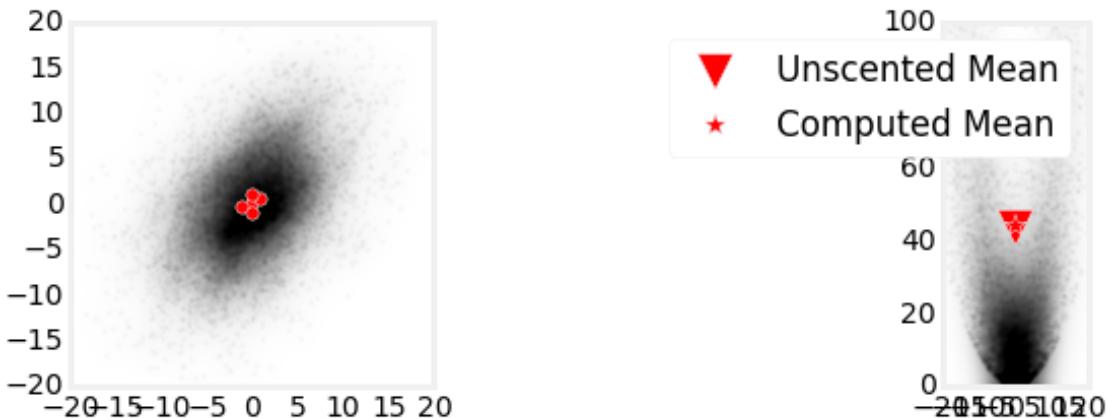
### pass through nonlinear function
sigmas_f = np.zeros((5, 2))
for i in range(5):
    sigmas_f[i] = f_nonlinear_xy(sigmas[i, 0], sigmas[i, 1])

### use unscented transform to get new mean and covariance
ukf_mean, ukf_cov = unscented_transform(sigmas_f, Wm, Wc)
```

```
#generate random points
np.random.seed(100)
xs, ys = multivariate_normal(mean=mean, cov=p, size=5000).T

plot_monte_carlo_mean(xs, ys, f, ukf_mean, 'Unscented Mean')
plt.xlim(-20, 20)
plt.ylim(0, 100)
plt.subplot(121)
plt.scatter(sigmas[:,0], sigmas[:,1], c='r', s=40);
```

Difference in mean x=-0.097, y=0.549



I find this result remarkable. Using only 5 points we were able to compute the mean with amazing accuracy. The error in x is only -0.097, and the error in y is 0.549. In contrast, a linearized approach (used by the predominant EKF filter) gave an error of over 43 in y. If you look at the code that generates the sigma points you'll see that it has no knowledge of the nonlinear function, only of the mean and covariance of our initial distribution. The same 5 sigma points would be generated if we had a completely different nonlinear function.

I will admit to choosing a nonlinear function that makes the performance of the unscented transform striking compared to the EKF. But the physical world is filled with very nonlinear behavior, and the UKF takes it in stride. I did not ‘work’ to find a function where the unscented transform happened to work well. You will see in the next chapter how more traditional techniques struggle with strong nonlinearities. This graph is the foundation of why I advise you to use the UKF or similar modern technique whenever possible.

## 10.4 The Unscented Kalman Filter

Now we can present the entire Unscented Kalman filter algorithm. Assume that there is a function  $f(x, dt)$  that performs the state transition for our filter - it predicts the next state given the current state. It is the analogue of the  $\mathbf{F}$  matrix of the Kalman filter. Also assume there is a measurement function  $h(x)$ . It takes the current state and computes what measurement that state corresponds to. It is the analogue of the  $\mathbf{H}$  matrix used by the Kalman filter.

### 10.4.1 Predict Step

As with the linear Kalman filter, the UKF's predict step computes the mean and covariance of the system for the next evolution using the process model. However, the computation is quite different.

First, we generate sigma points  $\mathcal{X}$  and their corresponding weights  $W^m, W^c$  according to some function:

$$\begin{aligned}\mathcal{X} &= \text{sigma\_function}(\mu, \Sigma) \\ W^m, W^c &= \text{weight\_function}(n, \text{parameters})\end{aligned}$$

We pass each sigma point through a state transition function  $f(\mathcal{X})$  which you define. This projects the sigma points forward in time according to the process model.

$$\mathcal{Y} = f(\mathcal{X})$$

Now we compute the predicted mean and covariance using the unscented transform on the transformed sigma points. I've dropped the subscript  $i$  for readability.

$$\begin{aligned}\bar{\mu} &= \sum w^m \mathcal{Y} \\ \bar{\Sigma} &= \sum w^c (\mathcal{Y} - \bar{\mu})(\mathcal{Y} - \bar{\mu})^\top + \mathbf{Q}\end{aligned}$$

This table compares the Kalman filter with the Unscented Kalman Filter equations.

Kalman	Unscented
$\bar{\mathbf{x}} = \mathbf{F}\mathbf{x}$	$\bar{\mu} = \sum w^m \mathcal{Y}$
$\bar{\mathbf{P}} = \mathbf{F}\mathbf{P}\mathbf{F}^\top + \mathbf{Q}$	$\sum w^c (\mathcal{Y} - \bar{\mu})(\mathcal{Y} - \bar{\mu})^\top + \mathbf{Q}$

### 10.4.2 Update Step

Now we can perform the update step of the filter. Recall that Kalman filters perform the update state in measurement space. So, the first thing we do is convert the sigma points from the predict step into measurements using a measurement function  $h(x)$  that you define.

$$\mathcal{Z} = h(\mathcal{Y})$$

Now we can compute the mean and covariance of these points using the unscented transform.

$$\begin{aligned}\mu_z &= \sum w^m \mathcal{Z} \\ \mathbf{P}_z &= \sum w^c (\mathcal{Z} - \bar{\mu})(\mathcal{Z} - \bar{\mu})^\top + \mathbf{R}\end{aligned}$$

The  $z$  subscript denotes that these are the mean and covariance for the measurements.

All that is left is to compute the residual and Kalman gain. The residual of the measurement  $\mathbf{z}$  is trivial to compute:

$$\mathbf{y} = \mathbf{z} - \mu_z$$

To compute the Kalman gain we first compute the [cross covariance](#) of the state and the measurements, which is defined as:

$$\mathbf{P}_{xz} = \sum w^c (\mathcal{X} - \mu)(\mathcal{Z} - \mu_z)^\top$$

And then the Kalman gain is defined as

$$\mathbf{K} = \mathbf{P}_{xz}\mathbf{P}_z^{-1}$$

If you think of the inverse as a kind of matrix division, you can see that the Kalman gain is a simple ratio which computes:

$$\mathbf{K} \approx \frac{\mathbf{P}_{xz}}{\mathbf{P}_z^{-1}} \approx \frac{\text{belief in state}}{\text{belief in measurement}}$$

Finally, we compute the new state estimate using the residual and Kalman gain:

$$\mathbf{x} = \bar{\mathbf{x}} + \mathbf{Ky}$$

and the new covariance is computed as:

$$\mathbf{P} = \bar{\Sigma} - \mathbf{KP}_z\mathbf{K}^\top$$

This step contains a few equations you have to take on faith, but you should be able to see how they relate to the linear Kalman filter equations. We convert the mean and covariance into measurement space, add the measurement error into the measurement covariance, compute the residual and kalman gain, compute the new state estimate as the old estimate plus the residual times the Kalman gain, and adjust the covariance for the information provided by the measurement. The linear algebra is slightly different from the linear Kalman filter, but the algorithm is the same Bayesian algorithm we have been implementing throughout the book.

This table comparing the Kalman filter equations with the UKF equations should help clarify the relationship between the filters.

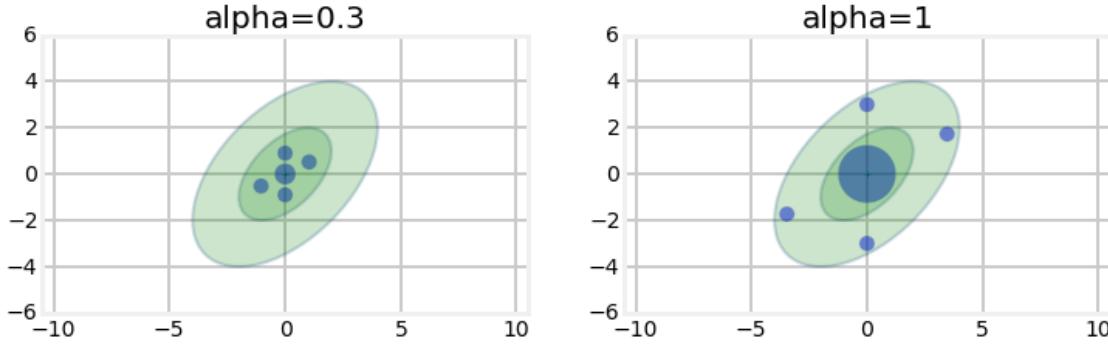
Kalman Filter	Unscented Kalman Filter
$\bar{\mathbf{x}} = \mathbf{Fx}$	$\bar{\mu} = \sum w^m \mathcal{Y}$
$\bar{\mathbf{P}} = \mathbf{FPF}^\top + \mathbf{Q}$	$\bar{\mathbf{P}} = \mathbf{FPF}^\top + \mathbf{Q}$
$\mathbf{y} = \mathbf{z} - \mathbf{H}\bar{\mathbf{x}}$	$\mathbf{y} = \mathbf{z} - \sum w^m h(\mathcal{Y})$
$\mathbf{S} = \mathbf{H}\bar{\mathbf{P}}\mathbf{H}^\top + \mathbf{R}$	$\mathbf{P}_z = \sum w^c (\mathcal{Z} - \bar{\mu})(\mathcal{Z} - \bar{\mu})^\top + \mathbf{R}$
$\mathbf{K} = \bar{\mathbf{P}}\mathbf{H}^\top \mathbf{S}^{-1}$	$\mathbf{K} = [\sum w^c (\chi - \mu)(\mathcal{Z} - \mu_z)^\top] \mathbf{P}_z^{-1}$
$\mathbf{x} = \bar{\mathbf{x}} + \mathbf{Ky}$	$\mathbf{x} = \bar{\mathbf{x}} + \mathbf{Ky}$
$\mathbf{P} = (\mathbf{I} - \mathbf{KH})\bar{\mathbf{P}}$	$\mathbf{P} = \bar{\Sigma} - \mathbf{KP}_z\mathbf{K}^\top$

## 10.5 Van der Merwe's Scaled Sigma Point Algorithm

There are several published ways for selecting the sigma points for the unscented Kalman filter, and you are free to invent your own. However, since 2005 or so research and industry have mostly settled on the version published by Rudolph Van der Merwe in his 2004 PhD dissertation [1] because it performs well with a variety of problems and it has a good tradeoff between performance and accuracy. It is a slight reformulation of the Scaled Unscented Transform published by Simon J. Julier [2].

Van der Merwe's formulation uses 3 parameters to control how the sigma points are distributed and weighted:  $\alpha$ ,  $\beta$ , and  $\kappa$ . Before we work through the equations, let's look at an example. I will plot the sigma points on top of a covariance ellipse showing the first and second standard deviations, and size the points based on the weights assigned to them.

In [9]: `ukf_internal.plot_sigma_points()`



3.91

We can see that the sigma points lie between the first and second deviation, and that the larger  $\alpha$  spreads the points out. Furthermore, the larger  $\alpha$  weighs the mean (center point) higher than the smaller  $\alpha$ , and weighs the rest of the sigma points less. This should fit our intuition - the further a point is from the mean the less we should weight it. We don't know how these weights and sigma points are selected yet, but the choices look reasonable.

### 10.5.1 Sigma Point Computation

Our first sigma point is always going to be the mean of our input. This is the sigma point displayed in the center of the ellipses in the diagram above. We will call this  $\mathbf{x}_0$ . So,

$$\mathbf{x}_0 = \mu$$

For notational convenience we define  $\lambda = \alpha^2(n + \kappa) - n$ , where  $n$  is the dimension of  $\mathbf{x}$ . Then the remaining sigma points are computed as

$$\mathbf{x}_i = \begin{cases} \mu + (\sqrt{(n + \lambda)\Sigma}) & \text{for } i=1 \dots n \\ \mu - (\sqrt{(n + \lambda)\Sigma})_{i-n} & \text{for } i=(n+1) \dots 2n \end{cases}$$

In other words, we scale the covariance matrix by a constant, take the square root of it, and then to ensure symmetry both add and subtract it from the mean. We will discuss how you take the square root of a matrix later.

### 10.5.2 Weight Computation

This formulation uses one set of weights for the means, and another set for the covariance. The weights for the mean of  $\mathbf{x}_0$  is computed as

$$W_0^m = \frac{\lambda}{n + \lambda}$$

The weight for the covariance of  $\mathbf{x}_0$  is

$$W_0^c = \frac{\lambda}{n + \lambda} + 1 - \alpha^2 + \beta$$

The weights for the rest of the sigma points  $\chi_1 \dots \chi_{2n}$  are the same for the mean and covariance. They are

$$W_i^m = W_i^c = \frac{1}{2(n + \lambda)} \quad i = 1..2n$$

It may not be obvious why this is ‘correct’, and indeed, it cannot be proven that this is ideal for all nonlinear problems. But you can see that we are choosing the sigma points proportional to the square root of the covariance matrix, and the square root of variance is standard deviation. So, the sigma points are spread roughly according to 1 standard deviation. However, there is an  $n$  term in there - the more dimensions there are the more the points will be spread out and weighed less.

### 10.5.3 Reasonable Choices for the Parameters

$\beta = 2$  is a good choice for Gaussian problems,  $\kappa = 3 - n$  where  $n$  is the dimension of  $\mathbf{x}$  is a good choice for  $\kappa$ , and  $0 \leq \alpha \leq 1$  is an appropriate choice for  $\alpha$ , where a larger value for  $\alpha$  spreads the sigma points further from the mean.

## 10.6 Using the UKF

Lets launch into solving some problems so you can gain confidence in how easy the UKF actually is. All of the equations are implemented in FilterPy. Later we will revisit how the UKF is implemented in Python.

Let’s start by solving a problem you already know how to do. Although the UKF was designed for nonlinear problems, it works fine on linear problems. In fact, it is guaranteed to get the same result as the linear Kalman filter for linear problems. We will write a solver for the linear problem of tracking using a constant velocity model in 2D. This will allows us to focus on what is the same (and most is the same!) and what is different with the UKF.

To design a linear Kalman filter you need to design the  $\mathbf{x}$ ,  $\mathbf{F}$ ,  $\mathbf{H}$ ,  $\mathbf{R}$ , and  $\mathbf{Q}$  matrices. We have done this many times already so let me present a design to you without a lot of comment. We want a constant velocity model, so we can define  $\mathbf{x}$  to be

$$\mathbf{x} = [x \quad \dot{x} \quad y \quad \dot{y}]^\top$$

With this ordering of state variables we can define our state transition model to be

$$\mathbf{F} = \begin{bmatrix} 1 & \Delta t & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & \Delta t \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

which implements the Newtonian equation

$$x_k = x_{k-1} + \dot{x}_{k-1} \Delta t$$

Our sensors provide position measurements but not velocity, so the measurement function is

$$\mathbf{H} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

Let’s say our sensor gives positions in meters with an error of  $1\sigma = 0.3$  meters in both the  $x$  and  $y$  coordinates. This gives us a measurement noise matrix of

$$\mathbf{R} = \begin{bmatrix} 0.3^2 & 0 \\ 0 & 0.3^2 \end{bmatrix}$$

Finally, let's assume that the process noise can be represented by the discrete white noise model - that is, that over each time period the acceleration is constant. We can use FilterPy's `Q_discrete_white_noise()` method to create this matrix for us, but for review the matrix is

$$\mathbf{Q} = \begin{bmatrix} \frac{1}{4}\Delta t^4 & \frac{1}{2}\Delta t^3 \\ \frac{1}{2}\Delta t^3 & \Delta t^2 \end{bmatrix} \sigma^2$$

Our implementation might look like this:

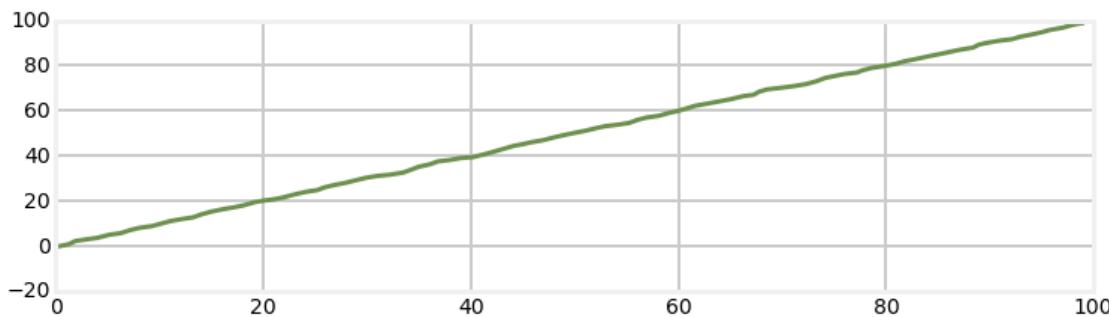
```
In [10]: from filterpy.kalman import KalmanFilter
        from filterpy.common import Q_discrete_white_noise
        from numpy import random
        from numpy.random import randn

        std_x, std_y = .3, .3
        dt = 1.0

        random.seed(1234)
        kf = KalmanFilter(4, 2)
        kf.x = np.array([0., 0., 0., 0.])
        kf.R = np.diag([std_x**2, std_y**2])
        kf.F = np.array([[1, dt, 0, 0],
                        [0, 1, 0, 0],
                        [0, 0, 1, dt],
                        [0, 0, 0, 1]])
        kf.H = np.array([[1, 0, 0, 0],
                        [0, 0, 1, 0]])

        kf.Q[0:2, 0:2] = Q_discrete_white_noise(2, dt=1, var=0.02)
        kf.Q[2:4, 2:4] = Q_discrete_white_noise(2, dt=1, var=0.02)

        zs = [np.array([i + randn()*std_x,
                       i + randn()*std_y]) for i in range(100)]
        xs, _, _, _ = kf.batch_filter(zs)
        plt.plot(xs[:, 0], xs[:, 2]);
```



This should hold no surprises for you. Now let's implement this filter as an Unscented Kalman filter. Again, this is purely for educational purposes; using a UKF for a linear filter confers no benefit.

The equations of the UKF are implemented for you with the `FilterPy` class `UnscentedKalmanFilter`; all you have to do is specify the matrices and the nonlinear functions  $f(x)$  and  $h(x)$ .  $f(x)$  implements the state transition function that is implemented by the matrix  $\mathbf{F}$  in the linear filter, and  $h(x)$  implements the measurement function implemented with the matrix  $\mathbf{H}$ . In nonlinear problems these functions are nonlinear, so we cannot use matrices to specify them.

For our linear problem we can define these functions to implement the linear equations. The function  $f(x)$  takes the signature `def f(x,dt)` and  $h(x)$  takes the signature `def h(x)`. Below is a reasonable implementation of these two functions. Each is expected to return a 1-D NumPy array containing the result.

```
In [11]: def f_cv(x, dt):
    """ state transition function for a
    constant velocity aircraft"""

    F = np.array([[1, dt, 0, 0],
                  [0, 1, 0, 0],
                  [0, 0, 1, dt],
                  [0, 0, 0, 1]])
    return np.dot(F, x)

def h_cv(x):
    return np.array([x[0], x[2]])
```

You need to tell the filter how to compute the sigma points and weights. We gave the Van der Merwe's scaled unscented transform version above, but there are many different choices. `FilterPy` uses a class named `SigmaPoints` which must implement two methods:

```
def sigma_points(self, x, P)
def weights(self)
```

`FilterPy` provides the class `MerweScaledSigmaPoints`, which implements the Van der Merwe algorithm. It derives from `SigmaPoints` and implements the aforementioned methods.

When you create the UKF you will pass in the  $f(x)$  and  $h(x)$  functions and the sigma point object as in this example:

```
from filterpy.kalman import MerweScaledSigmaPoints
from filterpy.kalman import UnscentedKalmanFilter as UKF

points = MerweScaledSigmaPoints(n=4, alpha=.1, beta=2., kappa=-1)
ukf = UKF(dim_x=4, dim_z=2, fx=f_cv, hx=h_cv, dt=dt, points=points)
```

The rest of the code is the same as for the linear kalman filter. Let's just implement it! I'll use the same measurements as used by the linear Kalman filter, and compute the standard deviation of the difference between the two solutions.

```
In [12]: from filterpy.kalman import UnscentedKalmanFilter as UKF

import numpy as np

sigmas = SigmaPoints(4, alpha=.1, beta=2., kappa=1.)
ukf = UKF(dim_x=4, dim_z=2, fx=f_cv,
```

```

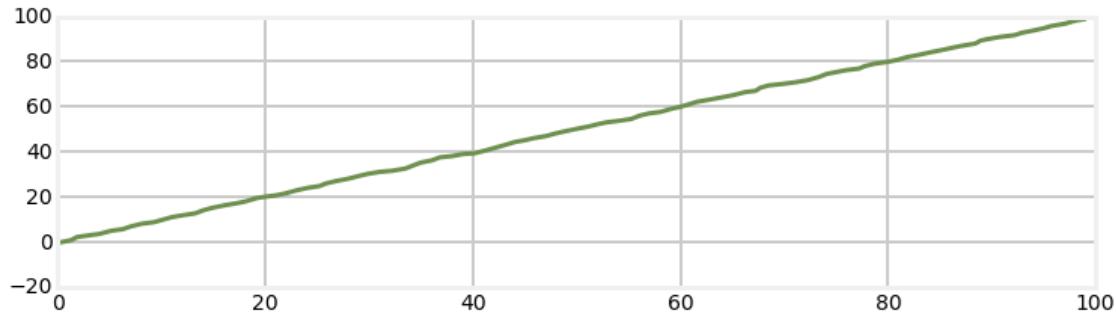
hx=h_cv, dt=dt, points=sigmas)
ukf.x = np.array([0., 0., 0., 0.])
ukf.R = np.diag([0.09, 0.09])
ukf.Q[0:2, 0:2] = Q_discrete_white_noise(2, dt=1, var=0.02)
ukf.Q[2:4, 2:4] = Q_discrete_white_noise(2, dt=1, var=0.02)

uxs = []
for z in zs:
    ukf.predict()
    ukf.update(z)
    uxs.append(ukf.x.copy())
uxs = np.array(uxs)

plt.plot(uxs[:, 0], uxs[:, 2])
print('UKF standard deviation {:.3f}'.format(np.std(uxs - xs)))

```

UKF standard deviation 0.013



This gave me a standard deviation of 0.013 which is quite small.

So far implementation of the UKF is not that different from the linear Kalman filter. Instead of implementing the state function and measurement function as the matrices  $\mathbf{F}$  and  $\mathbf{H}$  you supply functions  $\mathbf{f}(\mathbf{x})$  and  $\mathbf{h}(\mathbf{x})$ . The rest of the theory and implementation remains the same. Of course the FilterPy code for the UKF is quite different than the Kalman filter code, but from a designer's point of view the problem formulation and filter design is very similar.

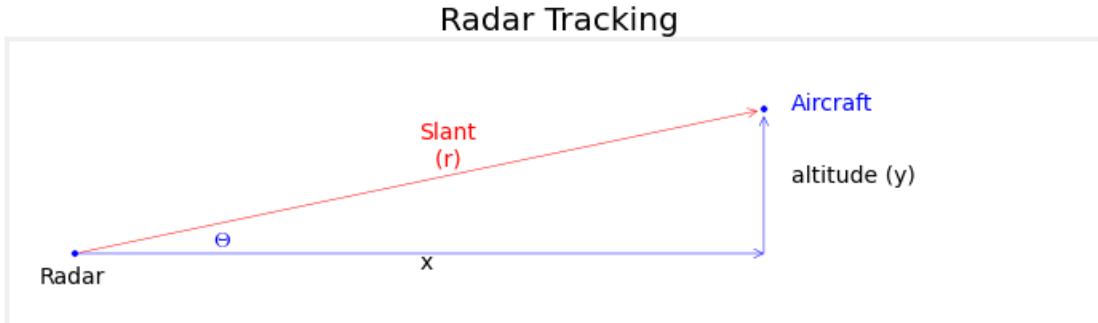
## 10.7 Tracking a Flying Airplane

Let's tackle our first nonlinear problem. I will keep the problem formulation very similar to the linear tracking problem above. For this problem we will write a filter to track a flying airplane using a stationary radar as the sensor. To keep the problem as close to the previous one as possible we will track in two dimensions. We will track one dimension on the ground and the altitude of the aircraft. The second dimension on the ground adds no difficulty or different information, so we can do this with no loss of generality.

Radar work by emitting a beam of radio waves and scanning for a return bounce. Anything in the beam's path will reflect some of the signal back to the radar. By timing how long it takes for the reflected signal to get back to the radar the system can compute the slant distance - the straight line distance from the radar installation to the object. We also get the bearing to the target. For this 2D problem that will be the angle above the ground plane.

We want to take the slant range measurement from the radar and compute the horizontal position (distance of aircraft from the radar measured over the ground) and altitude of the aircraft, as in the diagram below.

```
In [13]: import ekf_internal
ekf_internal.show_radar_chart()
```



We will assume that the aircraft is flying at a constant altitude. This leads us to design a three variable state vector:

$$\mathbf{x} = \begin{bmatrix} \text{distance} \\ \text{velocity} \\ \text{altitude} \end{bmatrix} = \begin{bmatrix} x \\ \dot{x} \\ y \end{bmatrix}$$

Our state transition function is linear

$$\bar{\mathbf{x}} = \begin{bmatrix} 1 & \Delta t & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \\ y \end{bmatrix}$$

which we can implement very much like the previous problem.

```
In [14]: def f_radar(x, dt):
    """ state transition function for a constant velocity
    aircraft with state vector [x, velocity, altitude] """
    F = np.array([[1, dt, 0],
                  [0, 1, 0],
                  [0, 0, 1]], dtype=float)
    return np.dot(F, x)
```

The next step is to design the measurement function. As in the linear Kalman filter the measurement function converts the filter's state into a measurement. For this problem we need to convert the position and velocity of the aircraft and into the bearing and range from the radar station.

To compute the range we use the Pythagorean theorem:

$$\text{range} = \sqrt{(x_{\text{ac}} - x_{\text{radar}})^2 + (z_{\text{ac}} - z_{\text{radar}})^2}$$

To compute the bearing we use the arctangent function.

$$\text{bearing} = \tan^{-1} \frac{z_{\text{ac}} - z_{\text{radar}}}{x_{\text{ac}} - x_{\text{radar}}}$$

As with the state transition function we need to define a Python function to compute this for the filter. I'll take advantage of the fact that a function can own a variable to store the radar's position. While this isn't necessary for this problem (we could hard code in the value), this gives the function more flexibility without requiring it to use a global variable.

```
In [15]: def h_radar(x):
    dx = x[0] - h_radar.radar_pos[0]
    dz = x[2] - h_radar.radar_pos[1]
    slant_range = math.sqrt(dx**2 + dz**2)
    bearing = math.atan2(dz, dx)
    return slant_range, bearing

h_radar.radar_pos = (0, 0)
```

There is a nonlinearity that we are not considering, the fact that angles are modular. Kalman filters operate by computing the differences between measurements. The difference between  $359^\circ$  and  $1^\circ$  is  $2^\circ$ , but a subtraction of the two values, as implemented by the filter, yields a value of  $358^\circ$ . This is exacerbated by the UKF which computes sums of weighted values in the unscented transform. For now we will place our sensors and targets in positions that avoid these nonlinear regions. Later in the chapter I will show you how to handle this problem.

We need to simulate the Radar station and the movement of the aircraft. By now this should be second nature for you, so I will present the code without much discussion.

```
In [16]: from numpy.linalg import norm
        from math import atan2

        class RadarStation(object):

            def __init__(self, pos, range_std, bearing_std):
                self.pos = np.asarray(pos)
                self.range_std = range_std
                self.bearing_std = bearing_std

            def reading_of(self, ac_pos):
                """ Returns (range, bearing) to aircraft.
                Bearing is in radians.
                """

                diff = np.subtract(ac_pos, self.pos)
                rng = norm(diff)
                brg = atan2(diff[1], diff[0])
                return rng, brg

            def noisy_reading(self, ac_pos):
                """ Compute range and bearing to aircraft with
                simulated noise"""
                rng, brg = self.reading_of(ac_pos)
```

```

        rng += randn() * self.range_std
        brg += randn() * self.bearing_std
        return rng, brg

    class ACSim(object):
        def __init__(self, pos, vel, vel_std):
            self.pos = np.asarray(pos, dtype=float)
            self.vel = np.asarray(vel, dtype=float)
            self.vel_std = vel_std

        def update(self, dt):
            """ Compute and returns next position. Incorporates
            random variation in velocity. """
            dx = self.vel*dt + (randn() * self.vel_std) * dt
            self.pos += dx
            return self.pos

```

Now let's put it all together. A military grade radar system can achieve 1 meter RMS range accuracy, and 1 mrad RMS for bearing [1]. We will assume a more modest 5 meter range accuracy, and  $0.5^\circ$  angular accuracy as this provides a more challenging data set for the filter.

The design of  $\mathbf{Q}$  requires some discussion. The state  $\mathbf{x}$  contains  $[x \quad \dot{x} \quad y]^\top$ . The first two elements are position (down range distance) and velocity, so we can use `Q_discrete_white_noise` noise to compute the values for the upper left hand side of  $\mathbf{Q}$ . The third element of  $\mathbf{x}$  is altitude, which we are assuming is independent of the down range distance. That leads us to a block design of  $\mathbf{Q}$  of:

$$\mathbf{Q} = \begin{bmatrix} \mathbf{Q}_x & 0 \\ 0 & Q_y \end{bmatrix}$$

I'll start with the aircraft positioned directly over the radar station flying at 100 m/s. A typical radar might update only once every 12 seconds and so we will use that for our update period.

```

In [17]: from filterpy.common import Q_discrete_white_noise
         import math
         from ukf_internal import plot_radar

         dt = 12. # 12 seconds between readings
         range_std = 5 # meters
         bearing_std = math.radians(0.5)
         ac_pos = (0., 1000.)
         ac_vel = (100., 0.)
         radar_pos = (0., 0.)
         h_radar.radar_pos = radar_pos

         points = SigmaPoints(n=3, alpha=.1, beta=2., kappa=0.)
         kf = UKF(3, 2, dt, fx=f_radar, hx=h_radar, points=points)

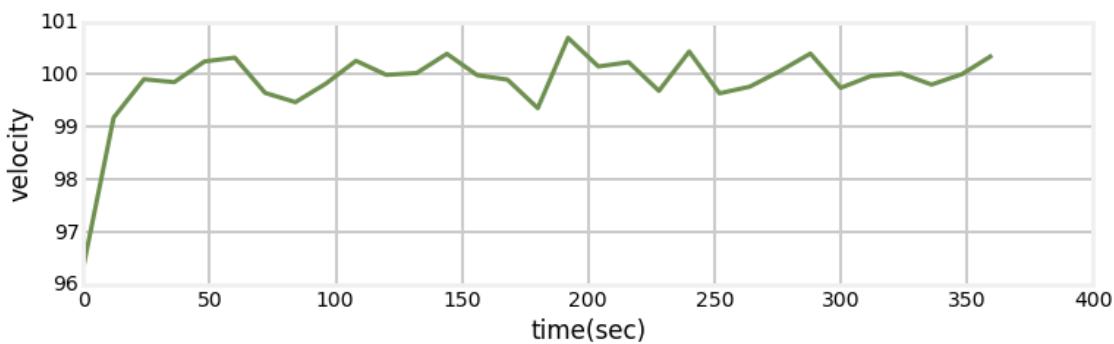
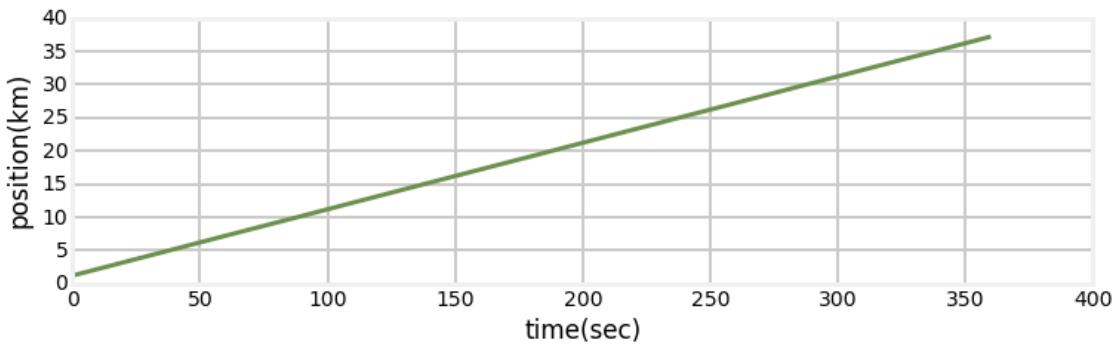
         kf.Q[0:2, 0:2] = Q_discrete_white_noise(2, dt=dt, var=0.1)
         kf.Q[2, 2] = 0.1

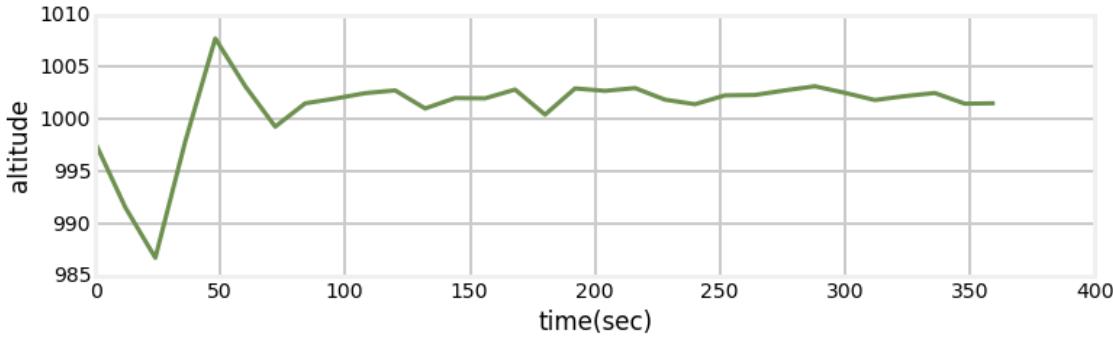
         kf.R = np.diag([range_std**2, bearing_std**2])
         kf.x = np.array([0., 90., 1100.])
         kf.P = np.diag([300**2, 30**2, 150**2])

```

```
random.seed(200)
pos = (0, 0)
radar = RadarStation(pos, range_std, bearing_std)
ac = ACSim(ac_pos, (100, 0), 0.02)

time = np.arange(0, 360 + dt, dt)
xs = []
for _ in time:
    ac.update(dt)
    r = radar.noisy_reading(ac.pos)
    kf.predict()
    kf.update([r[0], r[1]])
    xs.append(kf.x)
plot_radar(xs, time)
```





This may or may not impress you, but it impresses me! In the Extended Kalman filter chapter we will solve the same problem, but it will take significant amounts of mathematics to handle the same problem.

### 10.7.1 Tracking Manuevering Aircraft

The previous example produced extremely good results, but it also relied on an assumption of an aircraft flying at a constant speed with no change in altitude. I will relax that assumption by allowing the aircraft to change altitude. Here are the results if the aircraft starts climbing after one minute.

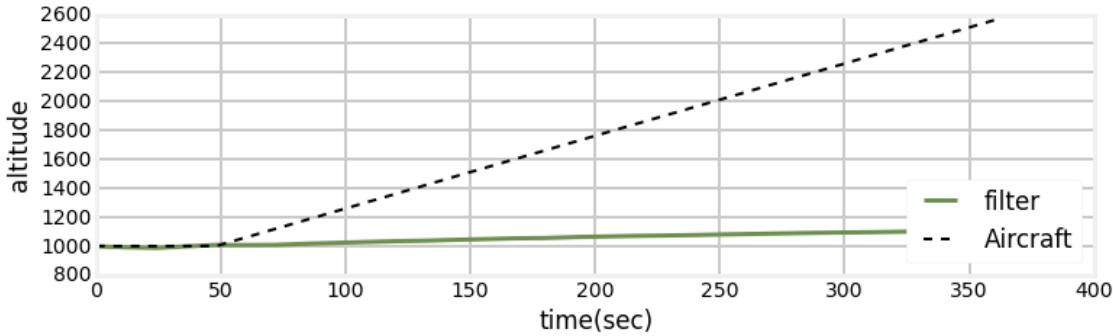
```
In [18]: from ukf_internal import plot_altitude

# reset aircraft position
kf.x = np.array([0., 90., 1100.])
kf.P = np.diag([300**2, 30**2, 150**2])
ac = ACSim(ac_pos, (100, 0), 0.02)

random.seed(200)
time = np.arange(0, 360 + dt, dt)
xs, ys = [], []
for t in time:
    if t >= 60:
        ac.vel[1] = 300/60 # 300 meters/minute climb
        ac.update(dt)
    r = radar.noisy_reading(ac.pos)
    ys.append(ac.pos[1])
    kf.predict()
    kf.update([r[0], r[1]])
    xs.append(kf.x)

plot_altitude(xs, time, ys)
print('Actual altitude: {:.1f}'.format(ac.pos[1]))
print('UKF altitude : {:.1f}'.format(xs[-1][2]))
```

Actual altitude: 2561.9  
 UKF altitude : 1107.2



The filter is completely unable to track the changing altitude. What do we have to change to allow the filter to track the aircraft?

I hope you answered add climb rate to the state, like so:

$$\mathbf{x} = \begin{bmatrix} \text{distance} \\ \text{velocity} \\ \text{altitude} \\ \text{climb rate} \end{bmatrix} = \begin{bmatrix} x \\ \dot{x} \\ y \\ \dot{y} \end{bmatrix}$$

This requires the following change to the state transition function, which is still linear.

$$\bar{\mathbf{x}} = \begin{bmatrix} 1 & \Delta t & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & \Delta t \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \\ y \\ \dot{y} \end{bmatrix}$$

The measurement function stays the same, but we will have to alter Q to account for the state dimensionality change.

```
In [19]: def f_cv_radar(x, dt):
    """ state transition function for a constant velocity
    aircraft"""
    F = np.array([[1, dt, 0, 0],
                  [0, 1, 0, 0],
                  [0, 0, 1, dt],
                  [0, 0, 0, 1]], dtype=float)
    return np.dot(F, x)

def cv_UKF(fx, hx, R_std):
    points = SigmaPoints(n=4, alpha=.1, beta=2., kappa=-1.)
    kf = UKF(4, len(R_std), dt, fx=fx, hx=hx, points=points)

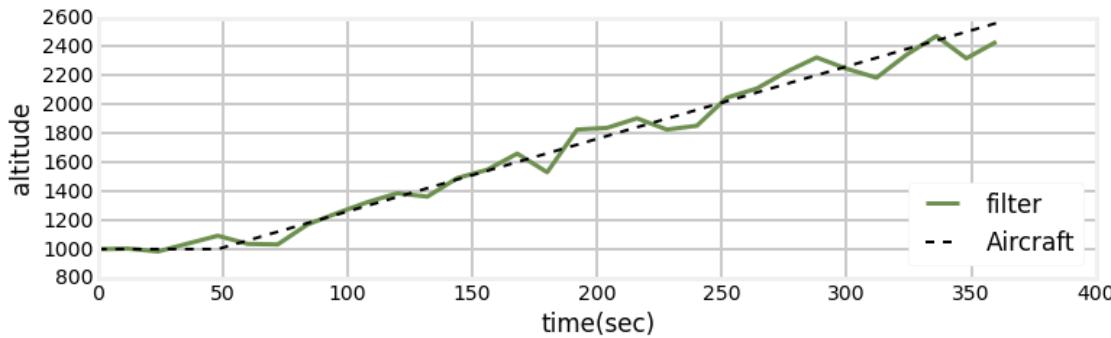
    kf.Q[0:2, 0:2] = Q_discrete_white_noise(2, dt=dt, var=0.1)
    kf.Q[2:4, 2:4] = Q_discrete_white_noise(2, dt=dt, var=0.1)
    kf.R = np.diag(R_std)
    kf.R = np.dot(kf.R, kf.R) # square to get variance
    kf.x = np.array([0., 90., 1100., 0.])
    kf.P = np.diag([300**2, 3**2, 150**2, 3**2])
    return kf
```

```
In [20]: random.seed(200)
ac = ACSim(ac_pos, (100, 0), 0.02)

kf_cv = cv_UKF(f_cv_radar, h_radar, R_std=[range_std, bearing_std])
time = np.arange(0, 360 + dt, dt)
xs, ys = [], []
for t in time:
    if t >= 60:
        ac.vel[1] = 300/60 # 300 meters/minute climb
    ac.update(dt)
    r = radar.noisy_reading(ac.pos)
    ys.append(ac.pos[1])
    kf_cv.predict()
    kf_cv.update([r[0], r[1]])
    xs.append(kf_cv.x)

plot_altitude(xs, time, ys)
print('Actual altitude: {:.1f}'.format(ac.pos[1]))
print('UKF altitude    : {:.1f}'.format(xs[-1][2]))
```

Actual altitude: 2561.9  
UKF altitude : 2432.9



We can see that a significant amount of noise has been introduced into the altitude, but we are now accurately tracking the altitude change.

### 10.7.2 Sensor Fusion

Now let's consider an example of sensor fusion. We have some type of Doppler system that produces a velocity estimate with 2 m/s RMS accuracy. I say "some type" because as with the radar I am not trying to teach you how to create an accurate filter for a Doppler system, where you have to account for the signal to noise ratio, atmospheric effects, the geometry of the system, and so on.

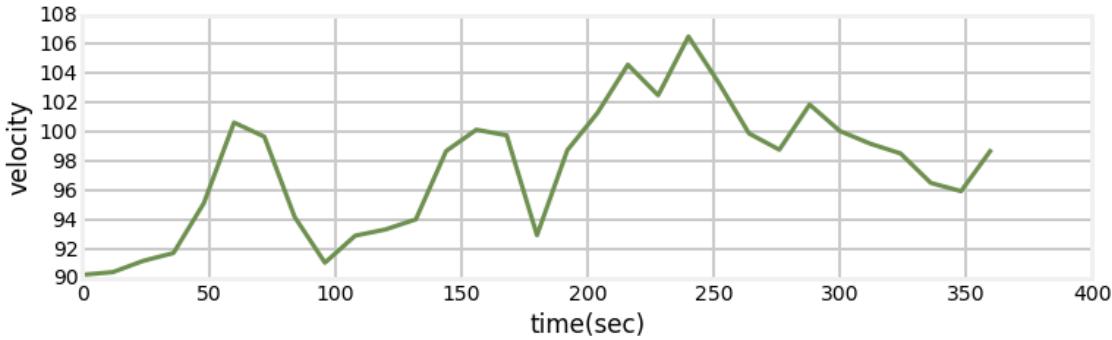
The accuracy of the radar system in the last examples allowed us to estimate velocities to within a m/s or so, so we have to degrade that accuracy to be able to easily see the effect of the sensor fusion. Let's change the range error to 500 meters and then compute the standard deviation of the computed velocity. I'll skip the first several measurements because the filter is converging during that time, causing artificially large deviations.

First, here is the code that computes the standard deviation of the filter without using the Doppler:

```
In [21]: range_std = 500.
bearing_std = math.degrees(0.5)
random.seed(200)
pos = (0, 0)
radar = RadarStation(pos, range_std, bearing_std)
ac = ACSim(ac_pos, (100, 0), 0.02)

kf_sf = cv_UKF(f_cv_radar, h_radar, R_std=[range_std, bearing_std])
time = np.arange(0, 360 + dt, dt)
xs = []
for _ in time:
    ac.update(dt)
    r = radar.noisy_reading(ac.pos)
    kf_sf.predict()
    kf_sf.update([r[0], r[1]])
    xs.append(kf_sf.x)

xs = np.asarray(xs)
plot_radar(xs, time, plot_x=False, plot_vel=True, plot_alt=False)
print('Velocity std {:.1f} m/s'.format(np.std(xs[10:, 1])))
```



Velocity std 3.4 m/s

To include the doppler we need to add the velocity in  $x$  and  $y$  into the measurement. The `ACSim` class stores the velocity in the data member `vel`. To perform the Kalman filter update we just need to call `update` with a list containing the bearing, distance, and velocity in  $x$  and  $y$ :

$$z = [slant\_range, bearing, \dot{x}, \dot{y}]$$

The measurement contains four values so the measurement function also needs to return four values. The slant range and bearing will be computed as before, but we do not need to compute the velocity in  $x$  and  $y$  as they are provided in the state estimate.

```
In [22]: def h_vel(x):
    dx = x[0] - h_vel.radar_pos[0]
    dz = x[2] - h_vel.radar_pos[1]
    slant_range = math.sqrt(dx**2 + dz**2)
    bearing = math.atan2(dz, dx)
    return slant_range, bearing, x[1], x[3]
```

With that implemented we can implement and execute our filter.

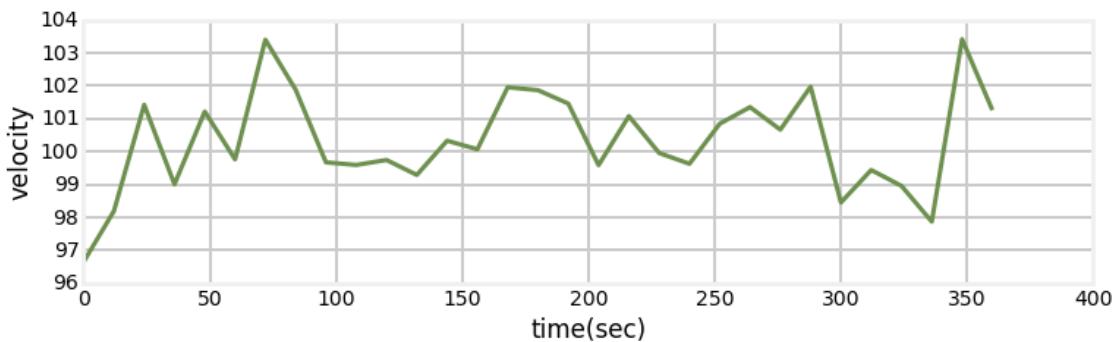
```
In [23]: h_radar.radar_pos = (0, 0)
h_vel.radar_pos = (0, 0)

range_std = 500.
bearing_std = math.degrees(0.5)
vel_std = 2.

random.seed(200)
ac = ACSim(ac_pos, (100, 0), 0.02)
radar = RadarStation((0, 0), range_std, bearing_std)

kf_sf2 = cv_UKF(f_cv_radar, h_vel,
                 R_std=[range_std, bearing_std, vel_std, vel_std])

time = np.arange(0, 360 + dt, dt)
xs = []
for t in time:
    ac.update(dt)
    r = radar.noisy_reading(ac.pos)
    # simulate the doppler velocity reading
    vx = ac.vel[0] + randn()*vel_std
    vz = ac.vel[1] + randn()*vel_std
    kf_sf2.predict()
    kf_sf2.update([r[0], r[1], vx, vz])
    xs.append(kf_sf2.x)
xs = np.asarray(xs)
plot_radar(xs, time, plot_x=False, plot_vel=True, plot_alt=False)
print('Velocity std {:.1f} m/s'.format(np.std(xs[10:,1])))
```



Velocity std 1.3 m/s

By incorporating the velocity sensor we were able to reduce the standard deviation from 3.5 m/s to 1.3 m/s. Sensor fusion is a large topic, and this is a rather simplistic implementation. In a typical navigation problem we have sensors that provide complementary information. For example, a GPS might provide somewhat accurate position updates once a second with poor velocity estimation while an inertial system might provide

very accurate velocity updates at 50Hz but terrible position estimates. The strengths and weaknesses of each sensor are orthogonal to each other. This leads to something called the Complementary filter, which uses the high update rate of the inertial sensor with the position accurate but slow estimates of the GPS to produce very high rate yet very accurate position and velocity estimates.

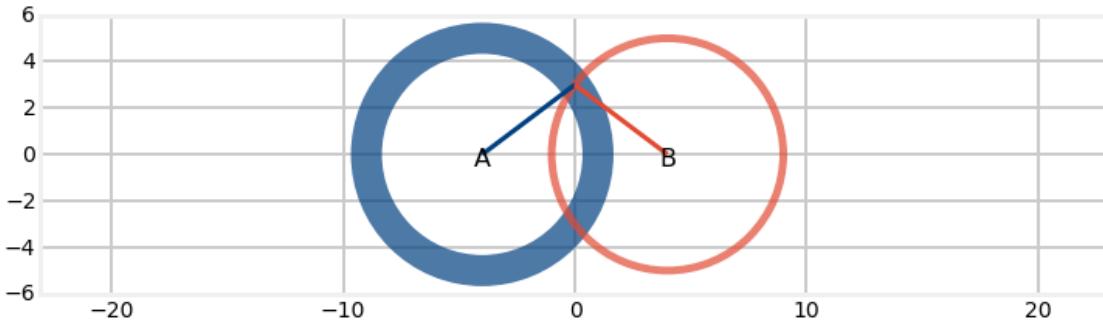
### 10.7.3 Multiple Position Sensors

The last sensor fusion problem was a toy example. Let's tackle a problem that is not so toy-like. Before GPS ships and aircraft navigated via various range and bearing systems such as VOR, LORAN, TACAN, DME, and so on. I do not intend to cover the intricacies of these systems - Wikipedia will fill in the basics if you are interested. In general these systems are beacons that allow you to extract the range and/or bearing to the beacon. For example, an aircraft might have two VOR receivers. The pilot tunes each receiver to a different VOR station. Each VOR receiver displays what is called the “radial” - the direction from the VOR station on the ground to the aircraft. Using a chart they can extend these two radials - the intersection point is the position of the aircraft.

That is a very manual approach with low accuracy. We can use a Kalman filter to filter the data and produce far more accurate position estimates. Let's work through that.

The problem is as follows. Assume we have two sensors, each which provides a bearing only measurement to the target, as in the chart below. The width of the circle is proportional to the amount of sensor noise.

In [24]: `ukf_internal.show_two_sensor_bearing()`



We can compute the bearing between a sensor and the target as follows:

```
def bearing(sensor, target):
    return math.atan2(target[1] - sensor[1], target[0] - sensor[0])
```

The filter receives a vector of 2 measurements during each update, one for each sensor. We can implement that as:

```
def measurement(A_pos, B_pos, pos):
    angle_a = bearing(A_pos, pos)
    angle_b = bearing(B_pos, pos)
    return [angle_a, angle_b]
```

In [25]: `from filterpy.kalman import UnscentedKalmanFilter as UKF
from filterpy.common import Q_discrete_white_noise`

```

sa_pos = [-400, 0]
sb_pos = [400, 0]

def bearing(sensor, target_pos):
    return math.atan2(target_pos[1] - sensor[1],
                      target_pos[0] - sensor[0])

def measurement(A_pos, B_pos, pos):
    angle_a = bearing(A_pos, pos)
    angle_b = bearing(B_pos, pos)
    return [angle_a, angle_b]

def fx_VOR(x, dt):
    x[0] += x[1] * dt
    x[2] += x[3] * dt
    return x

def hx_VOR(x):
    # measurement to A
    pos = (x[0], x[2])
    return measurement(sa_pos, sb_pos, pos)

def moving_target_filter(pos, std_noise, Q, dt=0.1, kappa=0.0):
    points = SigmaPoints(n=4, alpha=.1, beta=2., kappa=kappa)
    f = UKF(dim_x=4, dim_z=2, dt=dt,
            hx=hx_VOR, fx=fx_VOR, points=points)
    f.x = np.array([pos[0], 1., pos[1], 1.])

    q = Q_discrete_white_noise(2, dt, Q)
    f.Q[0:2, 0:2] = q
    f.Q[2:4, 2:4] = q
    f.R *= std_noise**2
    f.P *= 1000
    return f

def plot_straight_line_target(f, std_noise):
    xs = []
    txs = []

    for i in range(300):
        target_pos[0] += 1 + randn()*0.0001
        target_pos[1] += 1 + randn()*0.0001
        txs.append((target_pos[0], target_pos[1]))

        z = measurement(sa_pos, sb_pos, target_pos)
        z[0] += randn() * std_noise
        z[1] += randn() * std_noise

        f.predict()
        f.update(z)
        xs.append(f.x)

    xs = np.asarray(xs)

```

```

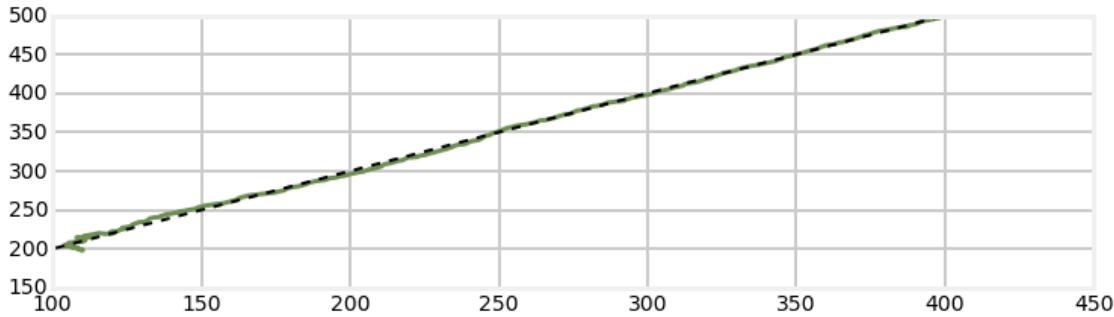
txs = np.asarray(txs)

plt.plot(xs[:, 0], xs[:, 2])
plt.plot(txs[:, 0], txa[:, 1], ls='--', lw=2, c='k')
plt.show()

np.random.seed(123)
target_pos = [100, 200]

std_noise = math.radians(0.5)
f = moving_target_filter(target_pos, std_noise, Q=1.0)
plot_straight_line_target(f, std_noise)

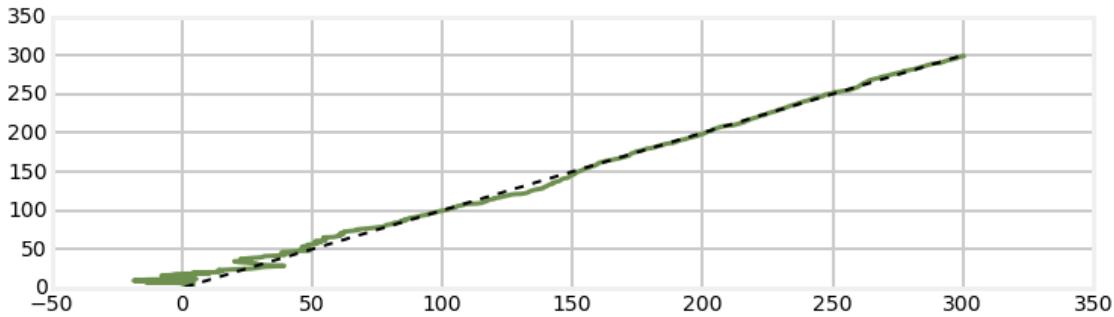
```



This looks quite good to me. There is a very large error at the beginning of the track, but the filter is able to settle down and start producing good data.

Let's revisit the nonlinearity of the angles. I will position the target between the two sensors with the same y-coordinate. This will cause a nonlinearity in the computation of the sigma means and the residuals because the mean angle will be near zero. As the angle goes below 0 the measurement function will compute a large positive angle of nearly  $2\pi$ . The residual between the prediction and measurement will thus be very large, nearly  $2\pi$  instead of nearly 0. This makes it impossible for the filter to perform accurately, as seen in the example below.

```
In [26]: target_pos = [0, 0]
f = moving_target_filter(target_pos, std_noise, Q=1.0)
plot_straight_line_target(f, std_noise)
```



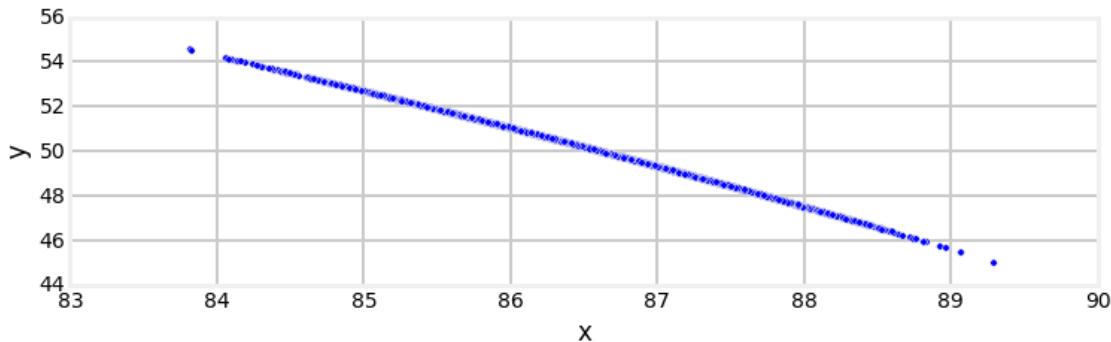
This is unacceptable performance. FilterPy's UKF code allows you to specify a function which computes the residuals in cases of nonlinear behavior like this,. The final example in this chapter explains how to handle this problem.

Finally, let's discuss the sensor fusion method that we used. We used the measurement from each bearing and had the Kalman filter's equations compute the world position from the measurements. This is equivalent to doing a weighted least squares solution. In the [Kalman Filter Math](#) chapter we discussed the limited accuracy of such a scheme and introduced an [Iterative Least Squares \(ILS\)](#) algorithm to produce greater accuracy. If you wanted to use that scheme you would write an ILS algorithm to solve the geometry of your problem, and pass the result of that calculation into the `update()` method as the measurement to be filtered. This imposes on you the need to compute the correct noise matrix for this computed positions, which may not be trivial. For example, the ILS is probably the most common algorithm used to turn GPS pseudoranges into a position, but the literature discusses a number of alternatives.

## 10.8 Effects of Sensor Error and Geometry

The geometry of the sensors relative to the tracked object imposes a physical limitation that can be extremely difficult to deal with when designing filters. If the radials of the VOR stations are nearly parallel to each other than a very small angular error translates into a very large distance error. What is worse, this behavior is nonlinear - the error in the x-axis vs the y-axis will vary depending on the actual bearing. For example, here is a scatter plot that shows the error distribution for a  $1^\circ$  standard deviation in the bearing measurement for a  $30^\circ$  bearing.

In [27]: `ukf_internal.plot_scatter_of_bearing_error()`



## 10.9 Exercise: Explain Filter Performance

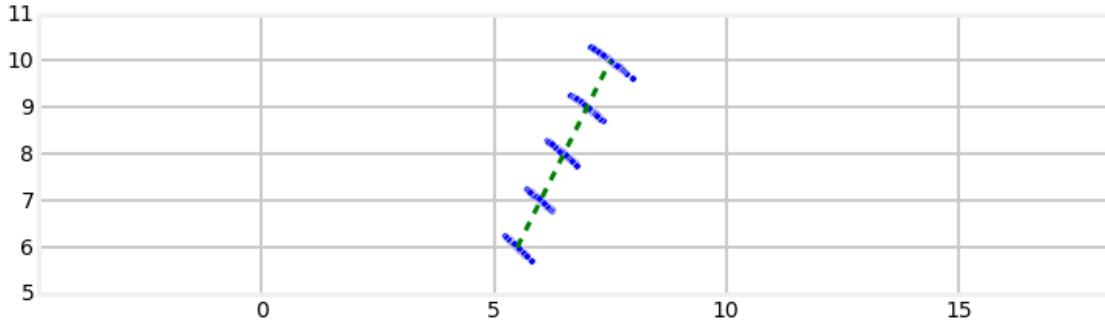
We can see that for even very small angular errors the (x, y) positional errors are very large. Explain how we got such relatively good performance out of the UKF in the target tracking problems above. Answer this for the one sensor problem as well as the multiple sensor problem.

### 10.9.1 Solution

This is **very** important to understand. Try very hard to answer this before reading the answer below. If you cannot answer this you really need to revisit some of the earlier Kalman filter material in the [Multidimensional Kalman Filter](#) chapter.

There are several factors contributing to our success. First, let's consider the case of having only one sensor. Any single measurement has an extreme range of possible positions. But, our target is moving, and the UKF is taking that into account. Let's plot the results of several measurements taken in a row for a moving target to form an intuition.

In [28]: `ukf_internal.plot_scatter_moving_target()`

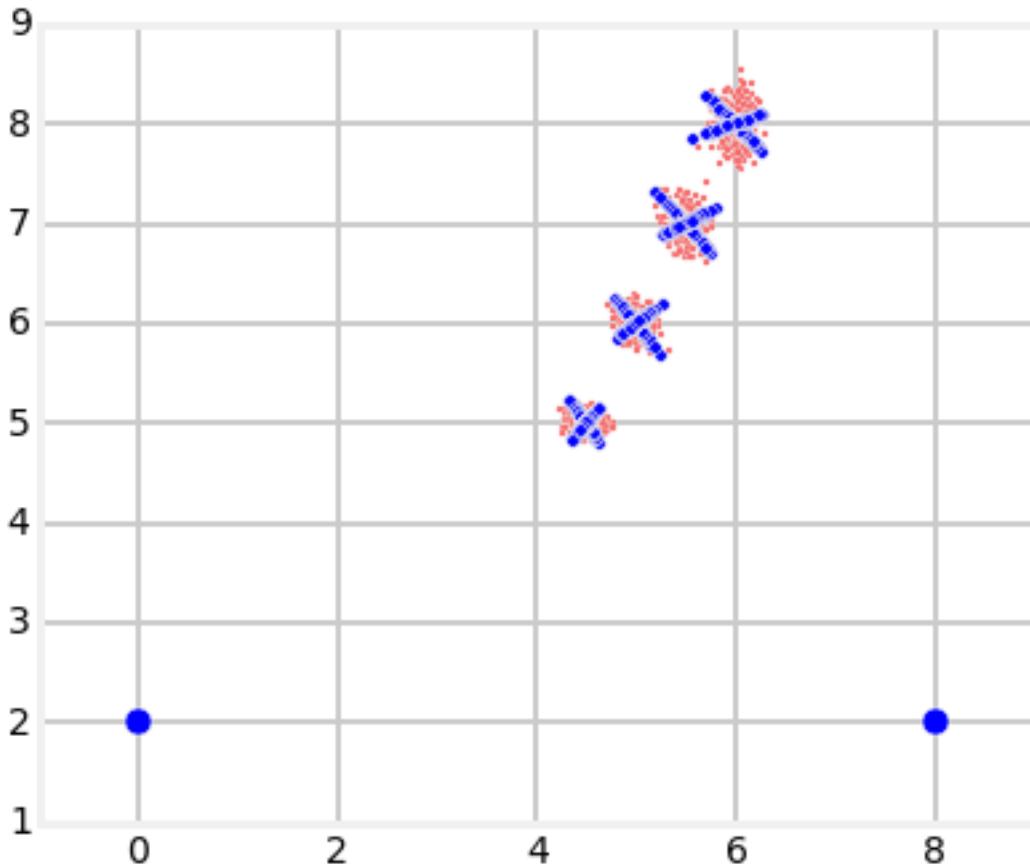


Each individual measurement has a very large position error. However, a plot of successive measurements shows a clear trend - the target is obviously moving towards the upper right. When a Kalman filter computes the Kalman gain it takes the distribution of errors into account by using the measurement function. In this example the error lies on an approximately  $45^\circ$  degree line, so the filter will discount errors in that direction. On the other hand, there is almost no error in measurement orthogonal to that, and again the Kalman gain will be taking that into account. The end result is the track can be computed even with this type of noise distribution. This graph makes it look easy because we have plotted 100 possible measurements for each position update. This makes the aircraft's movement obvious. In contrast, the Kalman filter only gets one measurement per update. Therefore the filter will not be able to generate as good a fit as the dotted green line implies.

The next interaction we must think about is the fact that the bearing gives us no distance information. Suppose we set the initial position to be 1,000 kilometers away from the sensor (vs the actual distance of 7.07 km) and make  $\mathbf{P}$  very small. At that distance a  $1^\circ$  error translates into a positional error of 17.5 km. The KF would never be able to converge onto the actual target position because the filter is incorrectly very certain about its position estimates and because there is no distance information provided in the measurements.

Now let's consider the effect of adding a second sensor. Again, I will start by plotting the measurements so we can form an intuition about the problem.

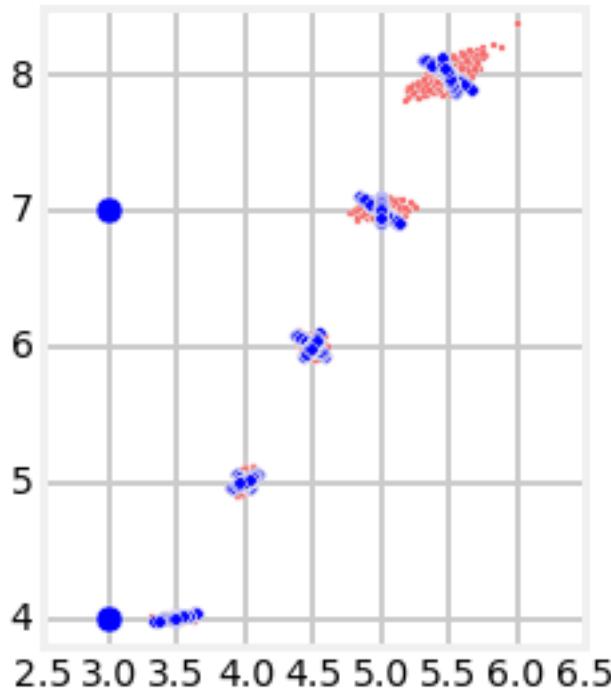
In [29]: `with figsize(y=5.):  
 ukf_internal.plot_iscts_two_sensors()`



I placed the sensors nearly orthogonal to the target’s initial position so we get these lovely ‘x’ shape intersections. I then computed the  $(x, y)$  coordinate corresponding to the two noisy bearing measurements and plotted them with red dots to show the distribution of the noisy measurements in  $x$  and  $y$ . We can see how the errors in  $x$  and  $y$  change as the target moves by the shape the scattered red dots make - as the target gets further away from the sensors, but nearer the  $y$  coordinate of sensor B the shape becomes strongly elliptical.

Next I will alter the sensor positions and rerun the simulation. Here the shape of the errors in  $x$  and  $y$  changes radically as the target position changes relative to the two sensors.

```
In [30]: with figsize(y=4.):
    ukf_internal.plot_iscts_two_sensors_changed_sensors()
```



## 10.10 Implementation of the UKF

Now let's implement the UKF equations with Python. As you have seen FilterPy already implements this code for you, but it is instructive to learn how to go from equations to code. Plus, if you encounter a problem and need to debug your code you will likely need to step through the FilterPy code with the debugger.

Implementing the UKF is quite straightforward. First, let's write the code to compute the mean and covariance given the sigma points.

We will store the sigma points and weights in matrices, like so:

$$\text{weights} = [w_1 \quad w_2 \quad \dots \quad w_{2n+1}]$$

$$\text{sigmas} = \begin{bmatrix} \mathcal{X}_{0,0} & \mathcal{X}_{0,1} & \mathcal{X}_{0,2} \\ \mathcal{X}_{1,0} & \mathcal{X}_{1,1} & \mathcal{X}_{1,2} \\ \vdots & \vdots & \vdots \\ \mathcal{X}_{2n+1,0} & \mathcal{X}_{2n+1,1} & \mathcal{X}_{2n+1,2} \end{bmatrix}$$

In other words, each column contains the  $2n + 1$  sigma points for one dimension in our problem. The  $0^{th}$  sigma point is always the mean, so first row of the sigmas is the mean of each of our state variables. The second through nth row contains the  $\mu + \sqrt{(n + \lambda)\Sigma}$  terms, and the  $n + 1$  to  $2n$  rows contains the  $\mu - \sqrt{(n + \lambda)\Sigma}$  terms. The choice to store the sigmas in row-column vs column row format is somewhat arbitrary; my choice makes the rest of the code a bit easier to code as I can refer to the ith sigma point for all dimensions as `sigmas[i]`.

### 10.10.1 Weights

Computing the weights in numpy is extremely simple. Recall that the Van der Merwe scaled sigma point implementation states:

$$\begin{aligned}\lambda &= \alpha^2(n + \kappa) - n \\ W_0^m &= \frac{\lambda}{n + \lambda} \\ W_0^c &= \frac{\lambda}{n + \lambda} + 1 - \alpha^2 + \beta \\ W_i^m = W_i^c &= \frac{1}{2(n + \lambda)} \quad i = 1..2n\end{aligned}$$

Our code for these is:

```
lambda_ = alpha**2 * (n + kappa) - n
Wc = np.full(2*n + 1, 1. / (2*(n+lambda_)))
Wm = np.full(2*n + 1, 1. / (2*(n+lambda_)))
Wc[0] = lambda_ / (n + lambda_) + (1. - alpha**2 + beta)
Wm[0] = lambda_ / (n + lambda_)
```

I use the underscore in `lambda_` because `lambda` is a reserved word in Python. A trailing underscore is the Pythonic workaround

### 10.10.2 Sigma Points

The equations for the sigma points are:

$$\begin{cases} \mathcal{X}_0 = \mu \\ \mathcal{X}_i = \mu + \left[ \sqrt{(n + \lambda)\Sigma} \right]_i, & \text{for } i=1..n \\ \mathcal{X}_i = \mu - \left[ \sqrt{(n + \lambda)\Sigma} \right]_{i-n}, & \text{for } i=(n+1)..2n \end{cases}$$

The Python is not difficult once we understand the  $\left[ \sqrt{(n + \lambda)\Sigma} \right]_i$  term.

The term  $\left[ \sqrt{(n + \lambda)\Sigma} \right]_i$  is a matrix because  $\Sigma$  is a matrix. The subscript  $i$  is choosing the column vector of the matrix. What is the square root of a matrix? There is no unique definition. One definition is that the square root of a matrix  $\Sigma$  is the matrix  $S$  that, when multiplied by itself, yields  $\Sigma$ .

$$\text{if } \Sigma = SS^\top$$

$$\text{then } S = \sqrt{\Sigma}$$

However there is an alternative definition, and we will chose it because it has numerical properties that makes it much easier for us to compute its value. We can define the square root as the matrix  $S$ , which when multiplied by its transpose, returns  $\Sigma$ :

$$\Sigma = SS^\top$$

This method is frequently chosen in computational linear algebra because this expression is easy to compute using something called the Cholesky decomposition [3]. SciPy provides this with the

`scipy.linalg.cholesky()` method. If your language of choice is Fortran, C, or C++, libraries such as LAPACK provide this routine. Matlab provides `chol()`.

```
sigmas = np.zeros((2*n+1, n))
U = scipy.linalg.cholesky((n+kappa)*P)

sigmas[0] = X
for k in range (n):
    sigmas[k+1] = X + U[k]
    sigmas[n+k+1] = X - U[k]
```

Now let's implement the unscented transform. Recall the equations

$$\mu = \sum_i w_i^m \mathcal{X}_i$$

$$\Sigma = \sum_i w_i^c (\mathcal{X}_i - \mu)(\mathcal{X}_i - \mu)^\top$$

We implement the sum of the means with

```
x = np.dot(Wm, sigmas)
```

If you are not a heavy user of NumPy this may look foreign to you. NumPy is not just a library that make linear algebra possible; under the hood it is written in C to achieve much faster speeds than Python can reach. A typical speedup is 20x to 100x. To get that speedup we must avoid using for loops, and instead use NumPy's built in functions to perform calculations. So, instead of writing a for loop to compute the sum, we call the built in `numpy.dot(x, y)` method. If passed a 1D array and a 2D array it will compute the sum of inner products:

```
In [31]: a = np.array([10, 100])
b = np.array([[1, 2, 3],
              [4, 5, 6]])
np.dot(a,b)
```

```
Out[31]: array([410, 520, 630])
```

All that is left is to compute  $\mathbf{P} = \sum_i w_i(\mathcal{X}_i - \mu)(\mathcal{X}_i - \mu)^\top + \mathbf{Q}$ :

```
kmax, n = sigmas.shape
P = zeros((n, n))
for k in range(kmax):
    y = sigmas[k] - x
    P += Wc[k] * np.outer(y, y)
P += Q
```

This introduces another feature of NumPy. The state variable `x` is one dimensional, as is `sigmas[k]`, so the difference `sigmas[k] - x` is also one dimensional. NumPy will not compute the transpose of a 1-D array; it considers the transpose of `[1,2,3]` to be `[1,2,3]`. So we call the function `np.outer(y,y)` which computes the value of  $\mathbf{y}\mathbf{y}^\top$  for a 1D array `y`. An alternative implementation could be:

```
y = (sigmas[k] - x).reshape(kmax, 1) # convert into 2D array
P += Wc[K] * np.dot(y, y.T)
```

However, that code is slower and not idiomatic, and we will not use it.

### 10.10.3 Predict Step

For the predict step, we will generate the weights and sigma points as specified above. We pass each sigma point through the function  $f$ .

$$\mathcal{Y} = f(\chi)$$

Then we compute the predicted mean and covariance using the unscented transform. In the code below you can see that I am assuming that this is a method in a class that stores the various matrices and vectors needed by the filter.

```
In [32]: def predict(self, sigma_points_fn):
    """ Performs the predict step of the UKF. On return,
    self.xp and self.Pp contain the predicted state (xp)
    and covariance (Pp). 'p' stands for prediction.
    """
    # calculate sigma points for given mean and covariance
    sigmas = sigma_points_fn(self.x, self.P)

    for i in range(self._num_sigmas):
        self.sigmas_f[i] = self.fx(sigmas[i], self._dt)

    self.xp, self.Pp = unscented_transform(
        self.sigmas_f, self.W, self.W, self.Q)
```

### 10.10.4 Update Step

The update step converts the sigmas into measurement space via the  $h(x)$  function.

$$\mathcal{Z} = h(\mathcal{Y})$$

The mean and covariance of those points is computed with the unscented transform. The residual and Kalman gain is then computed. The cross variance is computed as:

$$\mathbf{P}_{xz} = \sum w^c (\chi - \mu)(\mathcal{Z} - \mu_z)^\top$$

Finally, we compute the new state estimate using the residual and Kalman gain:

$$K = \mathbf{P}_{xz} \mathbf{P}_z^{-1} \mathbf{x} = \bar{\mathbf{x}} + \mathbf{K}\mathbf{y}$$

and the new covariance is computed as:

$$\mathbf{P} = \bar{\mathbf{P}} - \mathbf{K}\mathbf{P}_z\mathbf{K}^\top$$

This function can be implemented as follows, assuming it is a method of a class that stores the necessary matrices and data.

```
In [33]: def update(self, z):
    # rename for readability
    sigmas_f = self.sigmas_f
```

```

sigmas_h = self.sigmas_h

# transform sigma points into measurement space
for i in range(self._num_sigmas):
    sigmas_h[i] = self.hx(sigmas_f[i])

# mean and covariance of prediction passed through UT
zp, Pz = unscented_transform(sigmas_h, self.W, self.W, self.R)

# compute cross variance of the state and the measurements
Pxz = np.zeros((self._dim_x, self._dim_z))
for i in range(self._num_sigmas):
    Pxz += self.W[i] * np.outer(sigmas_f[i] - self.xp,
                                sigmas_h[i] - zp)

K = dot(Pxz, inv(Pz)) # Kalman gain

self.x = self.xp + dot(K, z-zp)
self.P = self.Pp - dot3(K, Pz, K.T)

```

### 10.10.5 FilterPy's Implementation

FilterPy has generalized the code somewhat. You can specify different sigma point algorithms, specify how to compute the residual of the state variables (you can not subtract angles because they are modular), provide a matrix square root function, and more. See the help for details.

<https://filterpy.readthedocs.org/#unscented-kalman-filter>

## 10.11 Batch Processing

The Kalman filter is designed as a recursive algorithm - as new measurements come in we immediately create a new estimate. But it is very common to have a set of data that have been already collected which we want to filter. Kalman filters can always be run in a batch mode, where all of the measurements are filtered at once.

Collect your measurements into an array or list.

```
zs = read_altitude_from_csv()
```

Then call the `batch_filter()` method.

```
Xs, Ps = ukf.batch_filter(zs)
```

The function takes the list/array of measurements, filters it, and returns an array of state estimates (Xs) and covariance matrices (Ps) for the entire data set.

Here is a complete example drawing from the radar tracking problem above.

```
In [34]: dt = 12. # 12 seconds between readings
range_std = 5 # meters
bearing_std = math.radians(0.5)

ac_pos = (0., 1000.)
```

```

ac_vel = (100., 0.)
radar_pos = (0., 0.)
h_radar.radar_pos = radar_pos

points = SigmaPoints(n=3, alpha=.1, beta=2., kappa=0.)
kf = UKF(3, 2, dt, fx=f_radar, hx=h_radar, points=points)

kf.Q[0:2 ,0:2] = Q_discrete_white_noise(2, dt=dt, var=0.1)
kf.Q[2, 2] = 0.1

kf.R = np.diag([range_std**2, bearing_std**2])
kf.x = np.array([0., 90., 1100.])
kf.P = np.diag([300**2, 30**2, 150**2])

radar = RadarStation((0, 0), range_std, bearing_std)
ac = ACSim(ac_pos, (100, 0), 0.02)

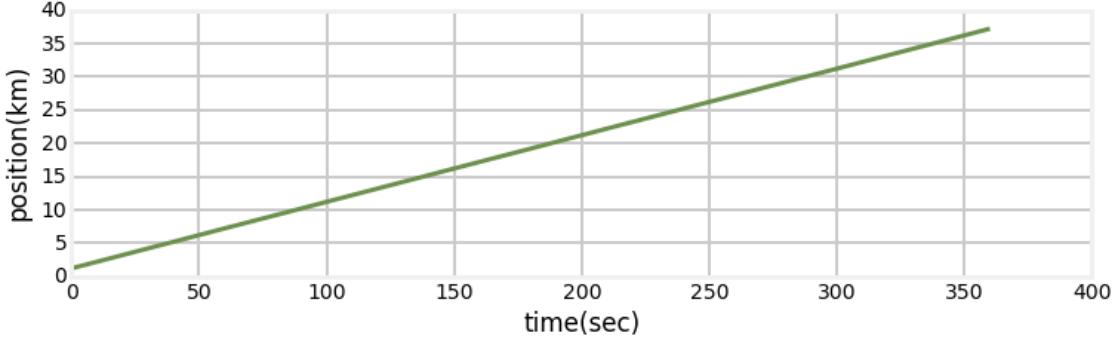
random.seed(200)

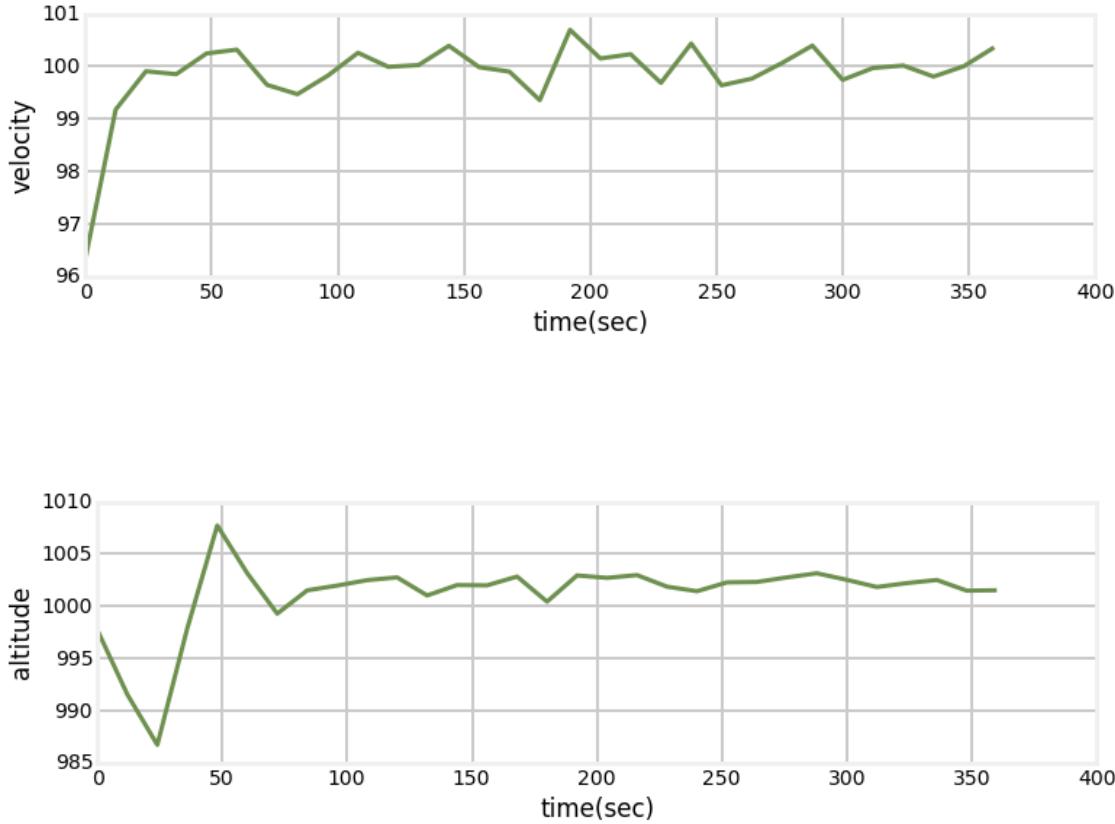
t = np.arange(0, 360 + dt, dt)
n = len(t)

zs = []
for i in range(len(t)):
    ac.update(dt)
    r = radar.noisy_reading(ac.pos)
    zs.append([r[0], r[1]])

xs, covs = kf.batch_filter(zs)
ukf_internal.plot_radar(xs, t)

```





## 10.12 Smoothing the Results

Let's assume that we are tracking a car. Suppose we get a noisy measurement that implies that the car is starting to turn to the left, but the state function has predicted that the car is moving straight. The Kalman filter has no choice but to move the state estimate somewhat towards the noisy measurement, as it cannot judge whether this is just a particularly noisy measurement or the true start of a turn.

However, if we are collecting data and post-processing it we have measurements after the questionable one that informs us if a turn was made or not. Suppose the subsequent measurements all continue turning left. We can then be sure that the measurement was not very noisy, but instead a turn was initiated.

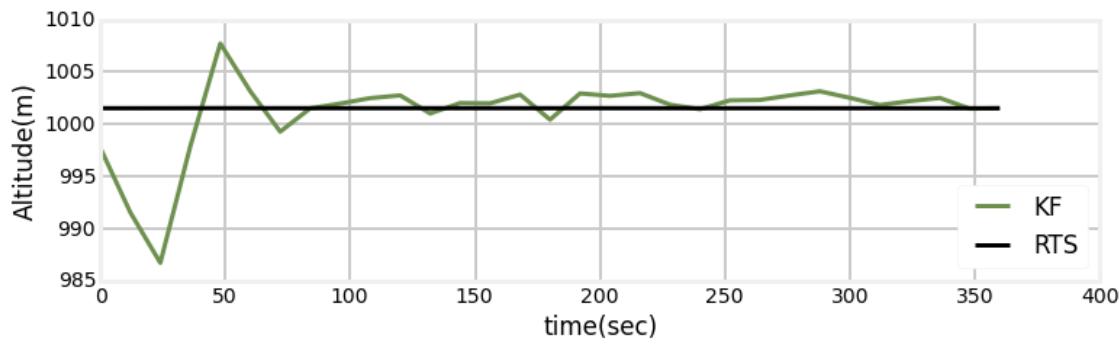
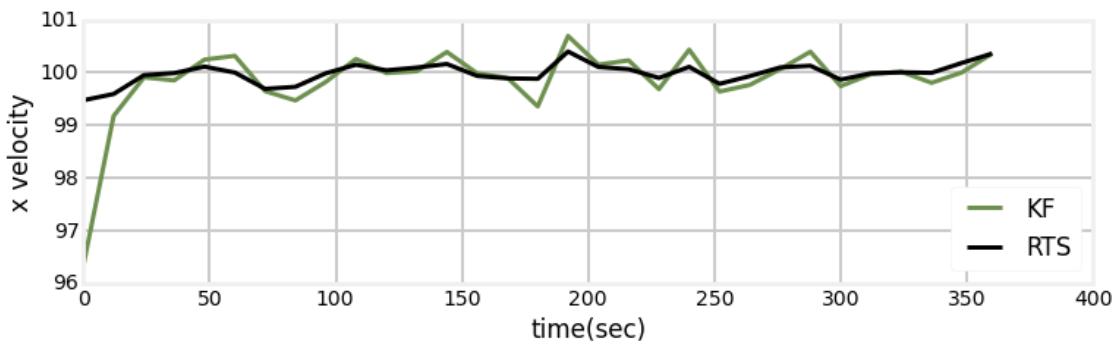
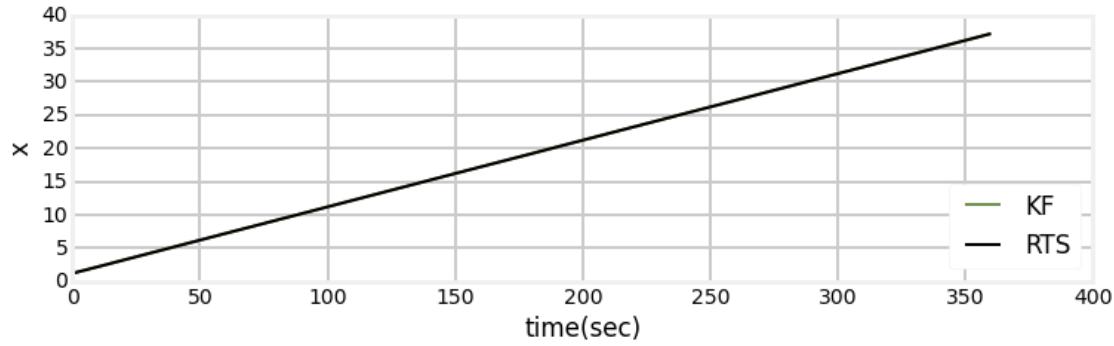
We will not develop the math or algorithm here, I will just show you how to call the algorithm in `FilterPy`. The algorithm that we have implemented is called an **RTS smoother**, after the three inventors of the algorithm: Rauch, Tung, and Striebel.

The routine is `UnscentedKalmanFilter.rts_smoother()`. Using it is trivial; we pass in the means and covariances computed from the `batch_filter` step, and receive back the smoothed means, covariances, and Kalman gain.

```
In [35]: Ms, P, K = kf.rts_smoother(xs, covs)
         ukf_internal.plot_rts_output(xs, Ms, t)
```

Difference in position in meters:

```
[-0.7717 -0.1195  0.0426 -1.1727 -1.0871]
```



From these charts we can see that the improvement in the position is small, but the improvement in the velocity is good, and spectacular for the altitude. The difference in the position are very small, so I printed the difference between the UKF and the smoothed results for the last 5 points. I recommend always usng the RTS smoother if you can post-process your data.

## 10.13 Choosing the Sigma Parameters

I have found the literature on choosing values for  $\alpha$ ,  $\beta$ , and  $\kappa$  to be rather lacking. Van der Merwe's dissertation contains the most information, but it is not exhaustive. So let's explore what they do.

Van der Merwe suggests using  $\beta = 2$  for Gaussian problems, and  $\kappa = 3 - n$ . So let's start there and vary  $\alpha$ . I will let  $n = 1$  minimize the size of the arrays we need to look at and to avoid having to compute the square root of matrices

```
In [36]: from ukf_internal import print_sigmas
print_sigmas(mean=0, cov=3, alpha=1)

sigmas: [ 0.  3. -3.]
mean weights: [ 0.6667  0.1667  0.1667]
cov weights: [ 2.6667  0.1667  0.1667]
lambda: 2
sum cov 3.0
```

So what is going on here? We can see that for a mean of 0 the algorithm choose sigma points of 0, 3, and -3, but why? Recall the equation for computing the sigma points:

$$\begin{aligned}\mathcal{X}_0 &= \mu \\ \mathcal{X}_i &= \mu \pm \sqrt{(n + \lambda)\Sigma}\end{aligned}$$

Here you will appreciate my choice of  $n = 1$  as it reduces everything to scalars, allowing us to avoid computing the square root of matrices. So, for our values the equation is

$$\begin{aligned}\mathcal{X}_0 &= 0 \\ \mathcal{X}_i &= 0 \pm \sqrt{(1 + 2) \times 3} \\ &= 0 \pm \sqrt{9} \\ &= \pm 3\end{aligned}$$

So as  $\alpha$  gets larger the sigma points get more spread out. Let's set it to an absurd value.

```
In [37]: print_sigmas(mean=0, cov=3, alpha=200)

sigmas: [ 0. 600. -600.]
mean weights: [ 1. 0. 0.]
cov weights: [-39996. 0. 0.]
lambda: 119999
sum cov -39996.0
```

We can see that the sigma point spread over 100 standard deviations. If our data was Gaussian we'd be incorporating data many standard deviations away from the mean; for nonlinear problems this is unlikely to produce good results. But suppose our distribution was not Gaussian, but instead had very fat tails? We might need to sample from those tails to get a good estimate, and hence it would make sense to make  $\kappa$  larger (not 200, which was absurdly large to make the change in the sigma points stark).

With a similar line of reasoning, suppose that our distribution has nearly no tails - the probability distribution looks more like an inverted parabola. In such a case we'd probably want to pull the sigma points in closer to the mean to avoid sampling in regions where there will never be real data.

Now let's look at the change in the weights. When we have  $k + n = 3$  the weights were 0.6667 for the mean, and 0.1667 for the two outlying sigma points. On the other hand, when  $\alpha = 200$  the mean weight shot up to 0.99999 and the outlier weights were set to 0.000004. Recall the equations for the weights:

$$W_0 = \frac{\lambda}{n + \lambda}$$

$$W_i = \frac{1}{2(n + \lambda)}$$

We can see that as  $\lambda$  gets larger the fraction for the weight of the mean ( $\lambda/(n + \lambda)$ ) approaches 1, and the fraction for the weights of the rest of the sigma points approaches 0. This is invariant on the size of your covariance. So as we sample further and further away from the mean we end up giving less weight to those samples, and if we sampled very close to the mean we'd give very similar weights to all.

However, the advice that Van der Merwe gives is to constrain  $\alpha$  in the range  $0 > \alpha \geq 1$ . He suggests  $10^{-3}$  as a good value. Lets try that.

```
In [38]: print_sigmas(mean=0, cov=13, alpha=.1, kappa=0)
```

```
sigmas: [ 0.      0.3606 -0.3606]
mean weights: [-99.   50.   50.]
cov weights: [-96.01  50.     50.  ]
lambda: -0.99
sum cov 3.99
```

I must admit to not fully understanding this advice.

```
In [39]: print_sigmas(mean=0, cov=3, alpha=.1)
```

```
sigmas: [ 0.    0.3 -0.3]
mean weights: [-32.3333  16.6667  16.6667]
cov weights: [-29.3433  16.6667  16.6667]
lambda: -0.97
sum cov 3.99
```

```
In [40]: print_sigmas(mean=0, cov=3, alpha=1)
```

```
sigmas: [ 0.   3.  -3.]
mean weights: [ 0.6667  0.1667  0.1667]
cov weights: [ 2.6667  0.1667  0.1667]
lambda: 2
sum cov 3.0
```

```
In [41]: print_sigmas(mean=0, cov=3, alpha=.1, beta=2, kappa=1)
```

```
sigmas: [ 0.      0.2449 -0.2449]
mean weights: [-49.   25.   25.]
cov weights: [-46.01  25.     25.  ]
lambda: -0.98
sum cov 3.99
```

## 10.14 Robot Localization - A Fully Worked Example

It is time to try a significant problem. Most books choose simple, textbook problems with simple answers, and you are left wondering how to implement a real world problem. This example will not teach you how to tackle any problem, but illustrates the type of things you will have to consider as you design and implement a filter.

We will consider the problem of robot localization. In this scenario we have a robot that is moving through a landscape with sensors that give range and bearings to various landmarks. This could be a self driving car using computer vision to identify trees, buildings, and other landmarks. It might be one of those small robots that vacuum your house, or a robot in a warehouse.

Our robot has 4 wheels configured the same as an automobile. It maneuvers by pivoting the front wheels. This causes the robot to pivot around the rear axle while moving forward. This is nonlinear behavior which we will have to model.

The robot has a sensor that gives it approximate range and bearing to known targets in the landscape. This is nonlinear because computing a position from a range and bearing requires square roots and trigonometry.

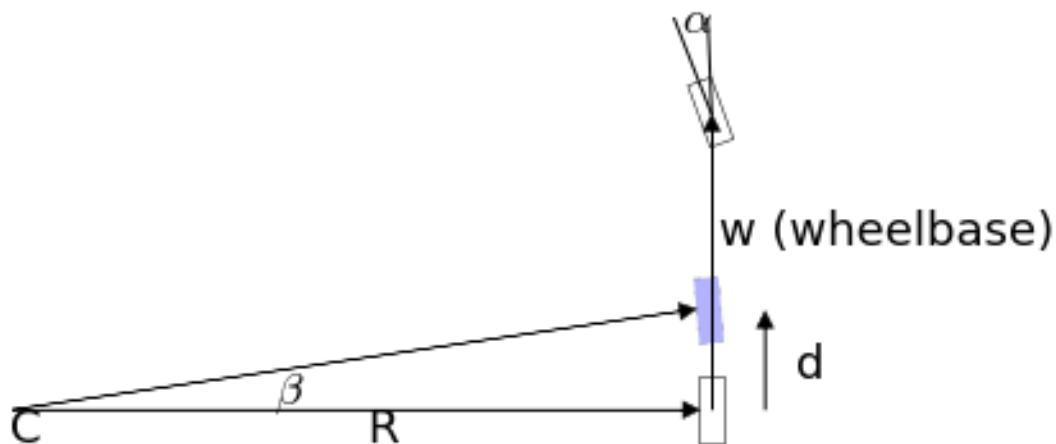
Both the process model and measurement models are nonlinear. The UKF accommodates both, so we provisionally conclude that the UKF is a viable choice for this problem.

### 10.14.1 Robot Motion Model

At a first approximation an automobile steers by pivoting the front tires while moving forward. The front of the car moves in the direction that the wheels are pointing while pivoting around the rear tires. This simple description is complicated by issues such as slippage due to friction, the differing behavior of the rubber tires at different speeds, and the need for the outside tire to travel a different radius than the inner tire. Accurately modeling steering requires a complicated set of differential equations.

For Kalman filtering, especially for lower speed robotic applications a simpler bicycle model has been found to perform well. This is a depiction of the model:

In [42]: `ekf_internal.plot_bicycle()`



Here we see the front tire is pointing in direction  $\alpha$  relative to the wheelbase. Over a short time period the car moves forward and the rear wheel ends up further ahead and slightly turned inward, as depicted with the blue shaded tire. Over such a short time frame we can approximate this as a turn around a radius  $R$ . We can compute the turn angle  $\beta$  with

$$\beta = \frac{d}{w} \tan(\alpha)$$

and the turning radius  $R$  is given by

$$R = \frac{d}{\beta}$$

where the distance the rear wheel travels given a forward velocity  $v$  is  $d = v\Delta t$ .

If we let  $\theta$  be our current orientation then we can compute the position  $C$  before the turn starts as

$$C_x = x - R \sin(\theta) C_y = y + R \cos(\theta)$$

After the move forward for time  $\Delta t$  the new position and orientation of the robot is

$$\begin{aligned} x &= C_x + R \sin(\theta + \beta) \\ y &= C_y - R \cos(\theta + \beta) \\ \theta &= \theta + \beta \end{aligned}$$

Once we substitute in for  $C$  we get

$$\begin{aligned} x &= x - R \sin(\theta) + R \sin(\theta + \beta) \\ y &= y + R \cos(\theta) - R \cos(\theta + \beta) \\ \theta &= \theta + \beta \end{aligned}$$

You do not need to understand this math in detail if you are not interested in steering models. The important thing to recognize is that our motion model is nonlinear, and we will need to deal with that with our Kalman filter.

### 10.14.2 Design the State Variables

For our robot we will maintain the position and orientation of the robot:

$$\mathbf{x} = [x \quad y \quad \theta]^T$$

I could include velocities into this model, but as you will see the math will already be quite challenging.

Our control input  $\mathbf{u}$  is the commanded velocity and steering angle

$$\mathbf{u} = [v \quad \alpha]^T$$

### 10.14.3 Design the System Model

In general we model our system as a nonlinear motion model plus noise.

$$\bar{x} = x + f(x, u) + \mathcal{N}(0, Q)$$

Using the motion model for a robot that we created above, we can write:

```
In [43]: from math import tan, sin, cos, sqrt

def move(x, u, dt, wheelbase):
    hdg = x[2]
    vel = u[0]
    steering_angle = u[1]
    dist = vel * dt

    if abs(steering_angle) > 0.001: # is robot turning?
        beta = (dist / wheelbase) * tan(steering_angle)
        r = wheelbase / tan(steering_angle) # radius

        sinh, sinhb = sin(hdg), sin(hdg + beta)
        cosh, coshb = cos(hdg), cos(hdg + beta)
        return x + np.array([-r*sinh + r*sinhb,
                             r*cosh - r*coshb, beta])
    else: # moving in straight line
        return x + np.array([dist*cos(hdg), dist*sin(hdg), 0])
```

We use this function to implement the state transition model function  $f(x)$ .

```
In [44]: def fx(x, dt, u):
          return move(x, u, dt, wheelbase)
```

I will design the UKF so that  $\Delta t$  is small. If the robot is moving slowly then this function should give a reasonably accurate prediction. In general, if  $\Delta t$  is large or your system's dynamics are very nonlinear this method will fail. In those cases you will need to implement this function using a more sophisticated numerical integration technique such as Runge Kutta. Numerical integration is covered briefly in the **Kalman Filter Math** chapter.

#### 10.14.4 Design the Measurement Model

Now we need to design our measurement model. For this problem we are assuming that we have a sensor that receives a noisy bearing and range to multiple known locations in the landscape. The measurement model must convert the state  $[x \ y \ \theta]^T$  into a range and bearing to the landmark. If  $p$  is the position of a landmark, the range  $r$  is

$$r = \sqrt{(p_x - x)^2 + (p_y - y)^2}$$

We assume that the sensor provides bearing relative to the orientation of the robot, so we must subtract the robot's orientation from the bearing to get the sensor reading, like so:

$$\phi = \tan^{-1}\left(\frac{p_y - y}{p_x - x}\right) - \theta$$

Thus our measurement function is

$$\begin{aligned} \mathbf{z} &= h(\mathbf{x}, \mathbf{P}) && + \mathcal{N}(0, R) \\ &= \begin{bmatrix} \sqrt{(p_x - x)^2 + (p_y - y)^2} \\ \tan^{-1}\left(\frac{p_y - y}{p_x - x}\right) - \theta \end{bmatrix} && + \mathcal{N}(0, R) \end{aligned}$$

I will not implement this in Python yet as there is a difficulty that will be discussed in the [Implementation](#) section below.

#### 10.14.5 Design Measurement Noise

This is quite straightforward as we need to specify measurement noise in measurement space, hence it is linear. It is reasonable to assume that the range and bearing measurement noise is independent, hence

$$\mathbf{R} = \begin{bmatrix} \sigma_{range}^2 & 0 \\ 0 & \sigma_{bearing}^2 \end{bmatrix}$$

#### 10.14.6 Implementation

Before we begin coding the main loop we have another issue to handle. The residual is defined as  $y = z - h(x)$ . Suppose  $z$  has a bearing of  $1^\circ$  and  $h(x)$  has a bearing of  $359^\circ$ . Subtracting them gives  $-358^\circ$ . this will throw off the computation of the Kalman gain because the correct angle difference in this case is  $2^\circ$ . So we will have to write code to correctly compute the bearing residual.

```
In [45]: def normalize_angle(x):
    x = x % (2 * np.pi)      # force in range [0, 2 pi)
    if x > np.pi:            # move to [-pi, pi)
        x -= 2 * np.pi
    return x
```

```
In [46]: print(np.rad2deg(normalize_angle(np.deg2rad(1-359))))
```

2.0

The state vector has the bearing at index 2, but the measurement vector has it at index 1, so we need to write functions to handle each. Another issue we face is that as the robot maneuvers different landmarks will be visible, so we need to handle a variable number of measurements. The function for the residual in the measurement will be passed an array of several measurements, one per landmark.

```
In [47]: def residual_h(a, b):
    y = a - b
    # data in format [dist_1, bearing_1, dist_2, bearing_2, ...]
    for i in range(0, len(y), 2):
        y[i + 1] = normalize_angle(y[i + 1])
    return y

def residual_x(a, b):
    y = a - b
    y[2] = normalize_angle(y[2])
    return y
```

We can now implement the measurement model. The equation is

$$h(\mathbf{x}, \mathbf{P}) = \begin{bmatrix} \sqrt{(p_x - x)^2 + (p_y - y)^2} \\ \tan^{-1}\left(\frac{p_y - y}{p_x - x}\right) - \theta \end{bmatrix}$$

The expression  $\tan^{-1}\left(\frac{p_y - y}{p_x - x}\right) - \theta$  can produce a result outside the range  $[-\pi, \pi]$ , so we should normalize the angle to that range.

The function will be passed an array of landmarks and needs to produce an array of measurements in the form [dist\_to\_1, bearing\_to\_1, dist\_to\_2, bearing\_to\_2, ...].

```
In [48]: def Hx(x, landmarks):
    """ takes a state variable and returns the measurement
    that would correspond to that state. """
    hx = []
    for lmark in landmarks:
        px, py = lmark
        dist = sqrt((px - x[0])**2 + (py - x[1])**2)
        angle = atan2(py - x[1], px - x[0])
        hx.extend([dist, normalize_angle(angle - x[2])])
    return np.array(hx)
```

Our difficulties are not over. The unscented transform computes the average of the state and measurement vectors, but each contains a bearing. There is no unique way to compute the average of a set of angles. For example, what is the average of  $359^\circ$  and  $3^\circ$ ? Intuition suggests the answer should be  $1^\circ$ , but a naive sum/count approach yields  $181^\circ$ .

One common approach is to take the arctan of the sum of the sins and cosines.

$$\bar{\theta} = \text{atan2}\left(\frac{\sum_{i=1}^n \sin \theta_i}{n}, \frac{\sum_{i=1}^n \cos \theta_i}{n}\right)$$

We have not used this feature yet, but the `UnscentedKalmanFilter.__init__()` method has an argument `x_mean_fn` for a function which computes the mean of the state, and `z_mean_fn` for a function which computes the mean of the measurement. We will code these function as:

```
In [49]: def state_mean(sigmas, Wm):
    x = np.zeros(3)

    sum_sin = np.sum(np.dot(np.sin(sigmas[:, 2]), Wm))
    sum_cos = np.sum(np.dot(np.cos(sigmas[:, 2]), Wm))
    x[0] = np.sum(np.dot(sigmas[:, 0], Wm))
    x[1] = np.sum(np.dot(sigmas[:, 1], Wm))
    x[2] = atan2(sum_sin, sum_cos)
    return x

def z_mean(sigmas, Wm):
    z_count = sigmas.shape[1]
    x = np.zeros(z_count)

    for z in range(0, z_count, 2):
        sum_sin = np.sum(np.dot(np.sin(sigmas[:, z+1]), Wm))
        sum_cos = np.sum(np.dot(np.cos(sigmas[:, z+1]), Wm))

        x[z] = np.sum(np.dot(sigmas[:, z], Wm))
        x[z+1] = atan2(sum_sin, sum_cos)
    return x
```

These functions take advantage of the fact that NumPy's trigometric functions operate on arrays, and `dot` performs element-wise multiplication. NumPy is implemented in C and Fortran, so `sum(dot(sin(x), w))` is much faster than writing the equivalent loop in Python.

With that done we are now ready to implement the UKF. I want to point out that when I designed this filter I did not just design all of functions above in one sitting, from scratch. I put together a basic UKF with predefined landmarks, verified it worked, then started filling in the pieces. “What if I see different landmarks?” That lead me to change the measurement function to accept an array of landmarks. “How do I deal with computing the residual of angles?” This led me to write the angle normalization code. “What is the mean of a set of angles?” I searched on the internet, found an article on Wikipedia, and implemented that algorithm. Do not be daunted. Design what you can, then ask questions and solve them, one by one.

You’ve seen the UKF implementation already, so I will not describe it in detail. There is one new thing here that is important to discuss. When we construct the sigma points and filter we have to provide it the functions that we have written to compute the residuals and means.

```
points = SigmaPoints(n=3, alpha=.00001, beta=2, kappa=0,
                     subtract=residual_x)

ukf = UKF(dim_x=3, dim_z=2, fx=fx, hx=Hx, dt=dt, points=points,
           x_mean_fn=state_mean, z_mean_fn=z_mean,
           residual_x=residual_x, residual_z=residual_h)
```

The rest of the code runs the simulation and plots the results. I create a variable `landmarks` that contains the coordinates of the landmarks. I update the simulated robot position 10 times a second, but run the UKF only once per second. This is for two reasons. First, we are not using Runge Kutta to integrate the differential equations of motion, so a narrow time step allows our simulation to be more accurate. Second, it is fairly normal in embedded systems to have limited processing speed. This forces you to run your Kalman filter only as frequently as absolutely needed.

```
In [50]: from filterpy.stats import plot_covariance_ellipse

dt = 1.0
wheelbase = 0.5

def run_localization(
    cmds, landmarks, sigma_vel, sigma_steer, sigma_range,
    sigma_bearing, ellipse_step=1, step=10):

    points = SigmaPoints(n=3, alpha=.00001, beta=2, kappa=0,
                         subtract=residual_x)
    ukf = UKF(dim_x=3, dim_z=2*len(landmarks), fx=fx, hx=Hx,
               dt=dt, points=points, x_mean_fn=state_mean,
               z_mean_fn=z_mean, residual_x=residual_x,
               residual_z=residual_h)

    ukf.x = np.array([2, 6, .3])
    ukf.P = np.diag([.1, .1, .05])
    ukf.R = np.diag([sigma_range**2,
                    sigma_bearing**2]*len(landmarks))
    ukf.Q = np.eye(3)*0.0001

    sim_pos = ukf.x.copy()

    # plot landmarks
    if len(landmarks) > 0:
        plt.scatter(landmarks[:, 0], landmarks[:, 1],
                    marker='s', s=60)
```

```

track = []
for i, u in enumerate(cmds):
    sim_pos = move(sim_pos, u, dt/step, wheelbase)
    track.append(sim_pos)

    if i % step == 0:
        ukf.predict(fx_args=u)

    if i % ellipse_step == 0:
        plot_covariance_ellipse(
            (ukf.x[0], ukf.x[1]), ukf.P[0:2, 0:2], std=6,
            facecolor='k', alpha=0.3)

x, y = sim_pos[0], sim_pos[1]
z = []
for lmark in landmarks:
    dx, dy = lmark[0] - x, lmark[1] - y
    d = sqrt(dx**2 + dy**2) + randn()*sigma_range
    bearing = atan2(lmark[1] - y, lmark[0] - x)
    a = (normalize_angle(bearing - sim_pos[2] +
                          randn()*sigma_bearing))
    z.extend([d, a])
ukf.update(z, hx_args=(landmarks,))

if i % ellipse_step == 0:
    plot_covariance_ellipse(
        (ukf.x[0], ukf.x[1]), ukf.P[0:2, 0:2], std=6,
        facecolor='g', alpha=0.8)

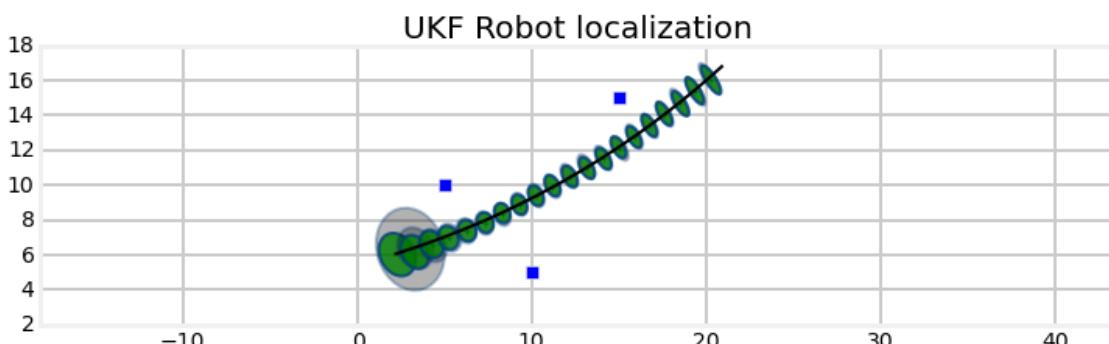
track = np.array(track)
plt.plot(track[:, 0], track[:, 1], color='k', lw=2)
plt.axis('equal')
plt.title("UKF Robot localization")
plt.show()
return ukf

```

```

In [51]: landmarks = np.array([[5, 10], [10, 5], [15, 15]])
cmds = [np.array([1.1, .01])] * 200
ukf = run_localization(
    cmd, landmarks, sigma_vel=0.1, sigma_steer=np.radians(1),
    sigma_range=0.3, sigma_bearing=0.1)
print('Final P:', ukf.P.diagonal())

```



```
Final P: [ 0.0092  0.0187  0.0007]
```

The rest of the code runs the simulation and plots the results. I create a variable `landmarks` that contains the coordinates of the landmarks. I update the simulated robot position 10 times a second, but run the UKF only once. This is for two reasons. First, we are not using Runge Kutta to integrate the differential equations of motion, so a narrow time step allows our simulation to be more accurate. Second, it is fairly normal in embedded systems to have limited processing speed. This forces you to run your Kalman filter only as frequently as absolutely needed.

#### 10.14.7 Steering the Robot

The steering simulation in the run above is not realistic. The velocity and steering angles never changed, which doesn't pose much of a problem for the Kalman filter. We could implement a complicated PID controlled robot simulation, but I will just generate varying steering commands using NumPy's `linspace` method. I'll also add more landmarks as the robot will be traveling much further than in the first example.

```
In [52]: landmarks = np.array([[5, 10], [10, 5], [15, 15], [20, 5],
                           [0, 30], [50, 30], [40, 10]])
dt = 0.1
wheelbase = 0.5
sigma_range=0.3
sigma_bearing=0.1

def turn(v, t0, t1, steps):
    return [[v, a] for a in np.linspace(
        np.radians(t0), np.radians(t1), steps)]]

# accelerate from a stop
cmds = [[v, .0] for v in np.linspace(0.001, 1.1, 30)]
cmds.extend([cmds[-1]]*50)

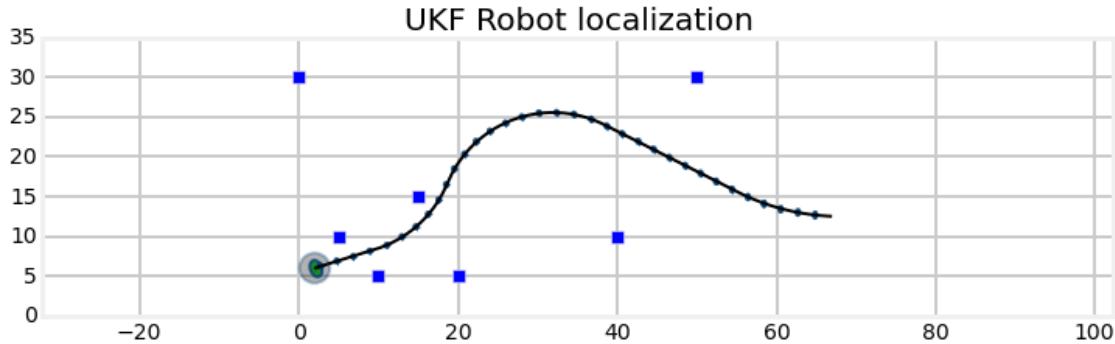
# turn left
v = cmds[-1][0]
cmds.extend(turn(v, 0, 2, 15))
cmds.extend([cmds[-1]]*100)

#turn right
cmds.extend(turn(v, 2, -2, 15))
cmds.extend([cmds[-1]]*200)

cmds.extend(turn(v, -2, 0, 15))
cmds.extend([cmds[-1]]*150)

cmds.extend(turn(v, 0, 1, 25))
cmds.extend([cmds[-1]]*100)

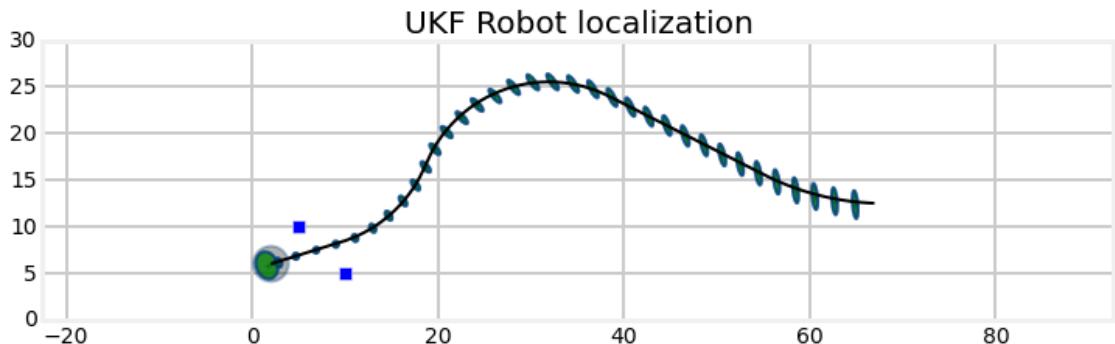
In [53]: ukf = run_localization(
            cmds, landmarks, sigma_vel=0.1, sigma_steer=np.radians(1),
            sigma_range=0.3, sigma_bearing=0.1, step=1,
            ellipse_step=20)
print('final covariance', ukf.P.diagonal())
```



```
final covariance [ 0.0013  0.0043  0.0004]
```

You can see that the uncertainty becomes very small very quickly. The covariance ellipses are displaying the  $6\sigma$  covariance, yet the ellipses are so small they are hard to see. We can incorporate error into the answer by only supplying two landmarks near the start point. When we run this filter the errors increase as the robot gets further away.

```
In [54]: ukf = run_localization(
    cmd, landmarks[0:2], sigma_vel=0.1, sigma_steer=np.radians(1),
    sigma_range=0.3, sigma_bearing=0.1, step=1,
    ellipse_step=20)
print('final covariance', ukf.P.diagonal())
```



```
final covariance [ 0.0026  0.0657  0.0008]
```

## 10.15 Discussion

Your impression of this chapter probably depends on how many nonlinear Kalman filters you have implemented in the past. If this is your first exposure perhaps the computation of  $2n + 1$  sigma points and the subsequent writing of the  $f(x)$  and  $h(x)$  function struck you as a bit finicky. Indeed, I spent more time than

I'd care to admit getting everything working because of the need to handle the modular math of angles. On the other hand, if you have implemented an extended Kalman filter (EKF) you are perhaps bouncing gleefully in your seat. There is a small amount of tedium in writing the functions for the UKF, but the concepts are very basic. The EKF for the same problems requires some fairly difficult mathematics. In fact, for many problems we cannot find a closed form solution for the equations of the EKF, and we must retreat to some sort of iterated solution.

However, the advantage of the UKF over the EKF is not only the relative ease of implementation. It is somewhat premature to discuss this because you haven't learned the EKF yet, but the EKF linearizes the problem at one point and passes that point through a linear Kalman filter. In contrast, the UKF takes  $2n+1$  samples. Therefore the UKF is often more accurate than the EKF, especially when the problem is highly nonlinear. While it is not true that the UKF is guaranteed to always outperform the EKF, in practice it has been shown to perform at least as well, and usually much better than the EKF.

Hence my recommendation is to always start by implementing the UKF. If your filter has real world consequences if it diverges (people die, lots of money lost, power plant blows up) of course you will have to engage in a lot of sophisticated analysis and experimentation to choose the best filter. That is beyond the scope of this book, and you should be going to graduate school to learn this theory.

I have spoken of the UKF as the way to perform sigma point filters. This is not true. The specific version I chose is Julier's scaled unscented filter as parameterized by Van der Merwe in his 2004 dissertation. If you search for Julier, Van der Merwe, Uhlmann, and Wan you will find a family of similar sigma point filters that they developed. Each technique uses a different way of choosing and weighting the sigma points. But the choices don't stop there. For example, the SVD Kalman filter uses singular value decomposition (SVD) to find the approximate mean and covariance of the probability distribution. Think of this chapter as an introduction to the sigma point filters, rather than a definitive treatment of how they work.

## 10.16 References

- [1] Rudolph Van der Merwe. "Sigma-Point Kalman Filters for Probabilistic Inference in Dynamic State-Space Models" dissertation (2004).
- [2] Simon J. Julier. "The Scaled Unscented Transformation". Proceedings of the American Control Conference 6. IEEE. (2002)
- [1] <http://www.esdradar.com/brochures/Compact%20Tracking%2037250X.pdf>
- [2] Julier, Simon J.; Uhlmann, Jeffrey "A New Extension of the Kalman Filter to Nonlinear Systems". Proc. SPIE 3068, Signal Processing, Sensor Fusion, and Target Recognition VI, 182 (July 28, 1997)
- [3] Cholesky decomposition. Wikipedia. [http://en.wikipedia.org/wiki/Cholesky\\_decomposition](http://en.wikipedia.org/wiki/Cholesky_decomposition)

# Chapter 11

## The Extended Kalman Filter

At this point in the book we have developed the theory for the linear Kalman filter. Then, in the last two chapters we broached the topic of using Kalman filters for nonlinear problems. In this chapter we will learn the Extended Kalman filter (EKF). The EKF handles nonlinearity by linearizing the system at the point of the current estimate, and then the linear Kalman filter is used to filter this linearized system. It was one of the very first techniques used for nonlinear problems, and it remains the most common technique.

The EKF provides significant mathematical challenges to the designer of the filter; this is the most challenging chapter of the book. To be honest, I do everything I can to avoid the EKF in favor of other techniques that have been developed to filter nonlinear problems. However, the topic is unavoidable; all classic papers and a majority of current papers in the field use the EKF. Even if you do not use the EKF in your own work you will need to be familiar with the topic to be able to read the literature.

### 11.1 Linearizing the Kalman Filter

The Kalman filter uses linear equations, so it does not work with nonlinear problems. Problems can be nonlinear in two ways. First, the process model might be nonlinear. An object falling through the atmosphere encounters drag which reduces its acceleration. The amount of drag varies based on the velocity the object. The resulting behavior is nonlinear - it cannot be modeled with linear equations. Second, the measurements could be nonlinear. For example, a radar gives a range and bearing to a target. We use trigonometry, which is nonlinear, to compute the position of the target.

For the linear filter we have these equations for the process and measurement models:

$$\begin{aligned}\bar{\mathbf{x}} &= \mathbf{Ax} + \mathbf{Bu} + w_x \\ \mathbf{z} &= \mathbf{Hx} + w_z\end{aligned}$$

For the nonlinear model these equations must be modified to read:

$$\begin{aligned}\bar{\mathbf{x}} &= f(\mathbf{x}, \mathbf{u}) + w_x \\ \mathbf{z} &= h(\mathbf{x}) + w_z\end{aligned}$$

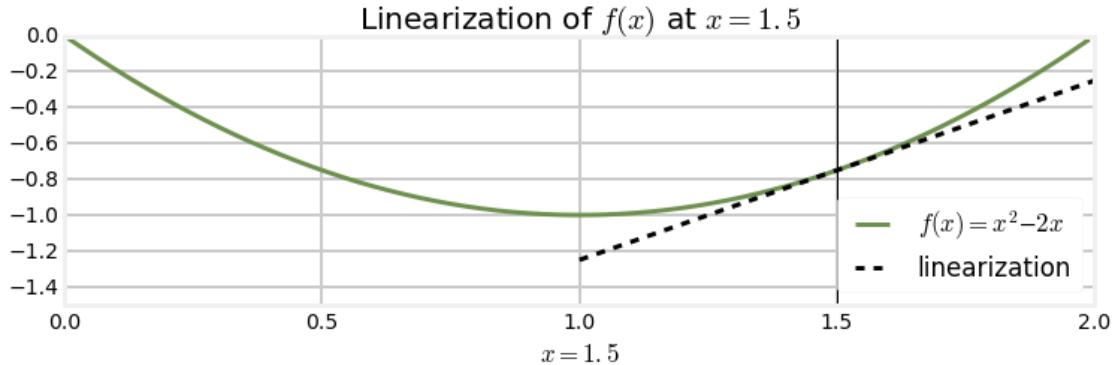
The linear expression  $\mathbf{Ax} + \mathbf{Bu}$  is replaced by a nonlinear function  $f(\mathbf{x}, \mathbf{u})$ , and the linear expression  $\mathbf{Hx}$  is replaced by a nonlinear function  $h(\mathbf{x})$ .

You might imagine that we proceed by finding a new set of Kalman filter equations that optimally solve these equations. But if you remember the charts in the **Nonlinear Filtering** chapter you'll recall that passing a Gaussian through a nonlinear function results in a probability distribution that is no longer Gaussian. So this will not work.

The EKF does not alter the Kalman filter's linear equations. Instead, it linearizes the nonlinear equations at the point of the current estimate, and uses this linearization in the linear Kalman filter.

Linearize means what it sounds like. We find a line that most closely matches the curve at a defined point. The graph below linearizes the parabola  $f(x) = x^2 - 2x$  at  $x = 1.5$ .

```
In [2]: import ekf_internal
ekf_internal.show_linearization()
```



If the curve above is the process model, then the dotted lines shows the linearization of that curve for the estimate  $x = 1.5$ .

We linearize systems by finding the slope of the curve at the given point:

$$\begin{aligned} f(x) &= x^2 - 2x \\ \frac{df}{dx} &= 2x - 2 \end{aligned}$$

and then finding its value at the evaluation point:

$$\begin{aligned} m &= f'(x = 1.5) \\ &= 2(1.5) - 2 \\ &= 1 \end{aligned}$$

Our math will be more complicated because we are working with systems of differential equations. We linearize  $f(\mathbf{x}, \mathbf{u})$ , and  $h(\mathbf{x})$  by taking the partial derivatives ( $\frac{\partial}{\partial \mathbf{x}}$ ) of each to evaluate  $\mathbf{A}$  and  $\mathbf{H}$  at the point  $\mathbf{x}_t$  and  $\mathbf{u}_t$ . This gives us the system dynamics matrix and measurement model matrix:

$$\begin{aligned} \mathbf{A} &= \left. \frac{\partial f(\mathbf{x}_t, \mathbf{u}_t)}{\partial \mathbf{x}} \right|_{\mathbf{x}_t, \mathbf{u}_t} \\ \mathbf{H} &= \left. \frac{\partial h(\mathbf{x}_t)}{\partial \mathbf{x}} \right|_{\mathbf{x}_t} \end{aligned}$$

Finally, we find the discrete state transition matrix  $\mathbf{F}$  by using the approximation of the Taylor-series expansion of  $e^{\mathbf{A}\Delta t}$ :

$$\mathbf{F} = e^{\mathbf{A}\Delta t} = \mathbf{I} + \mathbf{A}\Delta t + \frac{(\mathbf{A}\Delta t)^2}{2!} + \frac{(\mathbf{A}\Delta t)^3}{3!} + \dots$$

Alternatively, you can use one of the other techniques we learned in the **Kalman Math** chapter.

This leads to the following equations for the EKF. I placed them beside the equations for the linear Kalman filter, and put boxes around the only changes:

linear Kalman filter	EKF
	$\mathbf{A} = \left. \frac{\partial f(\mathbf{x}_t, \mathbf{u}_t)}{\partial \mathbf{x}} \right _{\mathbf{x}_t, \mathbf{u}_t}$
$\bar{\mathbf{x}} = \mathbf{F}\mathbf{x} + \mathbf{B}\mathbf{u}$	$\mathbf{F} = e^{\mathbf{A}\Delta t}$
$\bar{\mathbf{P}} = \mathbf{F}\mathbf{P}\mathbf{F}^\top + \mathbf{Q}$	$\bar{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$
	$\bar{\mathbf{P}} = \mathbf{F}\mathbf{P}\mathbf{F}^\top + \mathbf{Q}$
$\mathbf{y} = \mathbf{z} - \mathbf{H}\bar{\mathbf{x}}$	$\mathbf{H} = \left. \frac{\partial h(\mathbf{x}_t)}{\partial \mathbf{x}} \right _{\mathbf{x}_t}$
$\mathbf{K} = \bar{\mathbf{P}}\mathbf{H}^\top (\mathbf{H}\bar{\mathbf{P}}\mathbf{H}^\top + \mathbf{R})^{-1}$	$\mathbf{y} = \mathbf{z} - \mathbf{H}\bar{\mathbf{x}}$
$\mathbf{x} = \bar{\mathbf{x}} + \mathbf{K}\mathbf{y}$	$\mathbf{K} = \bar{\mathbf{P}}\mathbf{H}^\top (\mathbf{H}\bar{\mathbf{P}}\mathbf{H}^\top + \mathbf{R})^{-1}$
$\mathbf{P} = (\mathbf{I} - \mathbf{K}\mathbf{H})\bar{\mathbf{P}}$	$\mathbf{x} = \bar{\mathbf{x}} + \mathbf{K}\mathbf{y}$
	$\mathbf{P} = (\mathbf{I} - \mathbf{K}\mathbf{H})\bar{\mathbf{P}}$

We don't normally use  $\mathbf{F}\mathbf{x}$  to propagate the state for the EKF as the linearization causes inaccuracies. It is typical to compute  $\bar{\mathbf{x}}$  using a numerical integration technique such as Euler or Runge Kutta. Thus I wrote  $\bar{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$ .

I think the easiest way to understand the EKF is to start off with an example. After we do a few examples you may want to come back and reread this section.

## 11.2 Example: Tracking a Flying Airplane

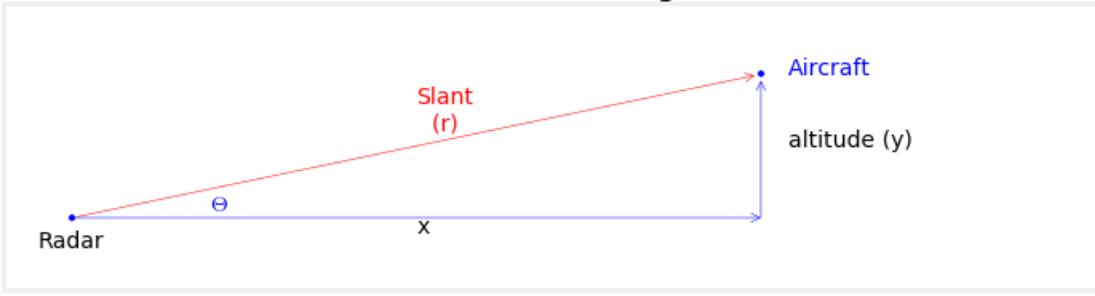
We will start by simulating tracking an airplane by using ground based radar. We implemented a UKF for this problem in the last chapter. Now we will implement an EKF for the same problem so we can compare both the filter performance and the level of effort required to implement the filter.

Radar work by emitting a beam of radio waves and scanning for a return bounce. Anything in the beam's path will reflects some of the signal back to the radar. By timing how long it takes for the reflected signal to get back to the radar the system can compute the slant distance - the straight line distance from the radar installation to the object.

For this example we want to take the slant range measurement from the radar and compute the horizontal position (distance of aircraft from the radar measured over the ground) and altitude of the aircraft, as in the diagram below.

```
In [3]: import ekf_internal
ekf_internal.show_radar_chart()
```

## Radar Tracking



This gives us the equality  $x = \sqrt{slant^2 - altitude^2}$ .

### 11.2.1 Design the State Variables

We want to track the position of an aircraft assuming a constant velocity and altitude, and measurements of the slant distance to the aircraft. That means we need 3 state variables - horizontal distance, horizontal velocity, and altitude:

$$\mathbf{x} = \begin{bmatrix} \text{distance} \\ \text{velocity} \\ \text{altitude} \end{bmatrix} = \begin{bmatrix} x \\ \dot{x} \\ y \end{bmatrix}$$

### 11.2.2 Design the Process Model

We assume a Newtonian, kinematic system for the aircraft. We've used this model in previous chapters, so by inspection you may recognize that we want

$$\mathbf{F} = \left[ \begin{array}{cc|c} 1 & \Delta t & 0 \\ 0 & 1 & 0 \\ \hline 0 & 0 & 1 \end{array} \right]$$

I've partitioned the matrix into blocks to show the upper left block is a constant velocity model for  $x$ , and the lower right block is a constant position model for  $y$ .

However, let's practice finding these matrix for a nonlinear system. We model nonlinear systems with a set of differential equations. We need an equation in the form

$$\dot{\mathbf{x}} = \mathbf{Ax} + \mathbf{w}$$

where  $\mathbf{w}$  is the system noise.

The variables  $x$  and  $y$  are independent so we can compute them separately. The differential equations for motion in one dimension are:

$$\begin{aligned} v &= \dot{x} \\ a &= \ddot{x} = 0 \end{aligned}$$

Now we put the differential equations into state-space form. If this was a second or greater order differential system we would have to first reduce them to an equivalent set of first degree equations. The equations are first order, so we put them in state space matrix form as

$$\begin{bmatrix} \dot{x} \\ \ddot{x} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}$$

where  $\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$ .

Recall that  $\mathbf{A}$  is the system dynamics matrix. It describes a set of linear differential equations. From it we must compute the state transition matrix  $\mathbf{F}$ .  $\mathbf{F}$  describes a discrete set of linear equations which compute  $\mathbf{x}$  for a discrete time step  $\Delta t$ .

and solve the following power series expansion of the matrix exponential to linearize the equations at  $t$ :

$$\mathbf{F}(\Delta t) = e^{\mathbf{A}\Delta t} = \mathbf{I} + \mathbf{A}\Delta t + \frac{(\mathbf{A}\Delta t)^2}{2!} + \frac{(\mathbf{A}\Delta t)^3}{3!} + \dots$$

$\mathbf{A}^2 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$ , so all higher powers of  $\mathbf{A}$  are also  $\mathbf{0}$ . Thus the power series expansion is:

$$\begin{aligned} \mathbf{F}(\Delta t) &= \mathbf{I} + \mathbf{A}\Delta t + \mathbf{0} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \Delta t \\ &= \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix} \end{aligned}$$

This give us

$$\begin{aligned} \bar{\mathbf{x}} &= \mathbf{F}\mathbf{x} \\ \bar{\mathbf{x}} &= \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix} \mathbf{x} \end{aligned}$$

This is the same result used by the kinematic equations! This exercise was unnecessary other than to illustrate linearizing differential equations. Subsequent examples will require you to use these techniques.

### 11.2.3 Design the Measurement Model

The measurement function for our filter needs to take the filter state  $\mathbf{x}$  and turn it into a measurement, which is the slant range distance. We use the Pythagorean theorem to derive

$$h(\mathbf{x}) = \sqrt{x^2 + y^2}$$

The relationship between the slant distance and the position on the ground is nonlinear due to the square root term. To use it in the EKF we must linearize it. As we discussed above, the best way to linearize an equation at a point is to find its slope, which we do by evaluating its partial derivative at a point:

$$\mathbf{H} = \left. \frac{\partial h(\mathbf{x})}{\partial \mathbf{x}} \right|_{\mathbf{x}_t}$$

The partial derivative of a matrix is called a Jacobian, and takes the form

$$\frac{\partial \mathbf{H}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial h_1}{\partial x_1} & \frac{\partial h_1}{\partial x_2} & \dots \\ \frac{\partial h_2}{\partial x_1} & \frac{\partial h_2}{\partial x_2} & \dots \\ \vdots & \vdots & \end{bmatrix}$$

In other words, each element in the matrix is the partial derivative of the function  $h$  with respect to the variables  $x$ . For our problem we have

$$\mathbf{H} = [\partial h / \partial x \quad \partial h / \partial \dot{x} \quad \partial h / \partial y]$$

where  $h(x) = \sqrt{x^2 + y^2}$ .

Solving each in turn:

$$\begin{aligned}\frac{\partial h}{\partial x} &= \frac{\partial}{\partial x} \sqrt{x^2 + y^2} \\ &= \frac{x}{\sqrt{x^2 + y^2}}\end{aligned}$$

and

$$\begin{aligned}\frac{\partial h}{\partial \dot{x}} &= \frac{\partial}{\partial \dot{x}} \sqrt{x^2 + y^2} \\ &= 0\end{aligned}$$

and

$$\begin{aligned}\frac{\partial h}{\partial y} &= \frac{\partial}{\partial y} \sqrt{x^2 + y^2} \\ &= \frac{y}{\sqrt{x^2 + y^2}}\end{aligned}$$

giving us

$$\mathbf{H} = \left[ \begin{array}{ccc} \frac{x}{\sqrt{x^2+y^2}} & 0 & \frac{y}{\sqrt{x^2+y^2}} \end{array} \right]$$

This may seem daunting, so step back and recognize that all of this math is doing something very simple. We have an equation for the slant range to the airplane which is nonlinear. The Kalman filter only works with linear equations, so we need to find a linear equation that approximates  $\mathbf{H}$ . As we discussed above, finding the slope of a nonlinear equation at a given point is a good approximation. For the Kalman filter, the ‘given point’ is the state variable  $\mathbf{x}$  so we need to take the derivative of the slant range with respect to  $\mathbf{x}$ .

To make this more concrete, let’s now write a Python function that computes the Jacobian of  $\mathbf{H}$  for this problem. The `ExtendedKalmanFilter` class will be using this to generate `ExtendedKalmanFilter.H` at each step of the process.

```
In [4]: from math import sqrt
def HJacobian_at(x):
    """ compute Jacobian of H matrix at x """

    horiz_dist = x[0]
    altitude   = x[2]
    denom = sqrt(horiz_dist**2 + altitude**2)
    return array ([[horiz_dist/denom, 0., altitude/denom]])
```

Finally, let’s provide the code for  $h(\mathbf{x})$

```
In [5]: def hx(x):
    """ compute measurement for slant range that
    would correspond to state x.
    """
    return (x[0]**2 + x[2]**2) ** 0.5
```

Now lets write a simulation for our radar.

```
In [6]: from numpy.random import randn
import math

class RadarSim(object):
    """ Simulates the radar signal returns from an object
    flying at a constant altitude and velocity in 1D.
    """

    def __init__(self, dt, pos, vel, alt):
        self.pos = pos
        self.vel = vel
        self.alt = alt
        self.dt = dt

    def get_range(self):
        """ Returns slant range to the object. Call once
        for each new measurement at dt time from last call.
        """

        # add some process noise to the system
        self.vel = self.vel + .1*randn()
        self.alt = self.alt + .1*randn()
        self.pos = self.pos + self.vel*self.dt

        # add measurement noise
        err = self.pos * 0.05*randn()
        slant_dist = math.sqrt(self.pos**2 + self.alt**2)

        return slant_dist + err
```

#### 11.2.4 Design Process and Measurement Noise

The radar returns the range distance. A good radar can achieve accuracy of  $\sigma_{range} = 5$  meters, so we will use that value. This gives us

$$\mathbf{R} = [\sigma_{range}^2] = [25]$$

The design of  $\mathbf{Q}$  requires some discussion. The state  $\mathbf{x} = [x \ \dot{x} \ y]^T$ . The first two elements are position (down range distance) and velocity, so we can use `Q_discrete_white_noise` noise to compute the values for the upper left hand side of  $\mathbf{Q}$ . The third element of  $\mathbf{x}$  is altitude, which we are assuming is independent of the down range distance. That leads us to a block design of  $\mathbf{Q}$  of:

$$\mathbf{Q} = \begin{bmatrix} \mathbf{Q}_x & 0 \\ 0 & \mathbf{Q}_y \end{bmatrix}$$

### 11.2.5 Implementation

The FilterPy library provides the class `ExtendedKalmanFilter`. It works very similar to the `KalmanFilter` class we have been using, except that it allows you to provide functions that compute the Jacobian of  $\mathbf{H}$  and the function  $h(\mathbf{x})$ . We have already written the code for these two functions, so let's get going.

We start by importing the filter and creating it. There are 3 variables in  $\mathbf{x}$  and only 1 measurement. At the same time we will create our radar simulator.

```
from filterpy.kalman import ExtendedKalmanFilter

rk = ExtendedKalmanFilter(dim_x=3, dim_z=1)
radar = RadarSim(dt, pos=0., vel=100., alt=1000.)
```

We will initialize the filter near the airplane's actual position

```
rk.x = array([radar.pos, radar.vel-10, radar.alt+100])
```

We assign the system matrix using the first term of the Taylor series expansion we computed above.

```
dt = 0.05
rk.F = eye(3) + array([[0, 1, 0],
                      [0, 0, 0],
                      [0, 0, 0]]) * dt
```

After assigning reasonable values to  $\mathbf{R}$ ,  $\mathbf{Q}$ , and  $\mathbf{P}$  we can run the filter with a simple loop

```
for i in range(int(20/dt)):
    z = radar.get_range()
    rk.update(array([z]), HJacobian_at, hx)
    rk.predict()
```

Adding some boilerplate code to save and plot the results we get:

```
In [7]: from filterpy.common import Q_discrete_white_noise
        from filterpy.kalman import ExtendedKalmanFilter
        from numpy import eye, array, asarray
        import numpy as np

        dt = 0.05
        rk = ExtendedKalmanFilter(dim_x=3, dim_z=1)
        radar = RadarSim(dt, pos=0., vel=100., alt=1000.)

        # make an imperfect starting guess
        rk.x = array([radar.pos-100, radar.vel+100, radar.alt+1000])

        rk.F = eye(3) + array([[0, 1, 0],
                              [0, 0, 0],
                              [0, 0, 0]]) * dt

        range_std = 5. # meters
        rk.R = np.diag([range_std**2])
        rk.Q[0:2, 0:2] = Q_discrete_white_noise(2, dt=dt, var=0.1)
```

```

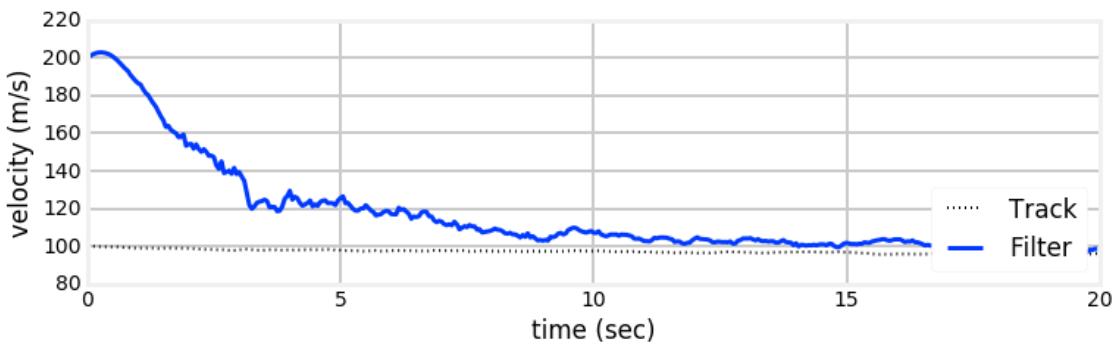
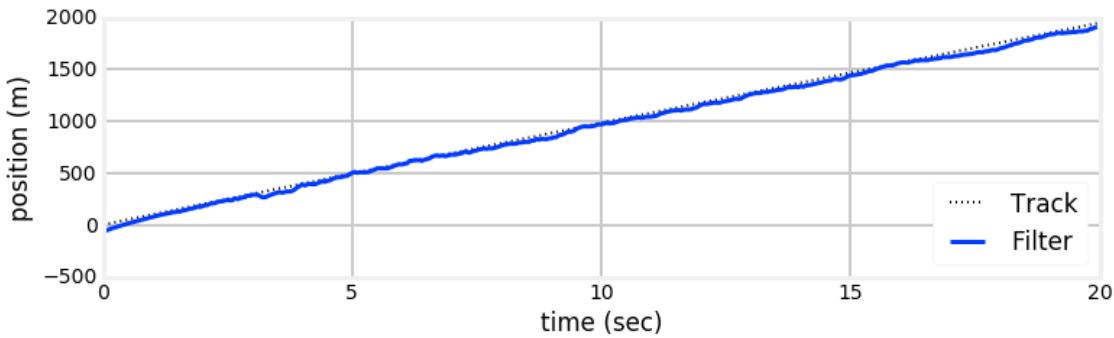
rk.Q[2,2] = 0.1
rk.P *= 50

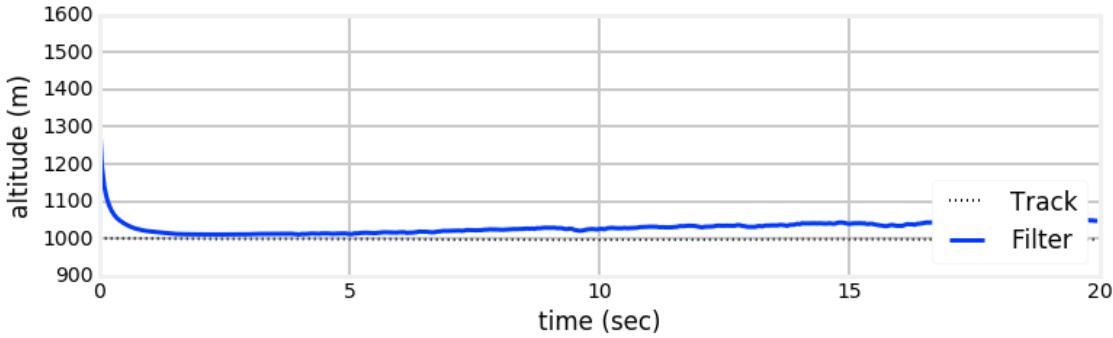
xs, track = [], []
for i in range(int(20/dt)):
    z = radar.get_range()
    track.append((radar.pos, radar.vel, radar.alt))

    rk.update(array([z]), HJacobian_at, hx)
    xs.append(rk.x)
    rk.predict()

xs = asarray(xs)
track = asarray(track)
time = np.arange(0, len(xs)*dt, dt)
ekf_internal.plot_radar(xs, track, time)

```





### 11.3 Using SymPy to compute Jacobians

Depending on your experience with derivatives you may have found the computation of the Jacobian difficult. Even if you found it easy, a slightly more difficult problem easily leads to very difficult computations.

As explained in Appendix A, we can use the SymPy package to compute the Jacobian for us.

```
In [8]: import sympy
sympy.init_printing(use_latex=True)

x, x_vel, y = sympy.symbols('x, x_vel y')

H = sympy.Matrix([sympy.sqrt(x**2 + y**2)])

state = sympy.Matrix([x, x_vel, y])
H.jacobian(state)
```

Out[8] :

$$\begin{bmatrix} \frac{x}{\sqrt{x^2+y^2}} & 0 & \frac{y}{\sqrt{x^2+y^2}} \end{bmatrix}$$

This result is the same as the result we computed above, and with much less effort on our part!

### 11.4 Robot Localization

So, time to try a real problem. I warn you that this is far from a simple problem. However, most books choose simple, textbook problems with simple answers, and you are left wondering how to implement a real world solution.

We will consider the problem of robot localization. We already implemented this in the **Unscented Kalman Filter** chapter, and I recommend you read that first. In this scenario we have a robot that is moving through a landscape with sensors that give range and bearings to various landmarks. This could be a self driving car using computer vision to identify trees, buildings, and other landmarks. It might be one of those small robots that vacuum your house, or a robot in a warehouse.

Our robot has 4 wheels configured the same as an automobile. It maneuvers by pivoting the front wheels. This causes the robot to pivot around the rear axle while moving forward. This is nonlinear behavior which we will have to model.

The robot has a sensor that gives it approximate range and bearing to known targets in the landscape. This is nonlinear because computing a position from a range and bearing requires square roots and trigonometry.

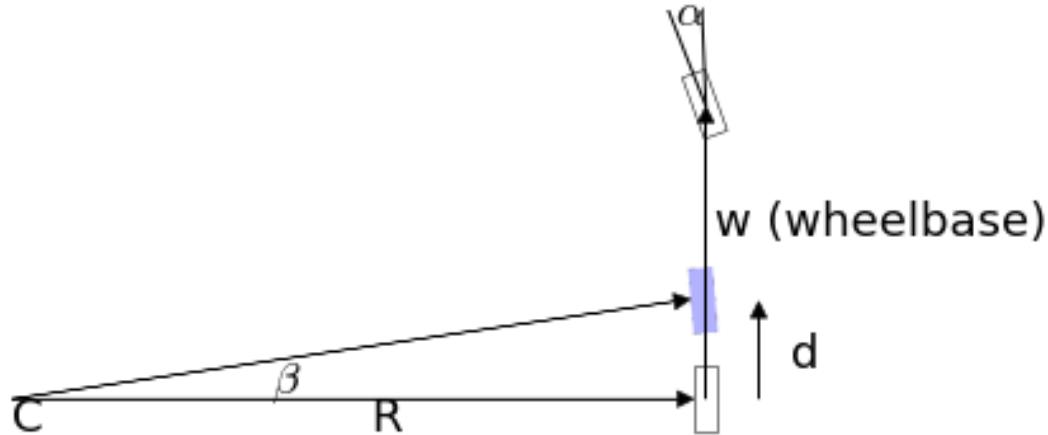
Both the process model and measurement models are nonlinear. The UKF accommodates both, so we provisionally conclude that the UKF is a viable choice for this problem.

### 11.4.1 Robot Motion Model

At a first approximation an automobile steers by pivoting the front tires while moving forward. The front of the car moves in the direction that the wheels are pointing while pivoting around the rear tires. This simple description is complicated by issues such as slippage due to friction, the differing behavior of the rubber tires at different speeds, and the need for the outside tire to travel a different radius than the inner tire. Accurately modeling steering requires a complicated set of differential equations.

For Kalman filtering, especially for lower speed robotic applications a simpler bicycle model has been found to perform well. This is a depiction of the model:

In [9]: `ekf_internal.plot_bicycle()`



In the **Unscented Kalman Filter** chapter we derived these equations describing for this model:

$$\begin{aligned}x &= x - R \sin(\theta) + R \sin(\theta + \beta) \\y &= y + R \cos(\theta) - R \cos(\theta + \beta) \\ \theta &= \theta + \beta\end{aligned}$$

where  $\theta$  is the robot's heading.

You do not need to understand this model in detail if you are not interested in steering models. The important thing to recognize is that our motion model is nonlinear, and we will need to deal with that with our Kalman filter.

### 11.4.2 Design the State Variables

For our robot we will maintain the position and orientation of the robot:

$$\mathbf{x} = \begin{bmatrix} x \\ y \\ \theta \end{bmatrix}$$

Our control input  $\mathbf{u}$  is the velocity  $v$  and steering angle  $\alpha$ :

$$\mathbf{u} = \begin{bmatrix} v \\ \alpha \end{bmatrix}$$

### 11.4.3 Design the System Model

In general we model our system as a nonlinear motion model plus noise.

$$\bar{x} = x + f(x, u) + \mathcal{N}(0, Q)$$

Using the motion model for a robot that we created above, we can expand this to

$$\overline{\begin{bmatrix} x \\ y \\ \theta \end{bmatrix}} = \begin{bmatrix} x \\ y \\ \theta \end{bmatrix} + \begin{bmatrix} -R \sin(\theta) + R \sin(\theta + \beta) \\ R \cos(\theta) - R \cos(\theta + \beta) \\ \beta \end{bmatrix}$$

We linearize this with a taylor expansion at  $x$ :

$$f(x, u) \approx \mathbf{x} + \frac{\partial f(\mathbf{x}, \mathbf{u})}{\partial x} \Big|_{\mathbf{x}, \mathbf{u}}$$

We replace  $f(x, u)$  with our state estimate  $\mathbf{x}$ , and the derivative is the Jacobian of  $f$ .

The Jacobian  $\mathbf{F}$  is

$$\mathbf{F} = \frac{\partial f(x, u)}{\partial x} = \begin{bmatrix} \frac{\partial \dot{x}}{\partial x} & \frac{\partial \dot{x}}{\partial y} & \frac{\partial \dot{x}}{\partial \theta} \\ \frac{\partial \dot{y}}{\partial x} & \frac{\partial \dot{y}}{\partial y} & \frac{\partial \dot{y}}{\partial \theta} \\ \frac{\partial \dot{\theta}}{\partial x} & \frac{\partial \dot{\theta}}{\partial y} & \frac{\partial \dot{\theta}}{\partial \theta} \end{bmatrix}$$

When we calculate these we get

$$\mathbf{F} = \begin{bmatrix} 1 & 0 & -R \cos(\theta) + R \cos(\theta + \beta) \\ 0 & 1 & -R \sin(\theta) + R \sin(\theta + \beta) \\ 0 & 0 & 1 \end{bmatrix}$$

We can double check our work with SymPy.

```
In [10]: import sympy
from sympy.abc import alpha, x, y, v, w, R, theta
from sympy import symbols, Matrix
sympy.init_printing(use_latex="mathjax", fontsize='16pt')
time = symbols('t')
d = v*time
```

```

beta = (d/w)*sympy.tan(alpha)
r = w/sympy.tan(alpha)

fxu = Matrix([[x-r*sympy.sin(theta) + r*sympy.sin(theta+beta)],
              [y+r*sympy.cos(theta)- r*sympy.cos(theta+beta)],
              [theta+beta]])
J = fxu.jacobian(Matrix([x, y, theta]))
J

```

Out[10]:

$$\begin{bmatrix} 1 & 0 & -\frac{w \cos(\theta)}{\tan(\alpha)} + \frac{w}{\tan(\alpha)} \cos\left(\frac{tv}{w} \tan(\alpha) + \theta\right) \\ 0 & 1 & -\frac{w \sin(\theta)}{\tan(\alpha)} + \frac{w}{\tan(\alpha)} \sin\left(\frac{tv}{w} \tan(\alpha) + \theta\right) \\ 0 & 0 & 1 \end{bmatrix}$$

That looks a bit complicated. We can use SymPy to substitute terms:

```

In [11]: # reduce common expressions
B, R = symbols('beta, R')
J = J.subs((d/w)*sympy.tan(alpha), B)
J.subs(w/sympy.tan(alpha), R)

```

Out[11]:

$$\begin{bmatrix} 1 & 0 & -R \cos(\theta) + R \cos(\beta + \theta) \\ 0 & 1 & -R \sin(\theta) + R \sin(\beta + \theta) \\ 0 & 0 & 1 \end{bmatrix}$$

In that form we can see that our computation of the Jacobian is correct.

Now we can turn our attention to the noise. Here, the noise is in our control input, so it is in control space. In other words, we command a specific velocity and steering angle, but we need to convert that into errors in  $x, y, \theta$ . In a real system this might vary depending on velocity, so it will need to be recomputed for every prediction. I will choose this as the noise model; for a real robot you will need to choose a model that accurately depicts the error in your system.

$$\mathbf{M} = \begin{bmatrix} \sigma_{vel}^2 & 0 \\ 0 & \sigma_\alpha^2 \end{bmatrix}$$

If this was a linear problem we would convert from control space to state space using the by now familiar  $\mathbf{FMF}^\top$  form. Since our motion model is nonlinear we do not try to find a closed form solution to this, but instead linearize it with a Jacobian which we will name  $\mathbf{V}$ .

$$\mathbf{V} = \frac{\partial f(x, u)}{\partial u} \begin{bmatrix} \frac{\partial \dot{x}}{\partial v} & \frac{\partial \dot{x}}{\partial \alpha} \\ \frac{\partial \dot{y}}{\partial v} & \frac{\partial \dot{y}}{\partial \alpha} \\ \frac{\partial \dot{\theta}}{\partial v} & \frac{\partial \dot{\theta}}{\partial \alpha} \end{bmatrix}$$

These partial derivatives become very difficult to work with. Let's compute them with SymPy.

```

In [12]: V = fxu.jacobian(Matrix([v, alpha]))
V = V.subs(sympy.tan(alpha)/w, 1/R)
V = V.subs(time*v/R, B)
V = V.subs(time*v, 'd')
V

```

Out[12] :

$$\begin{bmatrix} t \cos(\beta + \theta) & \frac{d \cos(\beta + \theta)}{\tan(\alpha)} (\tan^2(\alpha) + 1) - \frac{w \sin(\theta)}{\tan^2(\alpha)} (-\tan^2(\alpha) - 1) + \frac{w \sin(\beta + \theta)}{\tan^2(\alpha)} (-\tan^2(\alpha) - 1) \\ t \sin(\beta + \theta) & \frac{d \sin(\beta + \theta)}{\tan(\alpha)} (\tan^2(\alpha) + 1) + \frac{w \cos(\theta)}{\tan^2(\alpha)} (-\tan^2(\alpha) - 1) - \frac{w \cos(\beta + \theta)}{\tan^2(\alpha)} (-\tan^2(\alpha) - 1) \\ \frac{t}{R} & \frac{d}{w} (\tan^2(\alpha) + 1) \end{bmatrix}$$

This should give you an appreciation of how quickly the EKF become mathematically intractable.

This gives us the final form of our prediction equations:

$$\bar{\mathbf{x}} = \mathbf{x} + \begin{bmatrix} -R \sin(\theta) + R \sin(\theta + \beta) \\ R \cos(\theta) - R \cos(\theta + \beta) \\ \beta \end{bmatrix}$$

$$\bar{\mathbf{P}} = \mathbf{F} \mathbf{P} \mathbf{F}^T + \mathbf{V} \mathbf{M} \mathbf{V}^T$$

One final point. This form of linearization is not the only way to predict  $\mathbf{x}$ . For example, we could use a numerical integration technique like [Runge Kutta](#) to compute the position of the robot in the future. In fact, if the time step is relatively large you will have to do that. As I am sure you are realizing, things are not as cut and dried with the EKF as it was for the KF. For a real problem you have to very carefully model your system with differential equations and then determine the most appropriate way to solve that system. The correct approach depends on the accuracy you require, how nonlinear the equations are, your processor budget, and numerical stability concerns. These are all topics beyond the scope of this book.

#### 11.4.4 Design the Measurement Model

Now we need to design our measurement model. For this problem we are assuming that we have a sensor that receives a noisy bearing and range to multiple known locations in the landscape. The measurement model must convert the state  $[x \ y \ \theta]^T$  into a range and bearing to the landmark. Using  $p$  be the position of a landmark, the range  $r$  is

$$r = \sqrt{(p_x - x)^2 + (p_y - y)^2}$$

We assume that the sensor provides bearing relative to the orientation of the robot, so we must subtract the robot's orientation from the bearing to get the sensor reading, like so:

$$\phi = \arctan\left(\frac{p_y - y}{p_x - x}\right) - \theta$$

Thus our function is

$$\begin{aligned} \mathbf{x} &= h(x, p) && + \mathcal{N}(0, R) \\ &= \begin{bmatrix} \sqrt{(p_x - x)^2 + (p_y - y)^2} \\ \arctan\left(\frac{p_y - y}{p_x - x}\right) - \theta \end{bmatrix} && + \mathcal{N}(0, R) \end{aligned}$$

This is clearly nonlinear, so we need linearize  $h(x, p)$  at  $\mathbf{x}$  by taking its Jacobian. We compute that with SymPy below.

```
In [13]: px, py = symbols('p_x, p_y')
z = Matrix([[sympy.sqrt((px-x)**2 + (py-y)**2)],
            [sympy.atan2(py-y, px-x) - theta]])
z.jacobian(Matrix([x, y, theta]))
```

Out[13] :

$$\begin{bmatrix} \frac{-p_x+x}{\sqrt{(p_x-x)^2+(p_y-y)^2}} & \frac{-p_y+y}{\sqrt{(p_x-x)^2+(p_y-y)^2}} & 0 \\ -\frac{-p_y+y}{(p_x-x)^2+(p_y-y)^2} & -\frac{p_x-x}{(p_x-x)^2+(p_y-y)^2} & -1 \end{bmatrix}$$

Now we need to write that as a Python function. For example we might write:

In [14]: `from math import sqrt`

```
def H_of(x, landmark_pos):
    """ compute Jacobian of H matrix where h(x) computes
    the range and bearing to a landmark for state x """

    px = landmark_pos[0]
    py = landmark_pos[1]
    hyp = (px - x[0, 0])**2 + (py - x[1, 0])**2
    dist = sqrt(hyp)

    H = array(
        [[-(px - x[0, 0]) / dist, -(py - x[1, 0]) / dist, 0],
         [(py - x[1, 0]) / hyp, -(px - x[0, 0]) / hyp, -1]])
    return H
```

We also need to define a function that converts the system state into a measurement.

In [15]: `from math import atan2`

```
def Hx(x, landmark_pos):
    """ takes a state variable and returns the measurement
    that would correspond to that state.
    """
    px = landmark_pos[0]
    py = landmark_pos[1]
    dist = sqrt((px - x[0, 0])**2 + (py - x[1, 0])**2)

    Hx = array([[dist],
                [atan2(py - x[1, 0], px - x[0, 0]) - x[2, 0]]])
    return Hx
```

### 11.4.5 Design Measurement Noise

This is quite straightforward as we need to specify measurement noise in measurement space, hence it is linear. It is reasonable to assume that the range and bearing measurement noise is independent, hence

$$R = \begin{bmatrix} \sigma_{range}^2 & 0 \\ 0 & \sigma_{bearing}^2 \end{bmatrix}$$

### 11.4.6 Implementation

We will use FilterPy's `ExtendedKalmanFilter` class to implement the filter. Its `predict()` method uses the standard linear equations. Our process model is nonlinear, so we will have to override `predict()` with

our own version. I'll want to also use this class to simulate the robot, so I'll add a method `move()` that computes the position of the robot which both `predict()` and my simulation can call.

The matrices for the prediction step are quite large. While writing this code I made several errors before I finally got it working. I only found my errors by using SymPy's `evalf` function, which allows you to evaluate a SymPy `Matrix` with specific values for the variables. I decided to demonstrate this technique, and to eliminate a possible source of bugs, by using SymPy in the Kalman filter. You'll need to understand a couple of points.

First, `evalf` uses a dictionary to pass in the values you want to use. For example, if your matrix contains an `x` and `y`, you can write

```
M.evalf(subs={x:3, y:17})
```

to evaluate the matrix for `x=3` and `y=17`.

Second, `evalf` returns a `sympy.Matrix` object. Use `numpy.array(M).astype(float)` to convert it to a NumPy array. `numpy.array(M)` creates an array of type `object`, which is not what you want.

Here is the code for the EKF:

```
In [16]: from filterpy.kalman import ExtendedKalmanFilter as EKF
        from numpy import dot, array, sqrt
        class RobotEKF(EKF):
            def __init__(self, dt, wheelbase, std_vel, std_steer):
                EKF.__init__(self, 3, 2, 2)
                self.dt = dt
                self.wheelbase = wheelbase
                self.std_vel = std_vel
                self.std_steer = std_steer

                a, x, y, v, w, theta, time = symbols(
                    'a, x, y, v, w, theta, t')
                d = v*time
                beta = (d/w)*sympy.tan(a)
                r = w/sympy.tan(a)

                self.fxu = Matrix(
                    [[x-r*sympy.sin(theta)+r*sympy.sin(theta+beta)],
                     [y+r*sympy.cos(theta)-r*sympy.cos(theta+beta)],
                     [theta+beta]])

                self.F_j = self.fxu.jacobian(Matrix([x, y, theta]))
                self.V_j = self.fxu.jacobian(Matrix([v, a]))

                # save dictionary and it's variables for later use
                self.subs = {x: 0, y: 0, v:0, a:0,
                            time:dt, w:wheelbase, theta:0}
                self.x_x, self.x_y, = x, y
                self.v, self.a, self.theta = v, a, theta

            def predict(self, u=0):
                self.x = self.move(self.x, u, self.dt)

                self.subs[self.theta] = self.x[2, 0]
                self.subs[self.v] = u[0]
```

```

    self.subs[self.a] = u[1]

F = array(self.F_j.evalf(subs=self.subs)).astype(float)
V = array(self.V_j.evalf(subs=self.subs)).astype(float)

# covariance of motion noise in control space
M = array([[self.std_vel*u[0]**2, 0],
           [0, self.std_steer**2]])

self.P = dot(F, self.P).dot(F.T) + dot(V, M).dot(V.T)

def move(self, x, u, dt):
    hdg = x[2, 0]
    vel = u[0]
    steering_angle = u[1]
    dist = vel * dt

    if abs(steering_angle) > 0.001: # is robot turning?
        beta = (dist / self.wheelbase) * tan(steering_angle)
        r = self.wheelbase / tan(steering_angle) # radius

        dx = np.array([[-r*sin(hdg) + r*sin(hdg + beta)],
                      [r*cos(hdg) - r*cos(hdg + beta)],
                      [beta]])
    else: # moving in straight line
        dx = np.array([[dist*cos(hdg)],
                      [dist*sin(hdg)],
                      [0]])
    return x + dx

```

Now we have another issue to handle. The residual is notionally computed as  $y = z - h(x)$  but this will not work because our measurement contains an angle in it. Suppose  $z$  has a bearing of  $1^\circ$  and  $h(x)$  has a bearing of  $359^\circ$ . Naively subtracting them would yield a bearing difference of  $-358^\circ$ , which will throw off the computation of the Kalman gain. The correct angle difference in this case is  $2^\circ$ . So we will have to write code to correctly compute the bearing residual.

```
In [17]: def residual(a, b):
    """ compute residual (a-b) between measurements containing
    [range, bearing]. Bearing is normalized to [-pi, pi]"""
    y = a - b
    y[1] = y[1] % (2 * np.pi) # force in range [0, 2 pi)
    if y[1] > np.pi: # move to [-pi, pi)
        y[1] -= 2 * np.pi
    return y
```

The rest of the code runs the simulation and plots the results, and shouldn't need too much comment by now. I create a variable `landmarks` that contains the coordinates of the landmarks. I update the simulated robot position 10 times a second, but run the EKF only once. This is for two reasons. First, we are not using Runge Kutta to integrate the differential equations of motion, so a narrow time step allows our simulation to be more accurate. Second, it is fairly normal in embedded systems to have limited processing speed. This forces you to run your Kalman filter only as frequently as absolutely needed.

```
In [18]: from filterpy.stats import plot_covariance_ellipse
from math import sqrt, tan, cos, sin, atan2
```

```

import matplotlib.pyplot as plt

dt = 1.0

def z_landmark(lmark, sim_pos, std_rng, std_brg):
    x, y = sim_pos[0, 0], sim_pos[1, 0]
    d = np.sqrt((lmark[0] - x)**2 + (lmark[1] - y)**2)
    a = atan2(lmark[1] - y, lmark[0] - x) - sim_pos[2, 0]
    z = np.array([[d + randn()*std_rng],
                  [a + randn()*std_brg]])
    return z

def ekf_update(ekf, z, landmark):
    ekf.update(z, HJacobian=H_of, Hx=Hx,
               residual=residual,
               args=(landmark), hx_args=(landmark))

def run_localization(landmarks, std_vel, std_steer,
                     std_range, std_bearing,
                     step=10, ellipse_step=20, ylim=None):
    ekf = RobotEKF(dt, wheelbase=0.5, std_vel=std_vel,
                   std_steer=std_steer)
    ekf.x = array([[2, 6, .3]]).T # x, y, steer angle
    ekf.P = np.diag([.1, .1, .1])
    ekf.R = np.diag([std_range**2, std_bearing**2])

    sim_pos = ekf.x.copy() # simulated position
    # steering command (vel, steering angle radians)
    u = array([1.1, .01])

    plt.scatter(landmarks[:, 0], landmarks[:, 1],
                marker='s', s=60)

    track = []
    for i in range(200):
        sim_pos = ekf.move(sim_pos, u, dt/10.) # simulate robot
        track.append(sim_pos)

        if i % step == 0:
            ekf.predict(u=u)

        if i % ellipse_step == 0:
            plot_covariance_ellipse(
                (ekf.x[0,0], ekf.x[1,0]), ekf.P[0:2, 0:2],
                std=6, facecolor='k', alpha=0.3)

        x, y = sim_pos[0, 0], sim_pos[1, 0]
        for lmark in landmarks:
            z = z_landmark(lmark, sim_pos,
                            std_range, std_bearing)
            ekf_update(ekf, z, lmark)

        if i % ellipse_step == 0:

```

```

        plot_covariance_ellipse(
            (ekf.x[0,0], ekf.x[1,0]), ekf.P[0:2, 0:2],
            std=6, facecolor='g', alpha=0.8)
track = np.array(track)
plt.plot(track[:, 0], track[:,1], color='k', lw=2)
plt.axis('equal')
plt.title("EKF Robot localization")
if ylim is not None: plt.ylim(*ylim)
plt.show()
return ekf

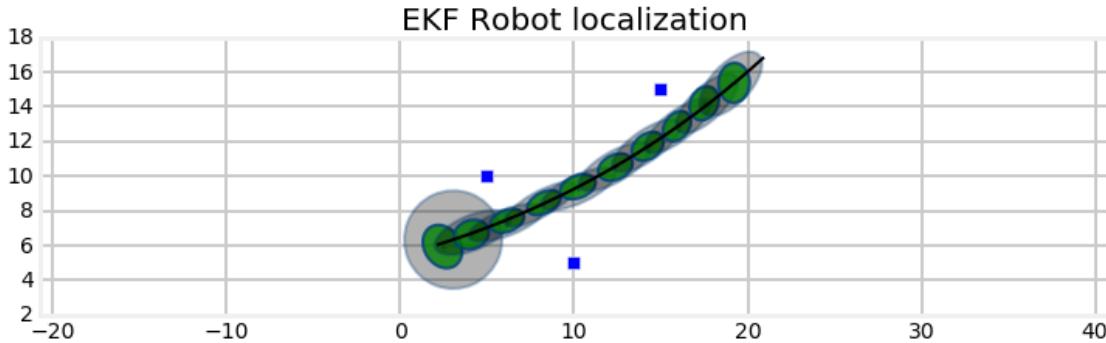
```

In [19]: landmarks = array([[5, 10], [10, 5], [15, 15]])

```

ekf = run_localization(
    landmarks, std_vel=0.1, std_steer=np.radians(1),
    std_range=0.3, std_bearing=0.1)
print('Final P:', ekf.P.diagonal())

```



Final P: [ 0.024 0.042 0.002]

I have plotted the landmarks as solid squares. The path of the robot is drawn with black line. The covariance ellipses for the predict step is light gray, and the covariances of the update are shown in green. To make them visible at this scale I have set the ellipse boundary at  $6\sigma$ .

From this we can see that there is a lot of uncertainty added by our motion model, and that most of the error is in the direction of motion. We can see that from the shape of the blue ellipses. After a few steps we can see that the filter incorporates the landmark measurements.

I used the same initial conditions and landmark locations in the UKF chapter. You can see both in the plot and in the printed final value for  $\mathbf{P}$  that the UKF achieves much better accuracy in terms of the error ellipse. The black solid line denotes the robot's actual path. Both perform roughly as well as far as their estimate for  $\mathbf{x}$  is concerned.

Now lets add another landmark.

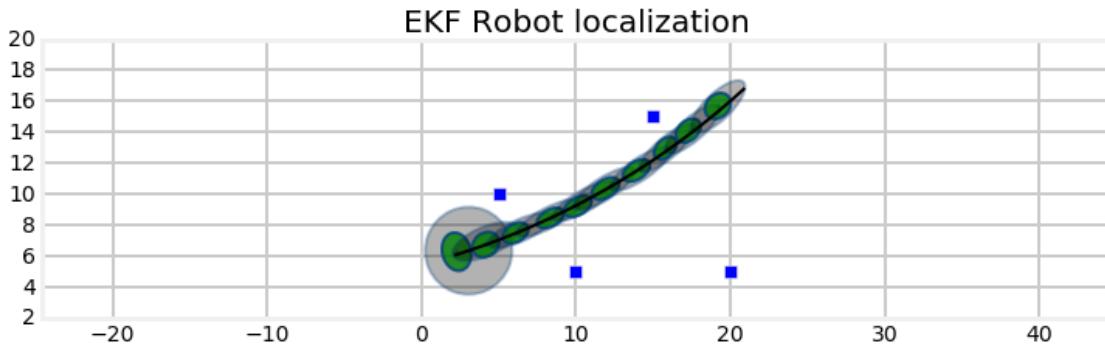
In [20]: landmarks = array([[5, 10], [10, 5], [15, 15], [20, 5]])

```

ekf = run_localization(
    landmarks, std_vel=0.1, std_steer=np.radians(1),

```

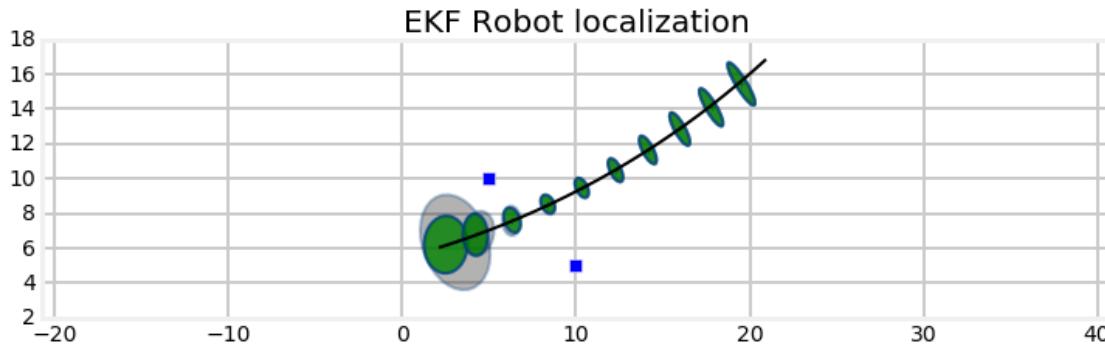
```
    std_range=0.3, std_bearing=0.1)
plt.show()
print('Final P:', ekf.P.diagonal())
```



```
Final P: [ 0.019  0.022  0.002]
```

The uncertainty in the estimates near the end of the track are smaller with the additional landmark. We can see the fantastic effect that multiple landmarks has on our uncertainty by only using the first two landmarks.

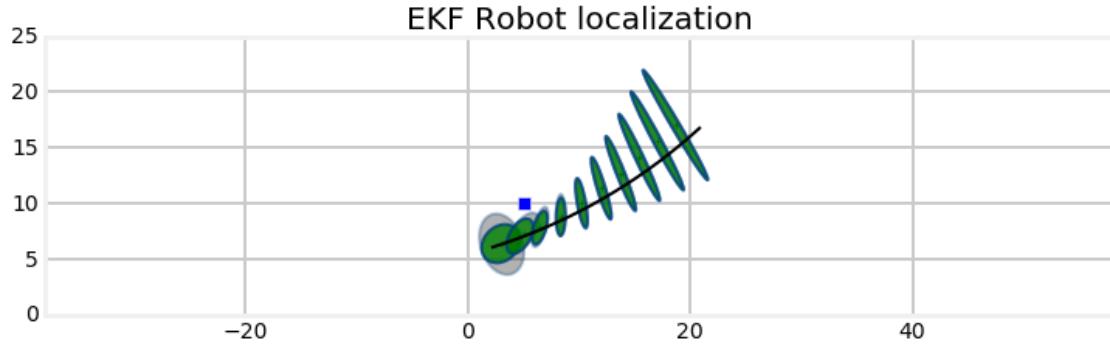
```
In [21]: ekf = run_localization(
    landmarks[0:2], std_vel=1.e-10, std_steer=1.e-10,
    std_range=1.4, std_bearing=.05)
print('Final P:', ekf.P.diagonal())
```



```
Final P: [ 0.02   0.046   0.    ]
```

The estimate quickly diverges from the robot's path after passing the landmarks. The covariance also grows quickly. Let's see what happens with only one landmark:

```
In [22]: ekf = run_localization(
    landmarks[0:1], std_vel=1.e-10, std_steer=1.e-10,
    std_range=1.4, std_bearing=.05)
print('Final P:', ekf.P.diagonal())
```

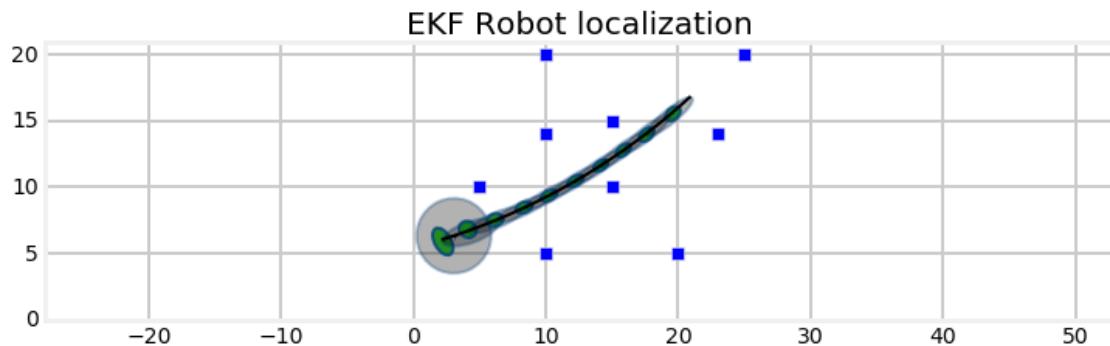


```
Final P: [ 0.284  0.74   0.003]
```

As you probably suspected, only one landmark produces a very bad result. Conversely, a large number of landmarks allows us to make very accurate estimates.

```
In [23]: landmarks = array([[5, 10], [10, 5], [15, 15], [20, 5], [15, 10],
                           [10, 14], [23, 14], [25, 20], [10, 20]])

ekf = run_localization(
    landmarks, std_vel=0.1, std_steer=np.radians(1),
    std_range=0.3, std_bearing=0.1, ylim=(0, 21))
print('Final P:', ekf.P.diagonal())
```



```
Final P: [ 0.009  0.008  0.001]
```

### 11.4.7 Discussion

I said that this was a ‘real’ problem, and in some ways it is. I’ve seen alternative presentations that used robot motion models that led to much easier Jacobians. On the other hand, my model of a automobile’s movement is itself simplistic in several ways. First, it uses a bicycle model. A real car has two sets of tires, and each travels on a different radius. The wheels do not grip the surface perfectly. I also assumed that the

robot responds instantaneously to the control input. Sebastian Thrun writes in [Probabilistic Robots](#) that simplified models are justified because the filters perform well when used to track real vehicles. The lesson here is that while you have to have a reasonably accurate nonlinear model, it does not need to be perfect to operate well. As a designer you will need to balance the fidelity of your model with the difficulty of the math and the computation required to implement the equations.

Another way in which this problem was simplistic is that we assumed that we knew the correspondance between the landmarks and measurements. But suppose we are using radar - how would we know that a specific signal return corresponded to a specific building in the local scene? This question hints at SLAM algorithms - simultaneous localization and mapping. SLAM is not the point of this book, so I will not elaborate on this topic.

## 11.5 UKF vs EKF

I implemented this tracking problem using the UKF in the previous chapter. The difference in implementation should be very clear. Computing the Jacobians for the state and measurement models was not trivial despite a rudimentary motion model. I am justified in using this model because the research resulting from the DARPA car challenges has shown that it works well in practice. A different problem could result in a Jacobian which is difficult or impossible to derive analytically. In contrast, the UKF only requires you to provide a function that computes the system motion model and another for the measurement model.

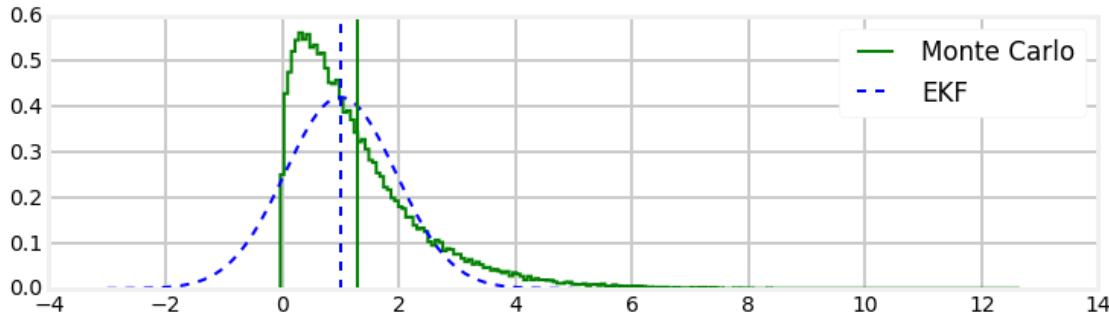
There are many cases where the Jacobian cannot be found analytically. The details are beyond the scope of this book, but you will have to use numerical methods to compute the Jacobian. That is a very nontrivial undertaking, and you will spend a significant portion of a master's degree at a STEM school learning techniques to handle such situations. Even then you'll likely only be able to solve problems related to your field - an aeronautical engineer learns a lot about Navier Stokes equations, but not much about modelling chemical reaction rates.

So, UKFs are easy. Are they accurate? In practice they often perform better than the EKF. You can find plenty of research papers that prove that the UKF outperforms the EKF in various problem domains. It's not hard to understand why this would be true. The EKF works by linearizing the system model and measurement model at a single point, and the UKF uses  $2n + 1$  points.

Let's look at a specific example. Take  $f(x) = x^3$  and pass a Gaussian distribution through it. I will compute an accurate answer using a monte carlo simulation. I generate 50,000 points randomly distributed according to the Gaussian, pass each through  $f(x)$ , then compute the mean and variance of the result.

First, let's see how the EKF fairs. The EKF linearizes the function by taking the derivative and evaluating it the mean  $x$  to get the slope tangent to the function at that point. This slope becomes the linear function that we use to transform the Gaussian. Here is a plot of that.

```
In [24]: import nonlinear_plots
nonlinear_plots.plot_ekf_vs_mc()
```

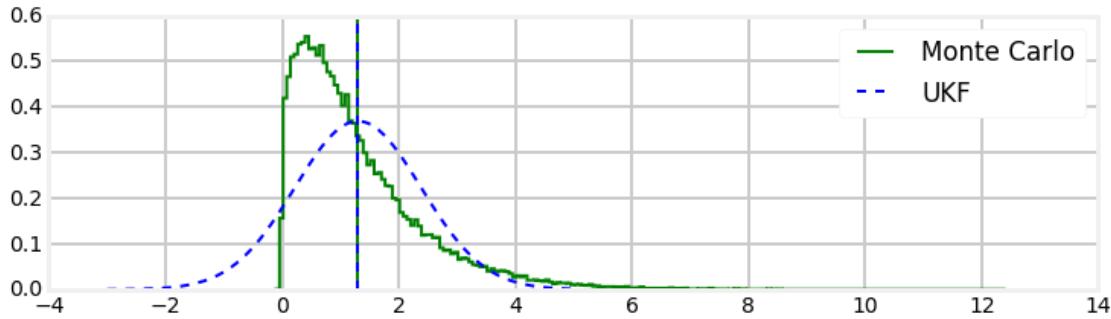


```
actual mean=1.29, std=1.13
EKF    mean=1.00, std=0.95
```

We can see from both the graph and the print out at the bottom that the EKF has introduced quite a bit of error.

In contrast, here is the performance of the UKF:

```
In [25]: nonlinear_plots.plot_ukf_vs_mc(alpha=0.001, beta=3., kappa=1.)
```



```
actual mean=1.30, std=1.13
UKF    mean=1.30, std=1.08
```

Here we can see that the computation of the UKF's mean is accurate to 2 decimal places. The standard deviation is slightly off, but you can also fine tune how the UKF computes the distribution by using the  $\alpha$ ,  $\beta$ , and  $\gamma$  parameters for generating the sigma points. Here I used  $\alpha = 0.001$ ,  $\beta = 3$ , and  $\gamma = 1$ . Feel free to modify them in the function call to see the result. You should be able to get better results than I did. However, avoid over-tuning the UKF for a specific test. It may perform better for your test case, but worse in general.

This is a contrived example, but as I said the literature is filled with detailed studies of real world problems that exhibit similar performance differences between the two filters.



# Chapter 12

## Particle Filters

### 12.1 Motivation

Here is our problem. We have moving objects that we want to track. Maybe the objects are fighter jets and missiles, or maybe we are tracking people playing cricket in a field. It doesn't really matter. Which of the filters that we have learned can handle this problem? Unfortunately, none of them are ideal. Let's think about the characteristics of this problem.

- **multimodal:** We want to track zero, one, or more than one object simultaneously.
- **occlusions:** One object can hide another, resulting in one measurement for multiple objects.
- **nonlinear behavior:** Aircraft are buffeted by winds, balls move in parabolas, and people collide into each other.
- **nonlinear measurements:** Radar gives us the distance to an object. Converting that to an (x,y,z) coordinate requires a square root, which is nonlinear.
- **non-Gaussian noise:** as objects move across a background the computer vision can mistake part of the background for the object.
- **continuous:** the object's position and velocity (i.e. the state space) can smoothly vary over time.
- **multivariate:** we want to track several attributes, such as position, velocity, turn rates, etc.
- **unknown process model:** we may not know the process model of the system

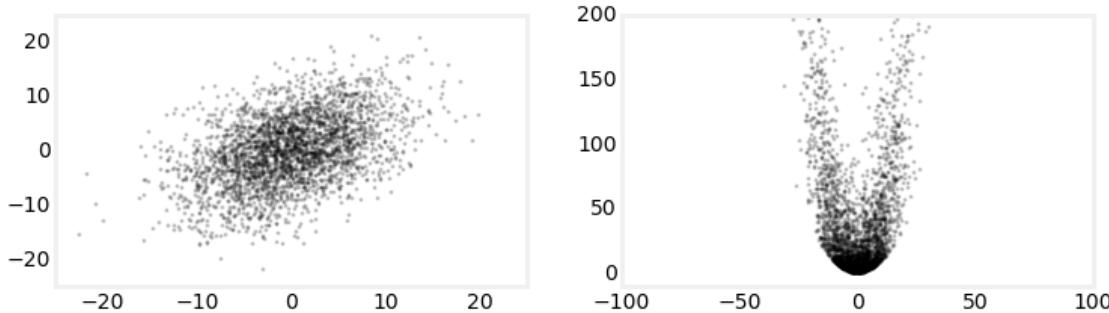
None of the filters we have learned work well with all of these constraints.

- **Discrete Bayes filter:** This has most of the attributes. It is multimodal, can handle nonlinear measurements, and can be extended to work with nonlinear behavior. However, it is discrete and univariate.
- **Kalman filter:** The Kalman filter produces optimal estimates for unimodal linear systems with Gaussian noise. None of these are true for our problem.
- **Unscented Kalman filter:** The UKF handles nonlinear, continuous, multivariate problems. However, it is not multimodal nor does it handle occlusions. It can handle noise that is modestly non-Gaussian, but does not do well with distributions that are very non-Gaussian or problems that are very nonlinear.
- **Extended Kalman filter:** The EKF has the same strengths and limitations as the UKF, except that is it even more sensitive to strong nonlinearities and non-Gaussian noise.

## 12.2 Monte Carlo Sampling

In the UKF chapter I generated a plot similar to this to illustrate the effects of nonlinear systems on Gaussians:

```
In [2]: import pf_internal
pf_internal.plot_monte_carlo_ukf()
```



The left plot shows 3,000 points normally distributed based on the Gaussian

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 32 & 15 \\ 15 & 40 \end{bmatrix}$$

The right plots shows these points passed through this set of equations:

$$\begin{aligned} x &= x + y \\ y &= 0.1x^2 + y^2 \end{aligned}$$

Using a finite number of randomly sampled points to compute a result is called a Monte Carlo (MC) method. The idea is simple. Generate enough points to get a representative sample of the problem, run the points through the system you are modeling, and then compute the results on the transformed points.

In a nutshell this is what particle filtering does. The Bayesian filter algorithm we have been using throughout the book is applied to thousands of particles, where each particle represents a possible state for the system. We extract the estimated state from the thousands of particles using weighted statistics of the particles.

## 12.3 Generic Particle Filter Algorithm

### 1. Randomly generate a bunch of particles

Particles can have position, heading, and/or whatever other state variable you need to estimate. Each has a weight (probability) indicating how likely it matches the actual state of the system. Initialize each with the same weight.

### 2. Predict next state of the particles

Move the particles based on how you predict the real system is behaving.

### 3. Update

Update the weighting of the particles based on the measurement. Particles that closely match the measurements are weighted higher than particles which don't match the measurements very well.

### 4. Resample

Discard highly improbable particle and replace them with copies of the more probable particles.

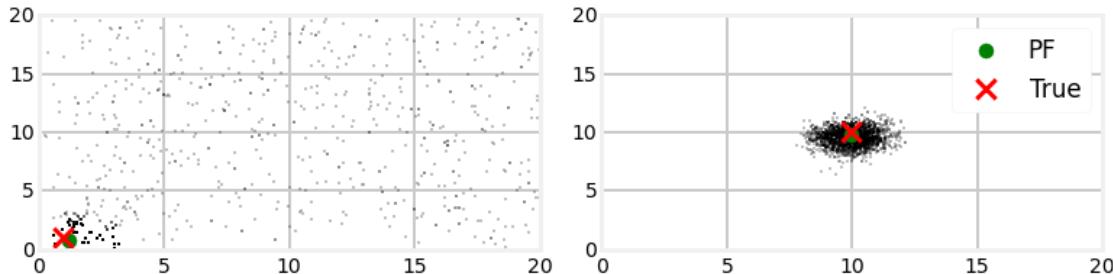
### 5. Compute Estimate

Optionally, compute weighted mean and covariance of the set of particles to get a state estimate.

This naive algorithm has practical difficulties which we will need to overcome, but this is the general idea. Let's see an example. I wrote a particle filter for the robot localization problem from the UKF and EKF chapters. The robot has steering and velocity control inputs. It has sensors that measures distance to visible landmarks. Both the sensors and control mechanism have noise in them, and we need to track the robot's position.

Here I run a particle filter and plotted the positions of the particles. The plot on the left is after one iteration, and on the right is after 10. The red 'X' shows the actual position of the robot, and the large circle is the computed weighted mean position.

In [3]: `pf_internal.show_two_pf_plots()`



If you are viewing this in a browser, this animation shows the entire sequence:

After the first iteration the particles are still largely randomly scattered around the map, but you can see that some have already collected near the robot's position. The computed mean is quite close to the robot's position. This is because each particle is weighted based on how closely it matches the measurement. The robot is near (1,1), so particles that are near (1, 1) will have a high weight because they closely match the measurements. Particles that are far from the robot will not match the measurements, and thus have a very low weight. The estimated position is computed as the weighted mean of positions of the particles. Particles near the robot contribute more to the computation so the estimate is quite accurate.

Several iterations later you can see that all the particles have clustered around the robot. This is due to the resampling step. Resampling discards particles that are very improbable (very low weight) and replaces them with particles with higher probability.

I haven't fully shown why this works nor fully explained the algorithms for particle weighting and resampling, but it should make intuitive sense. Make a bunch of random particles, move them so they 'kind of' follow the robot, weight them according to how well they match the measurements, only let the likely ones live. It seems like it should work, and it does.

## 12.4 Probability distributions via Monte Carlo

Suppose we want to know the area under the curve  $y = e^{\sin(x)}$  in the interval  $[0, \pi]$ . The area is computed with the definite integral  $\int_0^\pi e^{\sin(x)} dx$ . As an exercise, go ahead and find the answer; I'll wait.

If you are wise you did not take that challenge;  $e^{\sin(x)}$  cannot be integrated analytically. The world is filled with equations which we cannot integrate. For example, consider calculating the luminosity of an object. An object reflects some of the light that strike it. Some of the reflected light bounces off of other objects and restrikes the original object, increasing the luminosity. This creates a recursive integral. Good luck with that one.

However, integrals are trivial to compute using a Monte Carlo technique. To find the area under a curve create a bounding box that contains the curve in the desired interval. Generate randomly positioned point within the box, and compute the ratio of points that fall under the curve vs the total number of points. For example, if 40% of the points are under the curve and the area of the bounding box is 1, then the area under the curve is approximately 0.4. As you tend towards infinite points you can achieve any arbitrary precision. In practice, a few thousand points will give you a fairly accurate result.

You can use this technique to numerically integrate a function of any arbitrary difficulty. this includes non-integrable and noncontinuous functions. This technique was invented by Stanley Ulam at Los Alamos National Laboratory to allow him to perform computations for nuclear reactions which were unsolvable on paper.

Let's compute  $\pi$  by finding the area of a circle. We will define a circle with a radius of 1, and bound it in a square. The side of the square has length 2, so the area is 4. We generate a set of uniformly distributed random points within the box, and count how many fall inside the circle. The area of the circle is computed as the area of the box times the ratio of points inside the circle vs. the total number of points. Finally, we know that  $A = \pi r^2$ , so we compute  $\pi = A/r^2$ .

We start by creating the points.

```
N = 20000
pts = uniform(-1, 1, (N, 2))
```

A point is inside a circle if its distance from the center of the circle is less than or equal to the radius. We compute the distance with `numpy.linalg.norm`. If you aren't familiar with this, the norm is the magnitude of a vector. Since vectors start at (0, 0) calling norm will compute the point's distance from the origin.

```
dist = np.linalg.norm(pts, axis=1)
```

Next we compute which of these distances fit the criteria. This code returns a bool array that contains `True` if it meets the condition `dist <= 1`:

```
in_circle = dist <= 1
```

All that is left is to count the points inside the circle, compute pi, and plot the results. I've put it all in one cell so you can experiment with alternative values for `N`, the number of points.

```
In [4]: import matplotlib.pyplot as plt
import numpy as np
from numpy.random import uniform

N = 20000 # number of points
radius = 1
area = (2*radius)**2
```

```

pts = uniform(-1, 1, (N, 2))

# distance from (0,0)
dist = np.linalg.norm(pts, axis=1)
in_circle = dist <= 1

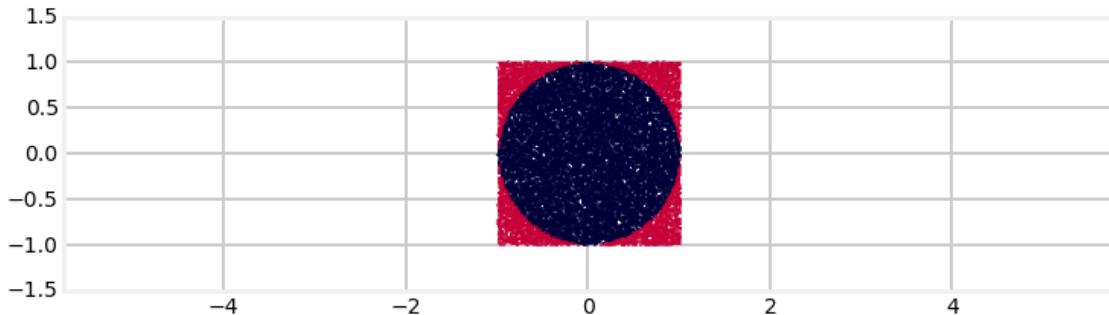
pts_in_circle = np.count_nonzero(in_circle)
pi = area * (pts_in_circle / N)

# plot results
plt.scatter(pts[in_circle,0], pts[in_circle,1],
            marker=',', edgecolor='k', s=1)
plt.scatter(pts[~in_circle,0], pts[~in_circle,1],
            marker=',', edgecolor='r', s=1)
plt.axis('equal')

print('mean pi(N={})= {:.4f}'.format(N, pi))
print('err pi(N={})= {:.4f}'.format(N, np.pi-pi))

mean pi(N=20000)= 3.1684
err pi(N=20000)= -0.0268

```



This insight leads us to the realization that we can use Monte Carlo to compute the probability density of any probability distribution. For example, suppose we have this Gaussian:

```
In [5]: from filterpy.stats import plot_gaussian_pdf
plot_gaussian_pdf(mean=2, variance=3);
```

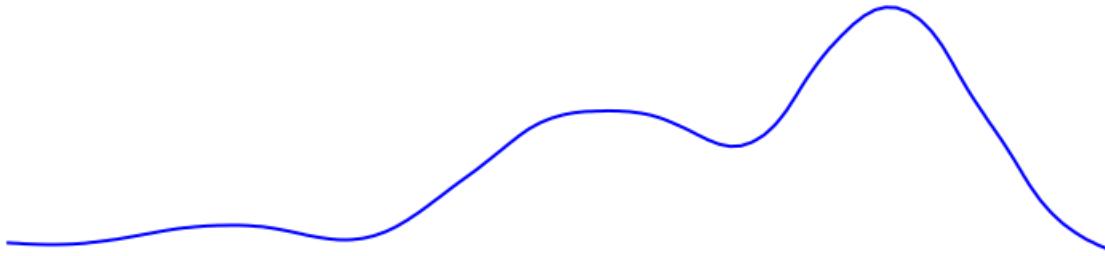


The probability density function (PDF) is the probability that the random value falls between 2 values. For example, we may want to know the probability of  $x$  being between 0 and 2 in the graph above. This is a continuous function, so we need to take the integral to find the area under the curve, as the area is equal to the probability for that range of values to occur.

$$P[a \leq X \leq b] = \int_a^b f_X(x) dx$$

It is easy to compute this integral for a Gaussian. But real life is not so easy. For example, the plot below shows a possible and realistic probability distribution. There is no way to analytically describe an arbitrary curve, let alone integrate it analytically.

In [6]: `pf_internal.plot_random_pd()`



We can use Monte Carlo methods to compute any integral. The PDF is computed with an integral, hence we can compute the PDF of this curve using Monte Carlo.

## 12.5 The Particle Filter

All of this brings us to the particle filter. Consider tracking a robot or a car in an urban environment. For consistency I will use the robot localization problem from the EKF and UKF chapters. In this problem we tracked a robot that had a sensor measures the range and bearing to landmarks.

Particle filters are a family of algorithms. I'm presenting a specific form of a particle filter that is intuitive to grasp and relates to the problems we have studied in this book. This will leave a few of the steps seeming a bit 'magical' since I haven't offered a full explanation. That will follow later in the chapter.

Taking insight from the discussion in the previous section we start by creating several thousand particles. Each particle has a position that represents a possible belief of where the robot is in the scene, and perhaps a heading and velocity. Suppose that we have no knowledge of the location of the robot. We would want to scatter the particles uniformly over the entire scene. If you think of all of the particles representing a probability distribution, locations where there are more particles represent a higher belief, and locations with fewer particles represents a lower belief. If there was a large clump of particles near a specific location that would imply that we were more certain that the robot is there.

Each particle needs a weight - ideally the probability that it represents the true position of the robot. This probability is rarely computable, so we only require it be proportional to that probability, which is computable. At initialization we have no reason to favor one particle over another, so we would typically

assign a weight of  $1/N$ , for  $N$  particles. We did the same thing in the discrete Bayes chapter. When we initialized the filter we assigned a probability of  $1/N$  to each hallway position. We use  $1/N$  so that the sum of all probabilities equals one.

The combination of particles and weights forms the probability distribution for our problem. Think back to the Discrete Bayes chapter. In that chapter we modeled positions in a hallway as discrete and uniformly spaced. This is very similar except the particles are randomly distributed in a continuous space rather than constrained to discrete locations. In this problem the robot can move on a plane of some arbitrary dimension, with the lower right corner at (0,0).

To track our robot we need to maintain states for x, y, and heading. We will store  $N$  particles in a  $(N, 3)$  shaped array. The three columns contain x, y, and heading, in that order.

If you are passively tracking something (no control input), then you would need to include velocity in the state and use that estimate to make the prediction. When using particle filters I minimize the dimensionality of the state. More dimensions requires exponentially more particles to form a good estimate.

Here is code to create either a uniform or Gaussian distribution of particles over a given region for this problem:

```
In [7]: from numpy.random import uniform

def create_uniform_particles(x_range, y_range, hdg_range, N):
    particles = np.empty((N, 3))
    particles[:, 0] = uniform(x_range[0], x_range[1], size=N)
    particles[:, 1] = uniform(y_range[0], y_range[1], size=N)
    particles[:, 2] = uniform(hdg_range[0], hdg_range[1], size=N)
    particles[:, 2] %= 2 * np.pi
    return particles

def create_gaussian_particles(mean, std, N):
    particles = np.empty((N, 3))
    particles[:, 0] = mean[0] + (randn(N) * std[0])
    particles[:, 1] = mean[1] + (randn(N) * std[1])
    particles[:, 2] = mean[2] + (randn(N) * std[2])
    particles[:, 2] %= 2 * np.pi
    return particles
```

For example:

```
In [8]: create_uniform_particles((0,1), (0,1), (0, 5), 4)
```

```
Out[8]: array([[ 0.772,  0.336,  3.319],
   [ 0.333,  0.34 ,  3.437],
   [ 0.6  ,  0.274,  3.995],
   [ 0.054,  0.022,  4.006]])
```

### 12.5.1 Predict Step

The predict step in the Bayes algorithm uses the process model to update the belief in the system state. How would we do that with particles? Each particle represents a possible position for the robot. Suppose we send a command to the robot to move 0.1 meters while turning by 0.007 radians. We could move each particle by this amount. If we did that we would soon run into a problem. The robot's controls are not perfect so it will not move exactly as commanded. Therefore we need to add noise to the particle's movements to have a reasonable chance of capturing the actual movement of the robot. If you do not model the uncertainty in

the system the particle filter will not correctly model the probability distribution of our belief in the robot's position.

```
In [9]: def predict(particles, u, std, dt=1.):
    """ move according to control input u (heading change, velocity)
    with noise Q (std heading change, std velocity) """
    N = len(particles)
    # update heading
    particles[:, 2] += u[0] + (randn(N) * std[0])
    particles[:, 2] %= 2 * np.pi

    # move in the (noisy) commanded direction
    dist = (u[1] * dt) + (randn(N) * std[1])
    particles[:, 0] += np.cos(particles[:, 2]) * dist
    particles[:, 1] += np.sin(particles[:, 2]) * dist
```

### 12.5.2 Update Step

Next we get a set of measurements - one for each landmark currently in view. How should these measurements be used to alter our probability distribution as modeled by the particles?

Think back to the **Discrete Bayes** chapter. In that chapter we modeled positions in a hallway as discrete and uniformly spaced. We assigned a probability to each position. When a new measurement came in we multiplied the current probability of that position by the probability that the measurement matched that location:

```
def update(map_, belief, z, prob_correct):
    scale = prob_correct / (1. - prob_correct)
    for i, val in enumerate(map_):
        if val == z:
            belief[i] *= scale
    normalize(belief)
```

We do the same with our particles. Each particle has a position and a weight which estimates how well it matches the measurement. Normalizing the weights so they sum to one turns them into a probability distribution. The particles those that are closest to the robot will generally have a higher weight than ones far from the robot.

```
In [10]: def update(particles, weights, z, R, landmarks):
    weights.fill(1.)
    for i, landmark in enumerate(landmarks):
        distance = np.linalg.norm(particles[:, 0:2] - landmark, axis=1)
        weights *= scipy.stats.norm(distance, R).pdf(z[i])

    weights += 1.e-300      # avoid round-off to zero
    weights /= sum(weights) # normalize
```

In the literature this part of the algorithm is called Sequential Importance Sampling, or SIS. The equation for the weights is called the importance density. I will give these theoretical underpinnings in a following section. For now I hope that this makes intuitive sense. If we weight the particles according to how well they match the measurements they are probably a good sample for the probability distribution of the system after incorporating the measurements. Theory proves this is so. Different problems will need to tackle this step in slightly different ways but this is the general idea.

### 12.5.3 Computing the State Estimate

In most applications you will want to know the estimated state after each update, but the filter consists of nothing but a collection of particles. Assuming that we are tracking one object (i.e. it is unimodal) we can compute the mean of the estimate as the sum of the weighted values of the particles.

$$\mu = \frac{1}{N} \sum_{i=1}^N w^i x^i$$

Here I adopt the notation  $x^i$  to indicate the  $i^{th}$  particle. A superscript is used because we often need to use subscripts to denote time steps the  $k^{th}$  or  $k+1^{th}$  particle, yielding the unwieldy  $x_{k+1}^i$ .

This computes both the mean and variance of the particles:

```
In [11]: def estimate(particles, weights):
    """"
    returns mean and variance of the weighted particles"""

    pos = particles[:, 0:2]
    mean = np.average(pos, weights=weights, axis=0)
    var = np.average((pos - mean)**2, weights=weights, axis=0)
    return mean, var
```

If we create a uniform distribution of points in a 1x1 square with equal weights we get a mean position very near the center of the square at (0.5, 0.5) and a small variance.

```
In [12]: particles = create_uniform_particles((0,1), (0,1), (0, 5), 1000)
weights = np.array([.25]*1000)
estimate(particles, weights)
```

```
Out[12]: (array([ 0.494,  0.514]), array([ 0.083,  0.085]))
```

### 12.5.4 Particle Resampling

The SIS algorithm suffers from the degeneracy problem. It starts with uniformly distributed particles with equal weights. There may only be a handful of particles near the robot. As the algorithm runs any particle that does not match the measurements will acquire an extremely low weight. Only the particles which are near the robot will have an appreciable weight. We could have 5,000 particles with only 3 contributing meaningfully to the state estimate! We say the filter has degenerated.

This problem is solved by some form of resampling of the particles. Particles with very small weights do not meaningfully describe the probability distribution of the robot.

The resampling algorithm discards particles with very low probability and replaces them with new particles with higher probability. It does that by duplicating particles with relatively high probability. The duplicates are slightly dispersed by the noise added in the predict step. This results in a set of points in which a large majority of the particles accurately represent the probability distribution.

There are many resampling algorithms. For now let's look at one of the simplest, simple random resampling, also called multinomial resampling. It samples from the current particle set  $N$  times, making a new set of particles from the sample. The probability of selecting any given particle should be proportional to its weight.

We accomplish this with NumPy's `cumsum` function. `cumsum` computes the cumulative sum of an array. That is, element one is the sum of elements zero and one, element two is the sum of elements zero and one, etc. Then we generate random numbers in the range of 0.0 to 1.0 and do a binary search to find the weight that most closely matches that number:

```
In [13]: def simple_resample(particles, weights):
    N = len(particles)
    cumulative_sum = np.cumsum(weights)
    cumulative_sum[-1] = 1. # avoid round-off error
    indexes = np.searchsorted(cumulative_sum, random(N))

    # resample according to indexes
    particles[:] = particles[indexes]
    weights[:] = weights[indexes]
    weights /= np.sum(weights) # normalize
```

We don't resample at every epoch. For example, if you received no new measurements you have not received any information from which the resample can benefit. We can determine when to resample by using something called the effective N, which approximately measures the number of particles which meaningfully contribute to the probability distribution. The equation for this is

$$\hat{N}_{eff} = \frac{1}{\sum w^2}$$

and we can implement this in Python with

```
In [14]: def neff(weights):
    return 1. / np.sum(np.square(weights))
```

If  $\hat{N}_{eff}$  falls below some threshold it is time to resample. A useful starting point is  $\frac{1}{2}N$ , but this varies by problem. It is also possible for  $\hat{N}_{eff} = N$ , which means the particle set has collapsed to one point (each has equal weight). It may not be theoretically pure, but if that happens I create a new distribution of particles in the hopes of generating particles with more diversity. If this happens to you often, you may need to increase the number of particles, or otherwise adjust your filter. We will talk more of this later.

## 12.6 SIR Filter - A Complete Example

There is more to learn, but we know enough to implement a full particle filter. We will implement the Sampling Importance Resampling filter, or SIR.

I need to introduce a more sophisticated resampling method than I gave above. FilterPy provides several resampling methods. I will describe them later. They take an array of weights and returns indexes particles that have been chosen for the resampling. We just need to write a function that performs the resampling from these indexes:

```
In [15]: from filterpy.monte_carlo import systematic_resample

def resample_from_index(particles, weights, indexes):
    particles[:] = particles[indexes]
    weights[:] = weights[indexes]
    weights /= np.sum(weights)
```

To implement the filter we need to create the particles and the landmarks. We then execute a loop, successively calling `predict`, `update`, resampling, and then computing the new state estimate with `estimate`.

```
In [16]: from numpy.linalg import norm
        from numpy.random import randn
```

```

import scipy.stats

def run_pf1(N, iters=18, sensor_std_err=.1,
            do_plot=True, plot_particles=False,
            xlim=(0, 20), ylim=(0, 20),
            initial_x=None):
    landmarks = np.array([[-1, 2], [5, 10], [12, 14], [18, 21]])
    NL = len(landmarks)

    # create particles and weights
    if initial_x is not None:
        particles = create_gaussian_particles(
            mean=initial_x, std=(5, 5, np.pi/4), N=N)
    else:
        particles = create_uniform_particles((0, 20), (0, 20), (0, 6.28), N)
    weights = np.zeros(N)

    if plot_particles:
        alpha = .20
        if N > 5000:
            alpha *= np.sqrt(5000)/np.sqrt(N)
        plt.scatter(particles[:, 0], particles[:, 1],
                    alpha=alpha, color='g')

    xs = []
    robot_pos = np.array([0., 0.])
    for x in range(iters):
        robot_pos += (1, 1)

        # distance from robot to each landmark
        zs = (norm(landmarks - robot_pos, axis=1) +
              (randn(NL) * sensor_std_err))

        # move diagonally forward to (x+1, x+1)
        predict(particles, u=(0.00, 1.414), std=(.2, .05))

        # incorporate measurements
        update(particles, weights, z=zs, R=sensor_std_err,
               landmarks=landmarks)

        # resample if too few effective particles
        if neff(weights) < N/2:
            indexes = systematic_resample(weights)
            resample_from_index(particles, weights, indexes)

        mu, var = estimate(particles, weights)
        xs.append(mu)

        if plot_particles:
            plt.scatter(particles[:, 0], particles[:, 1],
                        color='k', marker=',', s=1)
            p1 = plt.scatter(robot_pos[0], robot_pos[1], marker='+',
                            color='k', s=180, lw=3)
            p2 = plt.scatter(mu[0], mu[1], marker='s', color='r')

```

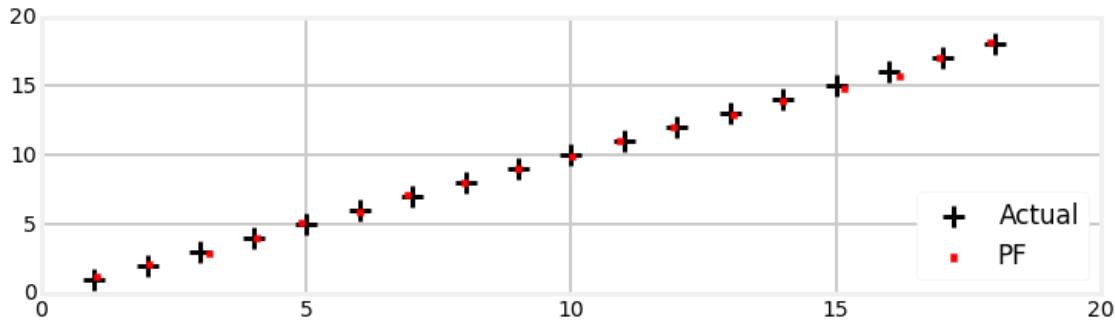
```

xs = np.array(xs)
#plt.plot(xs[:, 0], xs[:, 1])
plt.legend([p1, p2], ['Actual', 'PF'], loc=4, numpoints=1)
plt.xlim(*xlim)
plt.ylim(*ylim)
print('final position error, variance:\n\t', mu, var)

from numpy.random import seed
seed(2)
run_pf1(N=5000, plot_particles=False)

final position error, variance:
[ 17.904  18.116] [ 0.005  0.004]

```



Most of this code is devoted to initialization and plotting. The entirety of the particle filter processing consists of these lines:

```

# move diagonally forward to (x+1, x+1)
predict(particles, u=(0.00, 1.414), std=(.2, .05))

# incorporate measurements
update(particles, weights, z=zs, R=sensor_std_err,
       landmarks=landmarks)

# resample if too few effective particles
if neff(weights) < N/2:
    indexes = systematic_resample(weights)
    resample_from_index(particles, weights, indexes)

mu, var = estimate(particles, weights)

```

The first line predicts the position of the particles with the assumption that the robot is moving in a straight line ( $u[0] == 0$ ) and moving 1 unit in both the x and y axis ( $u[1]==1.414$ ). The standard deviation for the error in the turn is 0.2, and the standard deviation for the distance is 0.05. When this call returns the particles will all have been moved forward, but the weights are no longer correct as they have not been updated.

The next line incorporates the measurement into the filter. This does not alter the particle positions, it only alters the weights. If you recall the weight of the particle is computed as the probability that it matches the

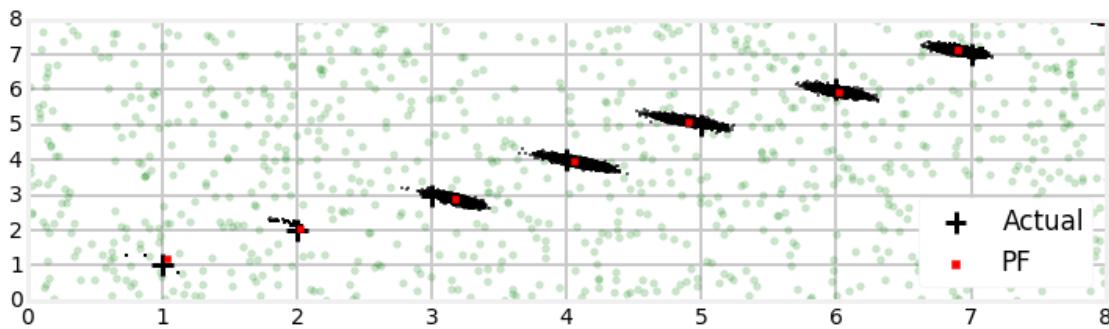
Gaussian of the sensor error model. The further the particle from the measured distance the less likely it is to be a good representation.

The final two lines example the effective particle count ( $\hat{N}_{eff}$ ). If it falls below  $\frac{1}{2}N$  we perform resampling to try to ensure our particles form a good representation of the actual probability distribution.

Now let's look at this with all the particles plotted. Seeing this happen interactively is much more useful, but this format still gives us useful information. I plotted the original random distribution of points in a very pale green and large circles to help distinguish them from the subsequent iterations where the particles are plotted with black pixels. The number of particles makes it hard to see the details, so I limited the number of iterations to 8 so we can zoom in and look more closely.

```
In [17]: seed(2)
        run_pf1(N=5000, iters=8, plot_particles=True,
                  xlim=(0,8), ylim=(0,8))

final position error, variance:
[ 7.987  7.983] [ 0.003  0.003]
```



From the plot it looks like there are only a few particles at the first two robot positions. This is not true; there are 5,000 particles, but due to resampling most are duplicates of each other. The reason for this is the Gaussian for the sensor is very narrow. This is called sample impoverishment and can lead to filter divergence. I'll address this in detail below. For now, looking at the second step at  $x=2$  we can see that the particles have dispersed a bit. This dispersion is due to the motion model noise. All particles are projected forward according to the control input  $u$ , but noise is added to each particle proportional to the error in the control mechanism in the robot. By the third step the particles have dispersed enough to make a convincing cloud of particles around the robot.

The shape of the particle cloud is an ellipse. This is not a coincidence. The sensors and robot control are both modeled as Gaussian, so the probability distribution of the system is also a Gaussian. The particle filter is a sampling of the probability distribution, so the cloud should be an ellipse.

It is important to recognize that the particle filter algorithm does not require the sensors or system to be Gaussian or linear. Because we represent the probability distribution with a cloud of particles we can handle any probability distribution and strongly nonlinear problems. There can be discontinuities and hard limits in the probability model. For example

### 12.6.1 Effect of Sensor Errors on the Filter

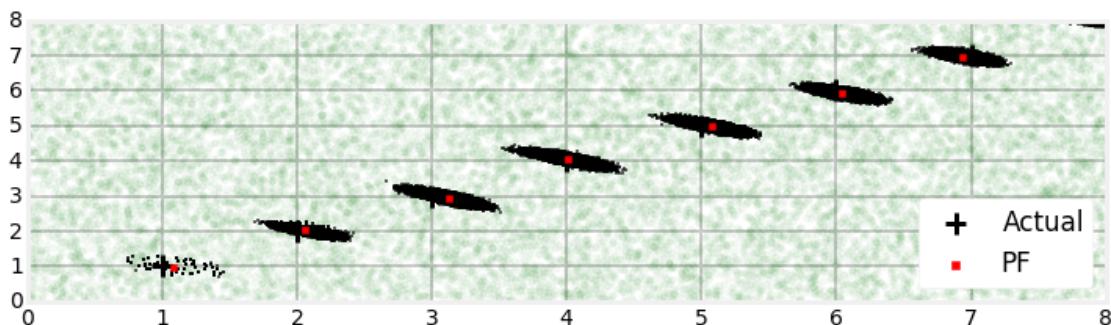
The first few iterations of the filter resulted in many duplicate particles. This happens because the model for the sensors is Gaussian, and we gave it a small standard deviation of  $\sigma = 0.1$ . This is counterintuitive at first. The Kalman filter performs better when the noise is smaller, yet the particle filter can perform worse.

We can reason why this is true. If  $\sigma = 0.1$ , the robot is at  $(1, 1)$  and a particle is at  $(2, 2)$  the particle is 14 standard deviations away from the robot. This gives it a near zero probability. It contributes nothing to the estimate of the mean, and it is extremely unlikely to survive after the resampling. If  $\sigma = 1.4$  then the particle is only  $1\sigma$  away and thus it will contribute to the estimate of the mean. During resampling it is likely to be copied one or more times.

This is very important to understand - a very accurate sensor can lead to poor performance of the filter because few of the particles will be a good sample of the probability distribution. There are a few fixes available to us. First, we can artificially increase the sensor noise standard deviation so the particle filter will accept more points as matching the robots probability distribution. This is non-optimal because some of those points will be a poor match. The real problem is that there aren't enough points being generated such that enough are near the robot. Increasing  $N$  usually fixes this problem. This decision is not cost free as increasing the number of particles significantly increase the computation time. Still, let's look at the result of using 100,000 particles.

```
In [18]: seed(2)
run_pf1(N=100000, iters=8, plot_particles=True,
        xlim=(0,8), ylim=(0,8))
```

```
final position error, variance:
[ 7.829  8.091] [ 0.003  0.003]
```



There are many more particles at  $x=1$ , and we have a convincing cloud at  $x=2$ . Clearly the filter is performing better, but at the cost of large memory usage and long run times.

Another approach is to be smarter about generating the initial particle cloud. Suppose we guess that the robot is near  $(0, 0)$ . This is not exact, as the simulation actually places the robot at  $(1, 1)$ , but it is close. If we create a normally distributed cloud near  $(0, 0)$  there is a much greater chance of the particles matching the robot's position.

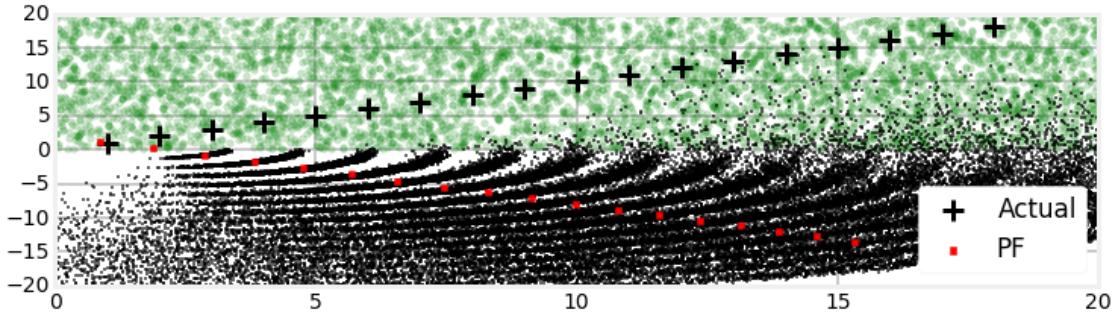
`run_pf1()` has an optional parameter `initial_x`. Use this to specify the initial position guess for the robot. The code then uses `create_gaussian_particles(mean, std, N)` to create particles distributed normally around the initial guess. We will use this in the next section.

### 12.6.2 Filter Degeneracy From Inadequate Samples

The filter as written is far from perfect. Here is how it performs with a different random seed.

```
In [19]: seed(6)
run_pf1(N=5000, plot_particles=True, ylim=(-20, 20))
```

```
final position error, variance:  
[ 15.312 -13.47 ] [ 47.065 47.03 ]
```

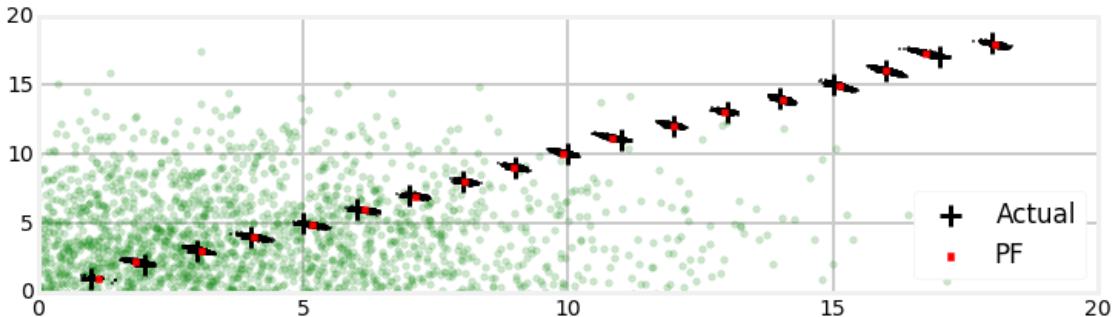


Here the initial sample of points did not generate any points near the robot. The particle filter does not create new points during the resample operation, so it ends up duplicating points which are not a representative sample of the probability distribution. As mentioned earlier this is called sample impoverishment. The problem quickly spirals out of control. The particles are not a good match for the landscape measurement so they become dispersed in a highly nonlinear, curved distribution, and the particle filter diverges from reality. No particles are available near the robot, so it cannot ever converge.

Let's make use of the `create_gaussian_particles()` method to try to generate more points near the robot. We can do this by using the `initial_x` parameter to specify a location to create the particles.

```
In [20]: seed(6)  
run_pf1(N=5000, plot_particles=True, initial_x=(1,1, np.pi/4))
```

```
final position error, variance:  
[ 18.056 17.892] [ 0.004 0.005]
```



This works great. You should always try to create particles near the initial position if you have any way to roughly estimate it. Do not be too careful - if you generate all the points very near a single position the particles may not be dispersed enough to capture the nonlinearities in the system. This is a fairly linear system, so we could get away with a smaller variance in the distribution. Clearly this depends on your problem. Increasing the number of particles is always a good way to get a better sample, but the processing cost may be a higher price than you are willing to pay.

## 12.7 Importance Sampling

I've hand waved a difficulty away which we must now confront. There is some probability distribution that describes the position and movement of our robot. We want to draw a sample of particles from that distribution and compute the integral using MC methods.

Our difficulty is that in many problems we don't know the distribution. For example, the tracked object might move very differently than we predicted with our state model. How can we draw a sample from a probability distribution that is unknown?

There is a theorem from statistics called importance sampling[1]. Remarkably it gives us a way to draw samples from a different and known probability distribution and use those to compute the properties of the unknown one. It's a fantastic theorem that brings joy to my heart.

The idea is simple, and we already used it. We draw samples from the known probability distribution, but weight the particles according to the distribution we are interested in. We can then compute properties such as the mean and variance by computing the weighted mean and weighted variance of the samples.

For the robot localization problem we drew samples from the probability distribution that we computed from our state model prediction step. In other words, we reasoned ‘the robot was there, it is perhaps moving at this direction and speed, hence it might be here’. Yet the robot might have done something completely different. It may have fell off a cliff or been hit by a mortar round. In each case the probability distribution is not correct. It seems like we are stymied, but we are not because we can use importance sampling. We drew particles from that likely incorrect probability distribution, then weighted them according to how well the particles match the measurements. That weighting is based on the true probability distribution, so according to the theory the resulting mean, variance, etc, will be correct. Magic!

How can that be true? I'll give you the math; you can safely skip this if you don't plan to go beyond the robot localization problem. However, other particle filter problems require different approaches to importance sampling, and a bit of math helps. Also, the literature and much of the content on the web uses the mathematical formulation in favor of my rather imprecise “imagine that...” exposition. If you want to understand the literature you will need to know the following equations.

We have some probability distribution  $\pi(x)$  which we want to take samples from. However, we don't know what  $\pi(x)$  is; instead we only know an alternative probability distribution  $q(x)$ . In the context of robot localization,  $\pi(x)$  is the probability distribution for the robot, but we don't know it, and  $q(x)$  is the probability distribution of our measurements, which we do know.

The expected value of a function  $f(x)$  with probability distribution  $\pi(x)$  is

$$\mathbb{E}[f(x)] = \int f(x)\pi(x) dx$$

We don't know  $\pi(x)$  so we cannot compute this integral. We do know an alternative distribution  $q(x)$  so we can add it into the integral without changing the value with

$$\mathbb{E}[f(x)] = \int f(x)\pi(x) \frac{q(x)}{q(x)} dx$$

Now we rearrange and group terms

$$\mathbb{E}[f(x)] = \int f(x)q(x) \cdot \frac{\pi(x)}{q(x)} dx$$

$q(x)$  is known to us, so we can compute  $\int f(x)q(x)$  using MC integration. That leaves us with  $\pi(x)/q(x)$ . That is a ratio, and we define it as a weight. This gives us

$$\mathbb{E}[f(x)] = \sum_{i=1}^N f(x^i)w(x^i)$$

Maybe that seems a little abstract. If we want to compute the mean of the particles we would compute

$$\mu = \sum_{i=1}^N x^i w^i$$

which is the equation I gave you earlier in the chapter.

It is required that the weights be proportional to the ratio  $\pi(x)/q(x)$ . We normally do not know the exact value, so in practice we normalize the weights by  $\sum w(x^i)$ .

When you formulate a particle filter algorithm you will have to implement this step depending on the particulars of your situation. For robot localization the best distribution to use for  $q(x)$  is the particle distribution from the `predict()` step of the filter. Let's look at the code again:

```
def update(particles, weights, z, R, landmarks):
    weights.fill(1.)
    for i, landmark in enumerate(landmarks):
        distance = np.linalg.norm(particles[:, 0:2] - landmark, axis=1)
        weights *= scipy.stats.norm(distance, R).pdf(z[i])

    weights += 1.e-300      # avoid round-off to zero
    weights /= sum(weights) # normalize
```

The reason for `self.weights.fill(1.)` might have confused you. In all the Bayesian filters up to this chapter we started with the probability distribution created by the `predict` step, and this appears to discard that information by setting all of the weights to 1. Well, we are discarding the weights, but we do not discard the particles. That is a direct result of applying importance sampling - we draw from the known distribution, but weight by the unknown distribution. In this case our known distribution is the uniform distribution - all are weighted equally.

Of course if you can compute the posterior probability distribution from the prior you should do so. If you cannot, then importance sampling gives you a way to solve this problem.

## 12.8 Resampling Methods

The resampling algorithm effects the performance of the filter. For example, suppose we resampled particles by picking particles at random. This would lead us to choosing many particles with a very low weight, and the resulting set of particles would be a terrible representation of the problem's probability distribution.

Research on the topic continues, but a handful of algorithms work well in practice across a wide variety of situations. We desire an algorithm that has several properties. It should preferentially select particles that have a higher probability. It should select a representative population of the higher probability particles to avoid sample impoverishment. It should include enough lower probability particles to give the filter a chance of detecting strongly nonlinear behavior.

FilterPy implements several of the popular algorithms. FilterPy doesn't know how your particle filter is implemented, so it doesn't make sense to have the resample algorithm generate the new samples. Instead, the algorithms create a `numpy.array` containing the indexes of the particles that are chosen. Your class or code needs to perform the resampling step. For example, I used this for the robot:

```
In [21]: def resample_from_index(particles, weights, indexes):
    particles[:] = particles[indexes]
    weights[:] = weights[indexes]
    weights /= np.sum(weights)
```

### 12.8.1 Multinomial Resampling

Multinomial resampling is the algorithm that I used while developing the robot localization example. The idea is simple. Compute the cumulative sum of the normalized weights. This gives you an array of increasing values from 0 to 1. Here is a plot which illustrates how this spaces out the weights. The colors are meaningless, they just make the divisions easier to see.

```
In [22]: from pf_internal import plot_cumsum
        print('cumulative sume is', np.cumsum([.1, .2, .1, .6]))
        plot_cumsum([.1, .2, .1, .6])
```

```
cumulative sume is [ 0.1  0.3  0.4  1. ]
```



To select a weight we generate a random number uniformly selected between 0 and 1 and use binary search to find its position inside the cumulative sum array. Large weights will be a further distance from its neighbors than low weights, so they will be more likely to be selected.

This is very easy to code using NumPy's ufunc support. `searchsorted` is NumPy's binary search algorithm. If you provide `is` with an array of search values it will return an array of answers; one answer for each search value.

```
In [23]: def multinomial_resample(weights):
    cumulative_sum = np.cumsum(weights)
    cumulative_sum[-1] = 1. # avoid round-off errors
    return np.searchsorted(cumulative_sum, random(len(weights)))
```

Here is an example:

```
In [24]: from pf_internal import plot_multinomial_resample
        plot_multinomial_resample([.1, .2, .3, .4, .2, .3, .1])
```



This is an  $O(n \log(n))$  algorithm. That is not terrible, but there are  $O(n)$  resampling algorithms with better properties with respect to the uniformity of the samples. I'm showing it because you can understand the other algorithms as variations on this one. There is a faster implementation of this multinomial resampling that uses the inverse of the CDF of the distribution. You can search on the internet if you are interested.

You may import the function from FilterPy using

```
from filterpy.monte_carlo import multinomial_resample
```

### 12.8.2 Residual Resampling

Residual resampling both improves the run time of multinomial resampling, and ensures that the sampling is uniform across the population of particles. It's fairly ingenious: the normalized weights are multiplied by  $N$ , and then the integer value of each weight is used to define how many samples of that particle will be taken. For example, if the weight of a particle is 0.0012 and  $N=3000$ , the scaled weight is 3.6, so 3 samples will be taken of that particle. This ensures that all higher weight particles are chosen at least once, and has  $O(N)$  running time.

However, this does not make  $N$  selections. To select the rest, we take the residual: the weights minus the integer part, which leaves the fractional part of the number. We then use a simpler sampling scheme such as multinomial, to select the rest of the particles based on the residual. In the example above the scaled weight was 3.6, so the residual will be 0.6 ( $3.6 - \text{int}(3.6)$ ). This residual is very large so the particle will be likely to be sampled again. This is reasonable because the larger the residual the larger the error in the round off, and thus the particle was relatively under sampled in the integer step.

```
In [25]: def residual_resample(weights):
    N = len(weights)
    indexes = np.zeros(N, 'i')

    # take int(N*w) copies of each weight
    num_copies = (N*np.asarray(weights)).astype(int)
    k = 0
    for i in range(N):
        for _ in range(num_copies[i]): # make n copies
            indexes[k] = i
            k += 1

    # use multinomial resample on the residual to fill up the rest.
    residual = w - num_copies      # get fractional part
    residual /= sum(residual)      # normalize
    cumulative_sum = np.cumsum(residual)
    cumulative_sum[-1] = 1. # ensures sum is exactly one
    indexes[k:N] = np.searchsorted(cumulative_sum, random(N-k))

    return indexes
```

You may be tempted to replace the inner for loop with a slice `indexes[k:k + num_copies[i]] = i`, but very short slices are comparatively slow, and the for loop usually runs faster.

Let's look at an example:

```
In [26]: from pf_internal import plot_residual_resample
plot_residual_resample([.1, .2, .3, .4, .2, .3, .1])
```

## residual resampling



You may import this from FilterPy using

```
from filterpy.monte_carlo import residual_resample
```

### 12.8.3 Stratified Resampling

This scheme aims to make selections relatively uniformly across the particles. It works by dividing the cumulative sum into  $N$  equal sections, and then selects one particle randomly from each section. This guarantees that each sample is between 0 and  $\frac{2}{N}$  apart.

The plot below illustrates this. The colored bars show the cumulative sum of the array, and the black lines show the  $N$  equal subdivisions. Particles, shown as black circles, are randomly placed in each subdivision.

```
In [27]: from pf_internal import plot_stratified_resample
plot_stratified_resample([.1, .2, .3, .4, .2, .3, .1])
```

## stratified resampling



The code to perform the stratification is quite straightforward.

```
In [28]: def stratified_resample(weights):
    N = len(weights)
    # make N subdivisions, choose a random position within each one
    positions = (random(N) + range(N)) / N

    indexes = np.zeros(N, 'i')
    cumulative_sum = np.cumsum(weights)
    i, j = 0, 0
    while i < N:
        if positions[i] < cumulative_sum[j]:
            indexes[i] = j
            i += 1
        else:
            j += 1
    return indexes
```

Import it from FilterPy with

```
from filterpy.monte_carlo import stratified_resample
```

### 12.8.4 Systematic Resampling

The last algorithm we will look at is systematic resampling. As with stratified resampling the space is divided into  $N$  divisions. We then choose a random offset to use for all of the divisions, ensuring that each sample is exactly  $\frac{1}{N}$  apart. It looks like this.

```
In [29]: from pf_internal import plot_systematic_resample
plot_systematic_resample([.1, .2, .3, .4, .2, .3, .1])
```



Having seen the earlier examples the code couldn't be simpler.

```
In [30]: def systematic_resample(weights):
    N = len(weights)

    # make N subdivisions, choose positions
    # with a consistent random offset
    positions = (np.arange(N) + random()) / N

    indexes = np.zeros(N, 'i')
    cumulative_sum = np.cumsum(weights)
    i, j = 0, 0
    while i < N:
        if positions[i] < cumulative_sum[j]:
            indexes[i] = j
            i += 1
        else:
            j += 1
    return indexes
```

Import from FilterPy with

```
python from filterpy.monte_carlo import systematic_resample
```

### 12.8.5 Choosing a Resampling Algorithm

Let's look at the four algorithms at once so they are easier to compare.

```
In [31]: a = [.1, .2, .3, .4, .2, .3, .1]
np.random.seed(4)
plot_multinomial_resample(a)
plot_residual_resample(a)
plot_systematic_resample(a)
plot_stratified_resample(a)
```

### multinomial resampling



### residual resampling



### systematic resampling



### stratified resampling



The performance of the multinomial resampling is quite bad. There is a very large weight that was not sampled at all. The largest weight only got one resample, yet the smallest weight was sample twice. Most tutorials on the net that I have read use multinomial resampling, and I am not sure why. Multinomial rarely used in the literature or for real problems. I recommend not using it unless you have a very good reason to do so.

The residual resampling algorithm does excellently at what it tries to do: ensure all the largest weights are resampled multiple times. It does do well at evenly distributing the samples across the particles - many reasonably large weights are not resampled at all.

Both systematic and stratified perform very well. Systematic sampling does an excellent job of ensuring we sample from all parts of the particle space while ensuring larger weights are proportionally resampled more often. Stratified resampling is not quite as uniform as systematic resampling, but it is a bit better at ensuring the higher weights get resampled more.

Plenty has been written on the theoretical performance of these algorithms, and feel free to read it. In practice I apply particle filters to problems that resist analytic efforts, and so I feel a bit dubious about the

validity of a specific analysis to these problems. In practice both the stratified and systematic algorithms perform well and very similarly across a variety of problems. I say try one, and if it works stick with it. If performance of the filter is critical try both, and perhaps see if there is literature published on your specific problem that will give you better guidance.

## 12.9 Summary

This chapter only touches the surface of what is a vast topic. My goal was not to teach you the field, but to expose you to practical Bayesian Monte Carlo techniques for filtering.

Particle filters are a type of ensemble filtering. Kalman filters represents state with a Gaussian. Measurements are applied to the Gaussian using Bayes Theorem, and the prediction is done using state-space methods. These techniques are applied to the Gaussian - the probability distribution.

In contrast, ensemble techniques represent a probability distribution using a discrete collection of points and associated probabilities. Measurements are applied to these points, not the Gaussian distribution. Likewise, the system model is applied to the points, not a Gaussian. We then compute the statistical properties of the resulting ensemble of points.

These choices have many trade-offs. The Kalman filter is very efficient, and is an optimal estimator if the assumptions of linearity and Gaussian noise are true. If the problem is nonlinear than we must linearize the problem. If the problem is multimodal (more than one object being tracked) then the Kalman filter cannot represent it. The Kalman filter requires that you know the state model. If you do not know how your system behaves the performance is poor.

In contrast, particle filters work with any arbitrary, non-analytic probability distribution. The ensemble of particles, if large enough, form an accurate approximation of the distribution. It performs wonderfully even in the presence of severe nonlinearities. Importance sampling allows us to compute probabilities even if we do not know the underlying probability distribution. Monte Carlo techniques replace the analytic integrals required by the other filters.

This power comes with a cost. The most obvious costs are the high computational and memory burdens the filter places on the computer. Less obvious is the fact that they are fickle. You have to be careful to avoid particle degeneracy and divergence. It can be very difficult to prove the correctness of your filter. If you are working with multimodal distributions you have further work to cluster the particles to determine the paths of the multiple objects. This can be very difficult when the objects are close to each other.

There are many different classes of particle filter; I only described the naive SIS algorithm, and followed that with a SIR algorithm that performs well. There are many more classes, and many examples of filters in each class. It would take a small book to describe them all.

When you read the literature on particle filters you will find that it is strewn with integrals. We perform computations on probability distributions using integrals, so using integrals gives the authors a powerful and compact notation. You must recognize that when you reduce these equations to code you will be representing the distributions with particles, and integrations are replaced with sums over the particles. If you keep in mind the core ideas in this chapter you material shouldn't be too daunting.

## 12.10 References

[1] Importance Sampling, Wikipedia. [https://en.wikipedia.org/wiki/Importance\\_sampling](https://en.wikipedia.org/wiki/Importance_sampling)



# Chapter 13

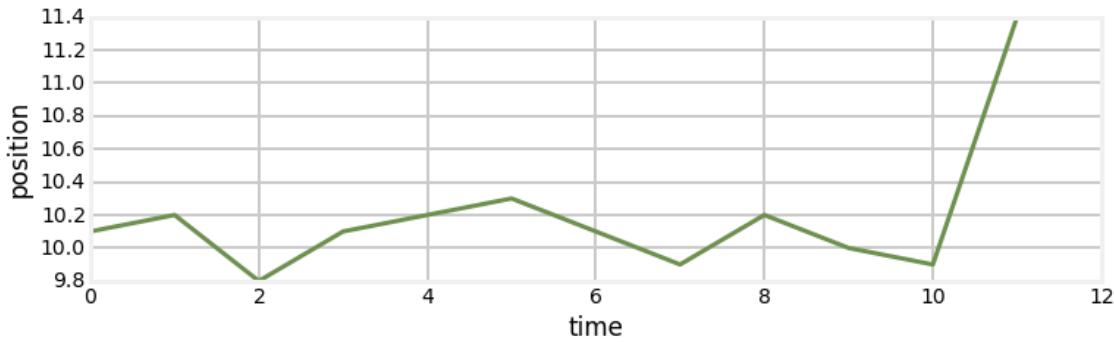
## Smoothing

### 13.1 Introduction

The performance of the Kalman filter is not optimal when you consider future data. For example, suppose we are tracking an aircraft, and the latest measurement deviates far from the current track, like so (I'll only consider 1 dimension for simplicity):

```
In [2]: import matplotlib.pyplot as plt
```

```
data = [10.1, 10.2, 9.8, 10.1, 10.2, 10.3,
        10.1, 9.9, 10.2, 10.0, 9.9, 11.4]
plt.plot(data)
plt.xlabel('time')
plt.ylabel('position');
```



After a period of near steady state, we have a very large change. Assume the change is past the limit of the aircraft's flight envelope. Nonetheless the Kalman filter incorporates that new measurement into the filter based on the current Kalman gain. It cannot reject the noise because the measurement could reflect the initiation of a turn. Granted it is unlikely that we are turning so abruptly, but it is impossible to say whether

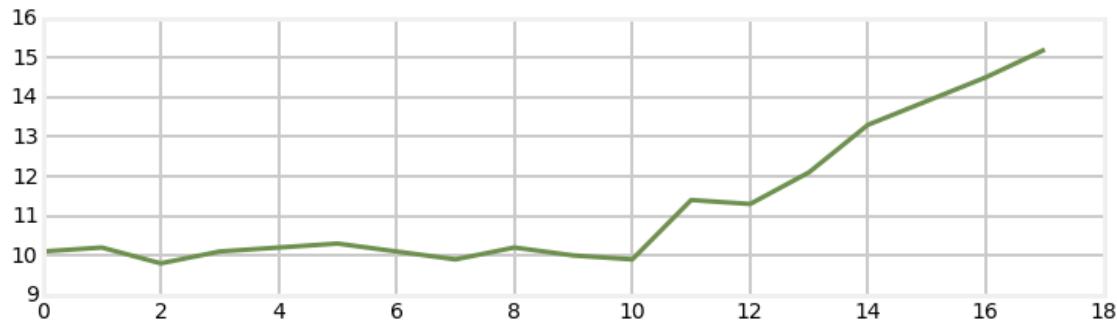
- The aircraft started a turn awhile ago, but the previous measurements were noisy and didn't show the change.

- The aircraft is turning, and this measurement is very noisy
- The measurement is very noisy and the aircraft has not turned
- The aircraft is turning in the opposite direction, and the measurement is extremely noisy

Now, suppose the following measurements are:

11.3 12.1 13.3 13.9 14.5 15.2

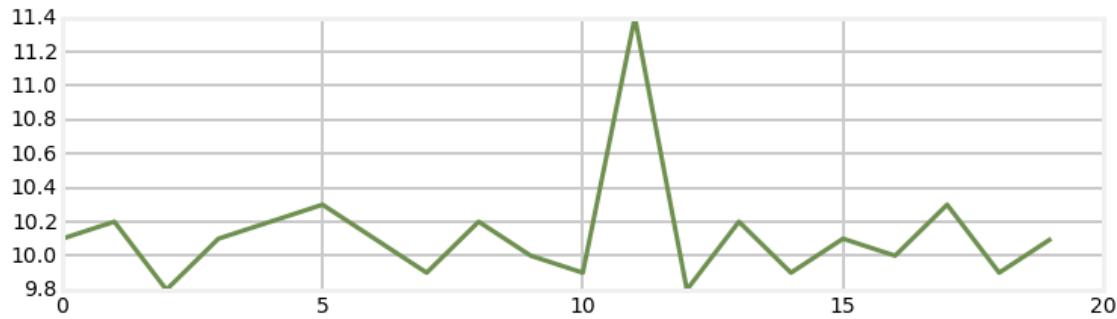
```
In [3]: data2 = [11.3, 12.1, 13.3, 13.9, 14.5, 15.2]
plt.plot(data + data2);
```



Given these future measurements we can infer that yes, the aircraft initiated a turn.

On the other hand, suppose these are the following measurements.

```
In [4]: data3 = [9.8, 10.2, 9.9, 10.1, 10.0, 10.3, 9.9, 10.1]
plt.plot(data + data3);
```



In this case we are led to conclude that the aircraft did not turn and that the outlying measurement was merely very noisy.

## 13.2 An Overview of How Smoothers Work

The Kalman filter is a recursive filter - it's estimate at step  $k$  is partially based on its estimate from step  $k-1$ . But this means that the estimate from step  $k-1$  is partially based on step  $k-2$ , and so on back to the

first measurement. Hence, the estimate at step  $k$  depends on all of the previous measurements, though to varying degrees.  $k-1$  has the most influence,  $k-2$  has the next most, and so on.

Smoothing filters incorporate future measurements into the estimate for step  $k$ . The measurement from  $k+1$  will have the most effect,  $k+2$  will have less effect,  $k+3$  less yet, and so on.

This topic is called smoothing, but I think that is a misleading name. I could smooth the data above by passing it through a low pass filter. The result would be smooth, but not necessarily accurate because a low pass filter will remove real variations just as much as it removes noise. In contrast, Kalman smoothers are optimal - they incorporate all available information to make the best estimate that is mathematically achievable.

### 13.3 Types of Smoothers

There are three classes of Kalman smoothers that produce better tracking in these situations.

- Fixed-Interval Smoothing

This is a batch processing based filter. This filter waits for all of the data to be collected before making any estimates. For example, you may be a scientist collecting data for an experiment, and don't need to know the result until the experiment is complete. A fixed-interval smoother will collect all the data, then estimate the state at each measurement using all available previous and future measurements. If it is possible for you to run your Kalman filter in batch mode it is always recommended to use one of these filters as it will provide much better results than the recursive forms of the filter from the previous chapters.

- Fixed-Lag Smoothing

Fixed-lag smoothers introduce latency into the output. Suppose we choose a lag of 4 steps. The filter will ingest the first 3 measurements but not output a filtered result. Then, when the 4th measurement comes in the filter will produce the output for measurement 1, taking measurements 1 through 4 into account. When the 5th measurement comes in, the filter will produce the result for measurement 2, taking measurements 2 through 5 into account. This is useful when you need recent data but can afford a bit of lag. For example, perhaps you are using machine vision to monitor a manufacturing process. If you can afford a few seconds delay in the estimate a fixed-lag smoother will allow you to produce very accurate and smooth results.

- Fixed-Point Smoothing

A fixed-point filter operates as a normal Kalman filter, but also produces an estimate for the state at some fixed time  $j$ . Before the time  $k$  reaches  $j$  the filter operates as a normal filter. Once  $k > j$  the filter estimates  $x_k$  and then also updates its estimate for  $x_j$  using all of the measurements between  $j \dots k$ . This can be useful to estimate initial parameters for a system, or for producing the best estimate for an event that happened at a specific time. For example, you may have a robot that took a photograph at time  $j$ . You can use a fixed-point smoother to get the best possible pose information for the camera at time  $j$  as the robot continues moving.

The choice of these filters depends on your needs and how much memory and processing time you can spare. Fixed-point smoothing requires storage of all measurements, and is very costly to compute because the output is for every time step is recomputed for every measurement. On the other hand, the filter does produce a decent output for the current measurement, so this filter can be used for real time applications.

Fixed-lag smoothing only requires you to store a window of data, and processing requirements are modest because only that window is processed for each new measurement. The drawback is that the filter's output always lags the input, and the smoothing is not as pronounced as is possible with fixed-interval smoothing.

Fixed-interval smoothing produces the most smoothed output at the cost of having to be batch processed. Most algorithms use some sort of forwards/backwards algorithm that is only twice as slow as a recursive Kalman filter.

## 13.4 Fixed-Point Smoothing

not done

## 13.5 Fixed-Interval Smoothing

There are many fixed-lag smoothers available in the literature. I have chosen to implement the smoother invented by Rauch, Tung, and Striebel because of its ease of implementation and efficiency of computation. It is also the smoother I have seen used most often in real applications. This smoother is commonly known as an RTS smoother.

Derivation of the RTS smoother runs to several pages of densely packed math, and to be honest I have never read it through. I'm certainly not going to inflict it on you. I don't think anyone but thesis writers really need to understand the derivation. Instead I will briefly present the algorithm, equations, and then move directly to implementation and demonstration of the smoother.

The RTS smoother works by first running the Kalman filter in a batch mode, computing the filter output for each step. Given the filter output for each measurement along with the covariance matrix corresponding to each output the RTS runs over the data backwards, incorporating it's knowledge of the future into the past measurements. When it reaches the first measurement it is done, and the filtered output incorporates all of the information in a maximally optimal form.

The equations for the RTS smoother are very straightforward and easy to implement. This derivation is for the linear Kalman filter. Similar derivations exist for the EKF and UKF. These steps are performed on the output of the batch processing, going backwards from the most recent in time back to the first estimate. Each iteration incorporates the knowledge of the future into the state estimate. Since the state estimate already incorporates all of the past measurements the result will be that each estimate will contain knowledge of all measurements in the past and future. Here is it very important to distinguish between past, present, and future so I have used subscripts to denote whether the data is from the future or not.

Predict Step

$$\mathbf{P} = \mathbf{F}\mathbf{P}_k\mathbf{F}^\top + \mathbf{Q}$$

Update Step

$$\begin{aligned}\mathbf{K}_k &= \mathbf{P}_k \mathbf{F} \mathbf{P}^{-1} \\ \mathbf{x}_k &= \mathbf{x}_k + \mathbf{K}_k (\mathbf{x}_{x+1} - \mathbf{F}\mathbf{x}_k) \\ \mathbf{P}_k &= \mathbf{P}_k + \mathbf{K}_k (\mathbf{P}_{K+1} - \mathbf{P}) \mathbf{K}_k^\top\end{aligned}$$

As always, the hardest part of the implementation is correctly accounting for the subscripts. A basic implementation without comments or error checking would be:

```
def rts_smoother(Xs, Ps, F, Q):
    n, dim_x, _ = Xs.shape
```

```

# smoother gain
K = zeros((n,dim_x, dim_x))
x, P = Xs.copy(), Ps.copy()

for k in range(n-2,-1,-1):
    P_pred = dot(F, P[k]).dot(F.T) + Q

    K[k] = dot(P[k], F.T).dot(inv(P_pred))
    x[k] += dot(K[k], x[k+1] - dot(F, x[k]))
    P[k] += dot(K[k], P[k+1] - P_pred).dot(K[k].T)

return (x, P, K)

```

This implementation mirrors the implementation provided in FilterPy. It assumes that the Kalman filter is being run externally in batch mode, and the results of the state and covariances are passed in via the `Xs` and `Ps` variable.

Here is an example.

```

In [5]: import numpy as np
        from numpy import random
        from numpy.random import randn
        import matplotlib.pyplot as plt
        from filterpy.kalman import KalmanFilter
        import book_plots as bp

def plot_rts(noise, Q=0.001, show_velocity=False):
    random.seed(123)
    fk = KalmanFilter(dim_x=2, dim_z=1)

    fk.x = np.array([0., 1.])      # state (x and dx)

    fk.F = np.array([[1., 1.],
                    [0., 1.]])    # state transition matrix

    fk.H = np.array([[1., 0.]])    # Measurement function
    fk.P = 10.                     # covariance matrix
    fk.R = noise                   # state uncertainty
    fk.Q = Q                       # process uncertainty

    # create noisy data
    zs = np.asarray([t + randn()*noise for t in range(40)])

    # filter data with Kalman filter, than run smoother on it
    mu, cov, _, _ = fk.batch_filter(zs)
    M,P,C = fk.rts_smoothen(mu, cov)

    # plot data
    if show_velocity:
        index = 1
        print('gu')
    else:
        index = 0
    if not show_velocity:
        bp.plot_measurements(zs, lw=1)
        plt.plot(M[:, index], c='b', label='RTS')

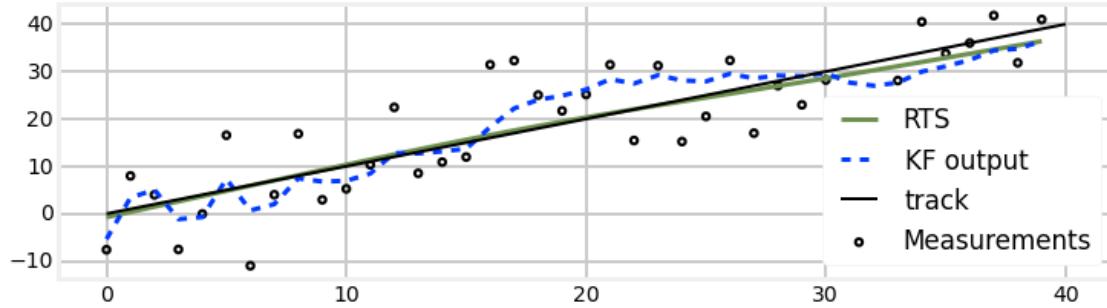
```

```

plt.plot(mu[:, index], c='g', ls='--', label='KF output')
if not show_velocity:
    N = len(zs)
    plt.plot([0, N], [0, N], 'k', lw=2, label='track')
plt.legend(loc=4)
plt.show()

plot_rts(7.)

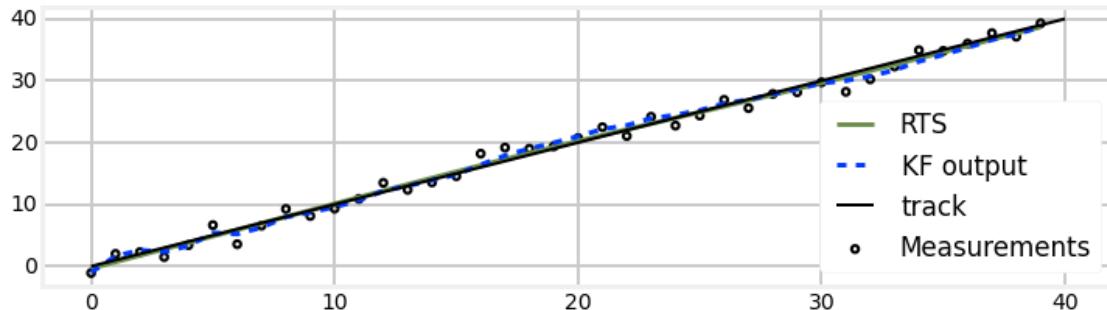
```



I've injected a lot of noise into the signal to allow you to visually distinguish the RTS output from the ideal output. In the graph above we can see that the Kalman filter, drawn as the green dotted line, is reasonably smooth compared to the input, but it still wanders from the ideal line when several measurements in a row are biased towards one side of the line. In contrast, the RTS output is both extremely smooth and very close to the ideal output.

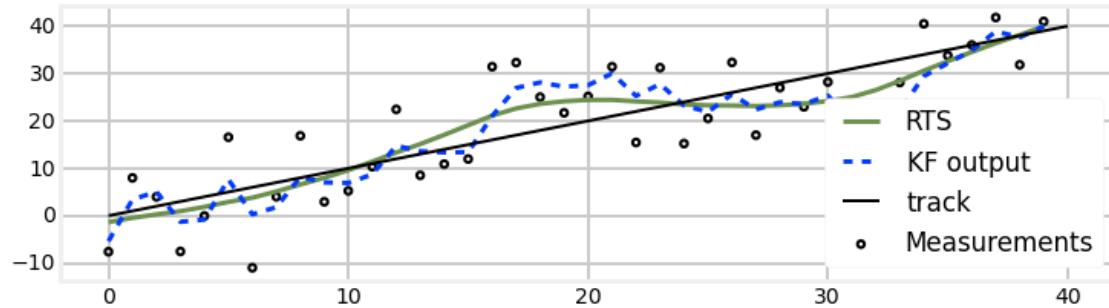
With a perhaps more reasonable amount of noise we can see that the RTS output nearly lies on the ideal output. The Kalman filter output, while much better, still varies by a far greater amount.

In [6]: `plot_rts(noise=1.)`



However, we must understand that this smoothing is predicated on the system model. We have told the filter that what we are tracking follows a constant velocity model with very low process error. When the filter looks ahead it sees that the future behavior closely matches a constant velocity so it is able to reject most of the noise in the signal. Suppose instead our system has a lot of process noise. For example, if we are tracking a light aircraft in gusty winds its velocity will change often, and the filter will be less able to distinguish between noise and erratic movement due to the wind. We can see this in the next graph.

In [7]: `plot_rts(noise=7., Q=.1)`

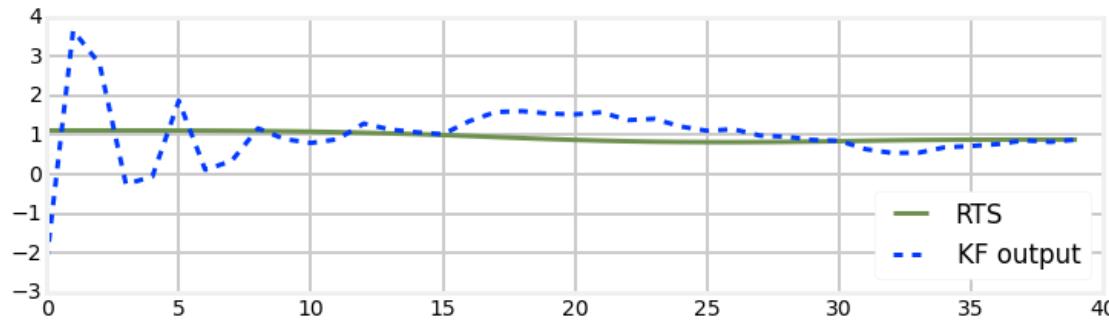


This underscores the fact that these filters are not smoothing the data in colloquial sense of the term. The filter is making an optimal estimate based on previous measurements, future measurements, and what you tell it about the behavior of the system and the noise in the system and measurements.

Let's wrap this up by looking at the velocity estimates of Kalman filter vs the RTS smoother.

In [8]: `plot_rts(7., show_velocity=True)`

gu



The improvement in the velocity, which is an hidden variable, is even more dramatic.

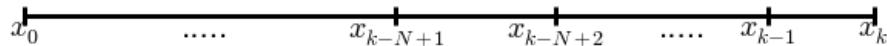
Since we are plotting velocities let's look at what the the ‘raw’ velocity is, which we can compute by subtracting subsequent measurements. We see below that the noise swamps the signal, causing the computed values to be essentially worthless. I show this to reiterate the importance of using Kalman filters to compute velocities, accelerations, and even higher order values. I use a Kalman filter even when my measurements are so accurate that I am willing to use them unfiltered if I am also interested in the velocities and/or accelerations. Even a very small error in a measurement gets magnified when computing velocities.

## 13.6 Fixed-Lag Smoothing

The RTS smoother presented above should always be your choice of algorithm if you can run in batch mode because it incorporates all available data into each estimate. Not all problems allow you to do that, but you

may still be interested in receiving smoothed values for previous estimates. The number line below illustrates this concept.

```
In [9]: from book_format import figsize
from smoothing_internal import *
with figsize(y=2):
    show_fixed_lag_numberline()
```



At step  $k$  we can estimate  $x_k$  using the normal Kalman filter equations. However, we can make a better estimate for  $x_{k-1}$  by using the measurement received for  $x_k$ . Likewise, we can make a better estimate for  $x_{k-2}$  by using the measurements received for  $x_{k-1}$  and  $x_k$ . We can extend this computation back for an arbitrary  $N$  steps.

Derivation for this math is beyond the scope of this book; Dan Simon's [Optimal State Estimation](#) [2] has a very good exposition if you are interested. The essence of the idea is that instead of having a state vector  $\mathbf{x}$  we make an augmented state containing

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_k \\ \mathbf{x}_{k-1} \\ \vdots \\ \mathbf{x}_{k-N+1} \end{bmatrix}$$

This yields a very large covariance matrix that contains the covariance between states at different steps. FilterPy's class `FixedLagSmoother` takes care of all of this computation for you, including creation of the augmented matrices. All you need to do is compose it as if you are using the `KalmanFilter` class and then call `smooth()`, which implements the predict and update steps of the algorithm.

Each call of `smooth` computes the estimate for the current measurement, but it also goes back and adjusts the previous  $N-1$  points as well. The smoothed values are contained in the list `FixedLagSmoother.xSmooth`. If you use `FixedLagSmoother.x` you will get the most recent estimate, but it is not smoothed and is no different from a standard Kalman filter output.

```
In [10]: from filterpy.kalman import FixedLagSmoother, KalmanFilter
import numpy.random as random

fls = FixedLagSmoother(dim_x=2, dim_z=1, N=8)

fls.x = np.array([0., .5])
fls.F = np.array([[1., 1.],
                 [0., 1.]))

fls.H = np.array([[1., 0.]])
```

```

fls.P *= 200
fls.R *= 5.
fls.Q *= 0.001

kf = KalmanFilter(dim_x=2, dim_z=1)
kf.x = np.array([0., .5])
kf.F = np.array([[1., 1.],
                 [0., 1.]])
kf.H = np.array([[1., 0.]])
kf.P *= 200
kf.R *= 5.
kf.Q *= 0.001

N = 4 # size of lag

nom = np.array([t/2. for t in range (0, 40)])
zs = np.array([t + random.randn()*5.1 for t in nom])

for z in zs:
    fls.smooth(z)

kf_x, _, _, _ = kf.batch_filter(zs)
x_smooth = np.array(fls.xSmooth)[:, 0]

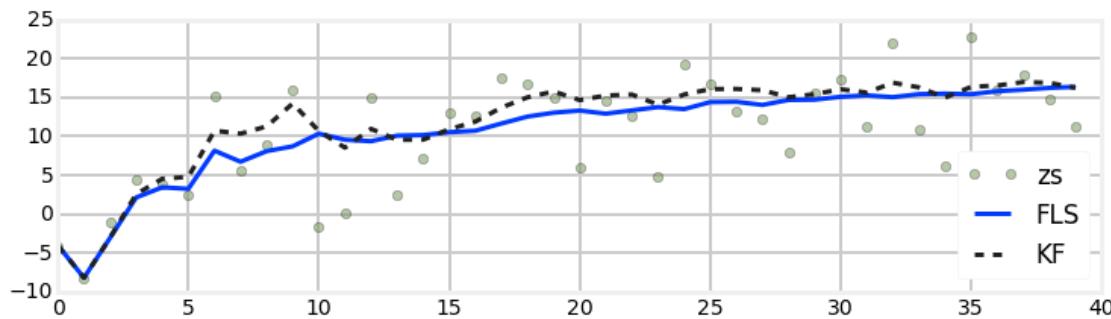
fls_res = abs(x_smooth - nom)
kf_res = abs(kf_x[:, 0] - nom)

plt.plot(zs, 'o', alpha=0.5, marker='o', label='zs')
plt.plot(x_smooth, label='FLS')
plt.plot(kf_x[:, 0], label='KF', ls='--')
plt.legend(loc=4)

print('standard deviation fixed-lag:', np.mean(fls_res))
print('standard deviation kalman:', np.mean(kf_res))

standard deviation fixed-lag: 2.61624417392
standard deviation kalman: 3.56221928461

```



Here I have set  $N=8$  which means that we will incorporate 8 future measurements into our estimates. This provides us with a very smooth estimate once the filter converges, at the cost of roughly 8x the amount of

computation of the standard Kalman filter. Feel free to experiment with larger and smaller values of  $N$ . I chose 8 somewhat at random, not due to any theoretical concerns.

### 13.7 References

- [1] H. Rauch, F. Tung, and C. Striebel. “Maximum likelihood estimates of linear dynamic systems,” AIAA Journal, **3**(8), pp. 1445-1450 (August 1965).
- [2] Dan Simon. “Optimal State Estimation,” John Wiley & Sons, 2006.  
<http://arc.aiaa.org/doi/abs/10.2514/3.3166>

# Chapter 14

## Adaptive Filtering

### 14.1 Introduction

So far we have considered the problem of tracking objects that are well behaved in relation to our process model. For example, we can use a constant velocity filter to track an object moving in a straight line. So long as the object moves in a straight line at a reasonably constant speed, or varies its track and/or velocity very slowly this filter will perform very well. Suppose instead that we are trying to track a maneuvering target, by which I mean an object with control inputs, such as a car along a road, an aircraft in flight, and so on. In these situations the filters perform quite poorly. Alternatively, consider a situation such as tracking a sailboat in the ocean. Even if we model the control inputs we have no way to model the wind or the ocean currents.

A first order approach to this problem is to make the process noise  $\mathbf{Q}$  larger to account for the unpredictability of the system dynamics. While this can work in the sense of providing a non-diverging filter, the result is typically far from optimal. The larger  $\mathbf{Q}$  results in the filter giving more emphasis to the noise in the measurements. We will see an example of this shortly.

In this chapter we will discuss the concept of an adaptive filter. This means about what it sounds like. The filter will adapt itself when it detects dynamics that the process model cannot account for. I will start with an example of the problem, and then discuss and implement various adaptive filters.

### 14.2 Maneuvering Targets

We need a simulation of a maneuvering target. I will implement a simple 2D model with steering inputs. You provide a new speed and/or direction, and it will modify its state to match.

```
In [2]: from math import sin, cos, radians

def angle_between(x, y):
    return min(y-x, y-x+360, y-x-360, key=abs)

class ManeuveringTarget(object):
    def __init__(self, x0, y0, v0, heading):
        self.x = x0
        self.y = y0
        self.vel = v0
        self.hdg = heading
```

```

        self.cmd_vel = v0
        self.cmd_hdg = heading
        self.vel_step = 0
        self.hdg_step = 0
        self.vel_delta = 0
        self.hdg_delta = 0

    def update(self):
        vx = self.vel * cos(radians(90-self.hdg))
        vy = self.vel * sin(radians(90-self.hdg))
        self.x += vx
        self.y += vy

        if self.hdg_step > 0:
            self.hdg_step -= 1
            self.hdg += self.hdg_delta

        if self.vel_step > 0:
            self.vel_step -= 1
            self.vel += self.vel_delta
        return (self.x, self.y)

    def set_commanded_heading(self, hdg_degrees, steps):
        self.cmd_hdg = hdg_degrees
        self.hdg_delta = angle_between(self.cmd_hdg,
                                       self.hdg) / steps
        if abs(self.hdg_delta) > 0:
            self.hdg_step = steps
        else:
            self.hdg_step = 0

    def set_commanded_speed(self, speed, steps):
        self.cmd_vel = speed
        self.vel_delta = (self.cmd_vel - self.vel) / steps
        if abs(self.vel_delta) > 0:
            self.vel_step = steps
        else:
            self.vel_step = 0

```

Now let's implement a simulated sensor with noise.

```
In [3]: from numpy.random import randn

class NoisySensor(object):
    def __init__(self, std_noise=1.):
        self.std = std_noise

    def sense(self, pos):
        """Pass in actual position as tuple (x, y).
        Returns position with noise added (x,y)"""

```

```

    return (pos[0] + (randn() * self.std),
           pos[1] + (randn() * self.std))
```

Now let's generate a track and plot it to test that everything is working. I'll put the data generation in a function so we can create paths of different lengths (why will be clear soon).

```

In [4]: import book_plots as bp
         import numpy as np
         import matplotlib.pyplot as plt

def generate_data(steady_count, std):
    t = ManeuveringTarget(x0=0, y0=0, v0=0.3, heading=0)
    xs, ys = [], []

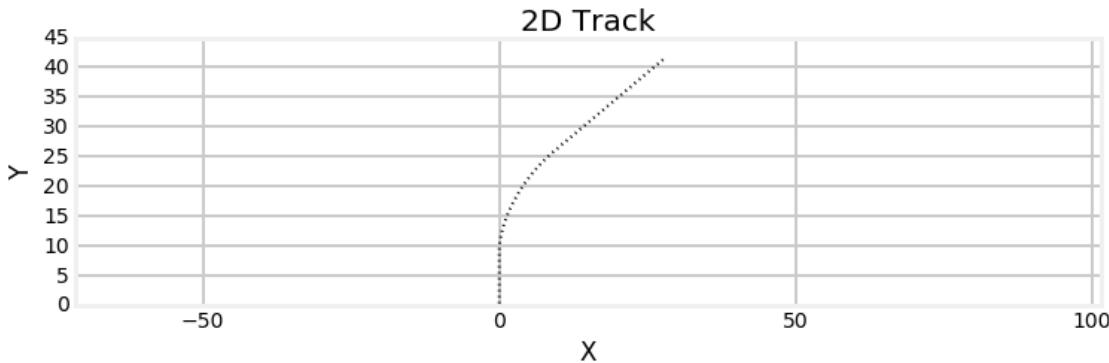
    for i in range(30):
        x, y = t.update()
        xs.append(x)
        ys.append(y)

    t.set_commanded_heading(310, 25)
    t.set_commanded_speed(1, 15)

    for i in range(steady_count):
        x, y = t.update()
        xs.append(x)
        ys.append(y)

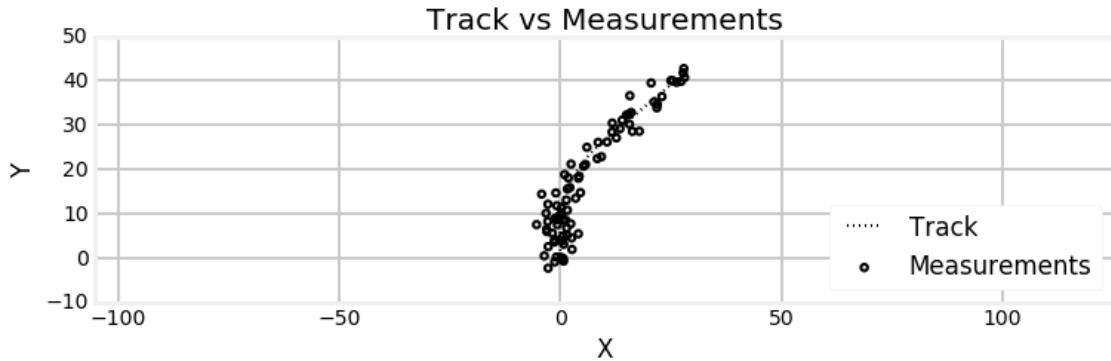
    ns = NoisySensor(std)
    pos = np.array(list(zip(xs, ys)))
    zs = np.array([ns.sense(p) for p in pos])
    return pos, zs

sensor_std = 2.
track, zs = generate_data(50, sensor_std)
bp.plot_track(*zip(*track))
plt.axis('equal')
bp.set_labels(title='2D Track', x='X', y='Y')
```



And here is the track vs the simulated sensor readings.

```
In [5]: bp.plot_track(*zip(*track))
bp.plot_measurements(*zip(*zs))
plt.axis('equal')
plt.legend(loc=4)
bp.set_labels(title='Track vs Measurements', x='X', y='Y')
```



This large amount of noise allows us to see the effect of various design choices more easily.

Now we can implement a Kalman filter to track this object. But let's make a simplification first. The  $x$  and  $y$  coordinates are independent, so we can track each independently. In the remainder of this chapter we will only track the  $x$  coordinate to keep the code and matrices as small as possible.

We start with a constant velocity filter.

```
In [6]: from filterpy.kalman import KalmanFilter
from filterpy.common import Q_discrete_white_noise

def make_cv_filter(dt, std):
    cvfilter = KalmanFilter(dim_x = 2, dim_z=1)
    cvfilter.x = np.array([0., 0.])
    cvfilter.P *= 3
    cvfilter.R *= std**2
    cvfilter.F = np.array([[1, dt],
                          [0, 1]], dtype=float)
    cvfilter.H = np.array([[1, 0]], dtype=float)
    cvfilter.Q = Q_discrete_white_noise(dim=2, dt=dt, var=0.02)

    return cvfilter

def initialize_filter(kf, std_R=None):
    """
    helper function - we will be reinitializing the filter
    many times.
    """

    kf.x.fill(0.)
    kf.P = np.eye(kf.dim_x) * .1
    if std_R is not None:
        kf.R = np.eye(kf.dim_z) * std_R
```

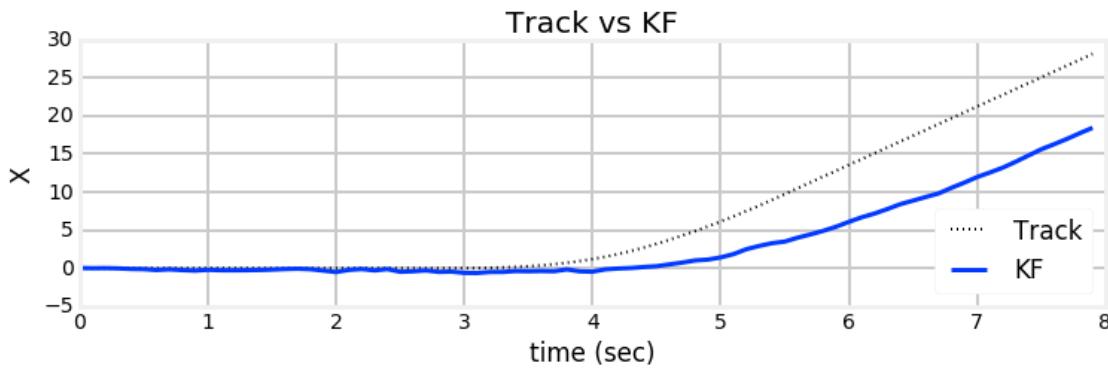
Now we run it:

```
In [7]: sensor_std = 2.
dt = 0.1

# initialize filter
cvfilter = make_cv_filter(dt, sensor_std)
initialize_filter(cvfilter)

# run it
z_xs = zs[:, 0]
kxs, _, _, _ = cvfilter.batch_filter(z_xs)

# plot results
t = np.arange(0, len(z_xs) * dt, dt)
bp.plot_track(t, track[:, 0])
bp.plot_filter(t, kxs[:, 0], label='KF')
bp.set_labels(title='Track vs KF', x='time (sec)', y='X');
plt.legend(loc=4);
```



We can see from the plot that the Kalman filter was unable to track the change in heading. Recall from the g-h Filter chapter that this is because the filter is not modeling acceleration, hence it will always lag the input. The filter will eventually catch up with the signal if the signal enters a steady state. Let's look at that.

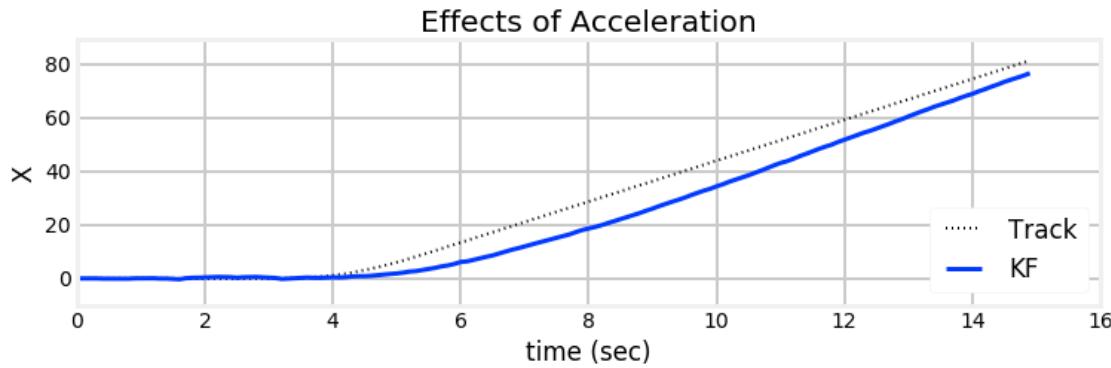
```
In [8]: # reinitialize filter
initialize_filter(cvfilter)

track2, zs2 = generate_data(120, sensor_std)
xs2 = track2[:, 0]
z_xs2 = zs2[:, 0]
t = np.arange(0, len(xs2) * dt, dt)

kxs2, _, _, _ = cvfilter.batch_filter(z_xs2)

bp.plot_track(t, xs2)
bp.plot_filter(t, kxs2[:, 0], label='KF')
plt.legend(loc=4)
```

```
bp.set_labels(title='Effects of Acceleration',
             x='time (sec)', y='X')
```

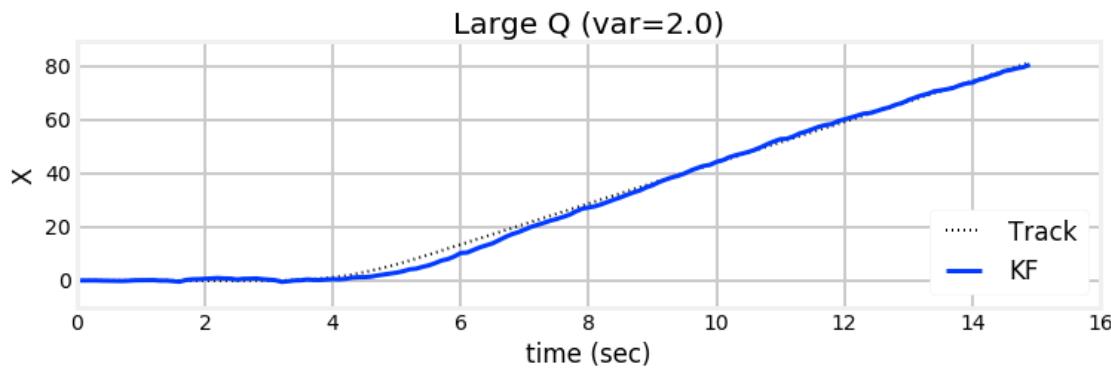


The underlying problem is that our process model is correct for the steady state sections, but incorrect for when the object is maneuvering. We can try to account for this by increasing the size of Q, like so.

```
In [9]: # reinitialize filter
    initialize_filter(cvfilter)
    cvfilter.Q = Q_discrete_white_noise(dim=2, dt=dt, var=2.0)

    # recompute track
    kxs2, _, _, _ = cvfilter.batch_filter(z_xs2)

    bp.plot_track(t, xs2)
    bp.plot_filter(t, kxs2[:, 0], label='KF')
    plt.legend(loc=4)
    bp.set_labels(title='Large Q (var=2.0)', x='time (sec)', y='X')
```

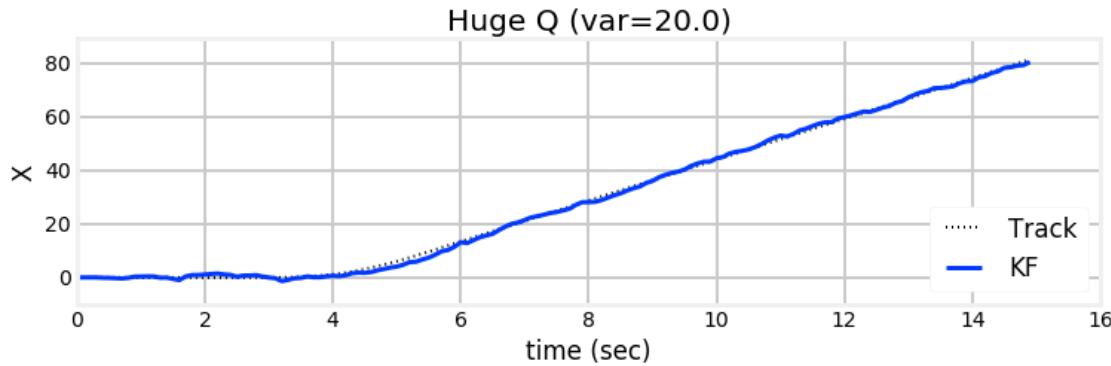


We can see that the filter reacquired the track more quickly, but at the cost of a lot of noise in the output. Furthermore, many tracking situations could not tolerate the amount of lag shown between seconds 4 and 8. We could reduce it further at the cost of very noisy output, like so:

```
In [10]: # reinitialize filter
    initialize_filter(cvfilter)
    cvfilter.Q = Q_discrete_white_noise(dim=2, dt=dt, var=20.0)

    # recompute track
    cvfilter.x.fill(0.)
    kxs2, _, _, _ = cvfilter.batch_filter(z_xs2)

    bp.plot_track(t, xs2)
    bp.plot_filter(t, kxs2[:, 0], label='KF')
    plt.legend(loc=4)
    bp.set_labels(title='Huge Q (var=20.0)', x='time (sec)', y='X')
```



Maneuvers imply acceleration, so let's implement a constant acceleration Kalman filter and see how it fairs with the same data.

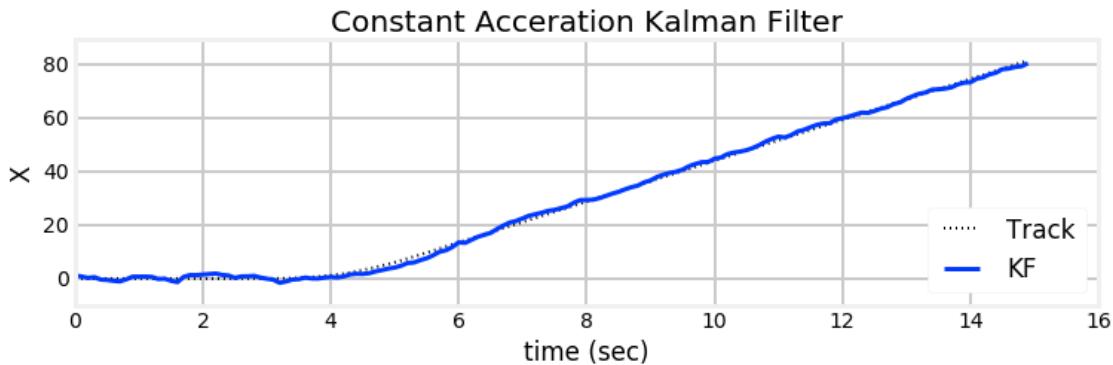
```
In [11]: def make_ca_filter(dt, std):
    cafILTER = KalmanFilter(dim_x=3, dim_z=1)
    cafILTER.x = np.array([0., 0., 0.])
    cafILTER.P *= 3
    cafILTER.R *= std
    cafILTER.Q = Q_discrete_white_noise(dim=3, dt=dt, var=0.02)
    cafILTER.F = np.array([[1, dt, 0.5*dt*dt],
                           [0, 1, dt],
                           [0, 0, 1]])
    cafILTER.H = np.array([[1., 0, 0]])
    return cafILTER

def initialize_const_accel(f):
    f.x = np.array([0., 0., 0.])
    f.P = np.eye(3) * 3
```

```
In [12]: cafILTER = make_ca_filter(dt, sensor_std)
initialize_const_accel(cafILTER)

kxs2, _, _, _ = cafILTER.batch_filter(z_xs2)
bp.plot_track(t, xs2)
bp.plot_filter(t, kxs2[:, 0], label='KF')
```

```
plt.legend(loc=4)
bp.set_labels(title='Constant Acceration Kalman Filter',
             x='time (sec)', y='X')
```



The constant acceleration model is able to track the maneuver with no lag, but at the cost of very noisy output during the steady state behavior. The noisy output is due to the filter being unable to distinguish between the beginning of an maneuver and noise in the signal. Noise in the signal implies an acceleration, and so the acceleration term of the filter tracks it.

It seems we cannot win. A constant velocity filter cannot react quickly when the target is accelerating, but a constant acceleration filter misinterprets noise during zero acceleration regimes as legitimate acceleration.

Yet there is an important insight here that will lead us to a solution. When the target is not maneuvering (the acceleration is zero) the constant velocity filter performs optimally. When the target is maneuvering the constant acceleration filter performs well, as does the constant velocity filter with an artificially large process noise  $\mathbf{Q}$ . If we make a filter that adapts itself to the behavior of the tracked object we could have the best of both worlds.

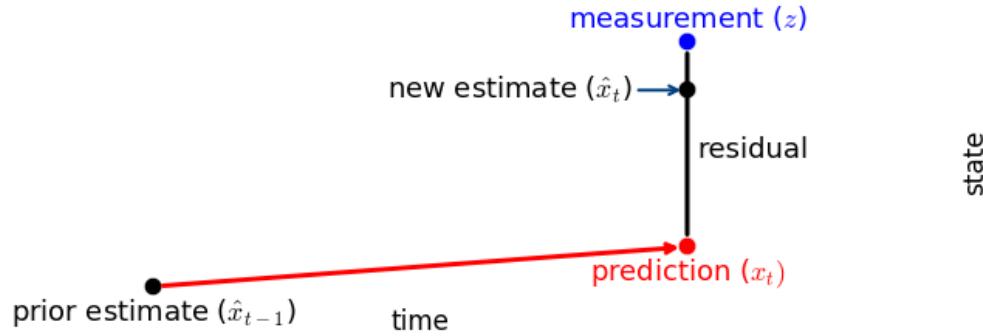
### 14.3 Detecting a Maneuver

Before we discuss how to create an adaptive filter we have to ask “how do we detect a maneuver?” We cannot reasonably adapt a filter to respond to maneuvers if we do not know when a maneuver is happening.

We have been defining maneuver as the time when the tracked object is accelerating, but in general we can say that the object is maneuvering with respect to the Kalman filter if its behavior is different than the process model being used by the filter.

What is the mathematical consequence of a maneuvering object for the filter? The object will be behaving differently than predicted by the filter, so the residual will be large. Recall that the residual is the difference between the current prediction of the filter and the measurement.

```
In [13]: from mkf_internal import show_residual_chart
show_residual_chart()
```



To confirm this, let's plot the residual for the filter during the maneuver. I will reduce the amount of noise in the data to make it easier to see the residual.

```
In [14]: from adaptive_internal import plot_track_and_residuals

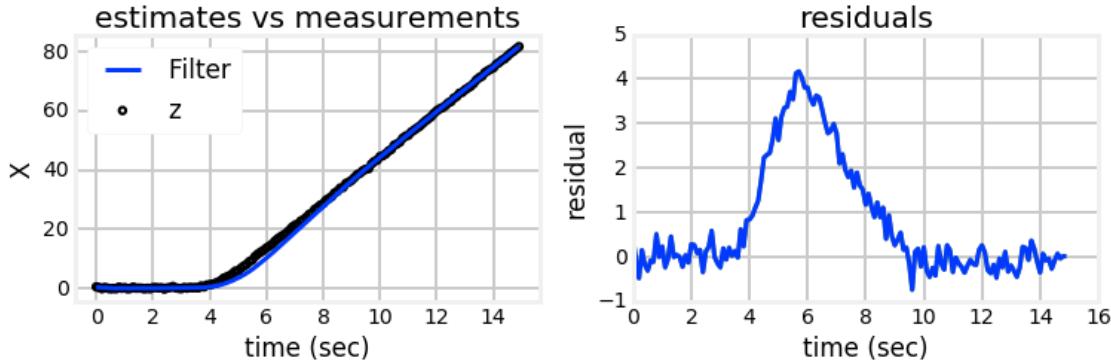
def show_residual_chart():
    dt = 0.1
    sensor_std = 0.2

    # initialize filter
    cvfilter= make_cv_filter(dt, sensor_std)
    initialize_filter(cvfilter)
    pos2, zs2 = generate_data(120, sensor_std)
    xs2 = pos2[:, 0]
    z_xs2 = zs2[:, 0]

    cvfilter.Q = Q_discrete_white_noise(dim=2, dt=dt, var=0.02)
    xs = []
    res = []
    for z in z_xs2:
        cvfilter.predict()
        cvfilter.update([z])
        xs.append(cvfilter.x[0])
        res.append(cvfilter.y[0])

    xs = np.asarray(xs)
    plot_track_and_residuals(t, xs, z_xs2, res)

show_residual_chart()
```



On the left I have plotted the noisy measurements against the Kalman filter output. On the right I display the residuals computed by the filter - the difference between the measurement and the predictions made by the Kalman filter. Let me emphasize this to make this clear. The plot on the right is not merely the difference between the two lines in the left plot. The left plot shows the difference between the measurements and the final Kalman filter output, whereas the right plot shows us the difference between the plot and the predictions of the process model.

That may seem like a subtle distinction, but from the plots you see it is not. The amount of deviation in the left plot when the maneuver starts is small, but the deviation in the right tells a very different story. If the tracked object was moving according to the process model the residual plot should bounce around 0.0. This is because the measurements will be obeying the equation

$$\text{measurement} = \text{process\_model}(t) + \text{noise}(t)$$

Once the target starts maneuvering the predictions of the target behavior will not match the behavior as the equation will be

$$\text{measurement} = \text{process\_model}(t) + \text{maneuver\_delta}(t) + \text{noise}(t)$$

Therefore if the residuals diverge from a mean of 0.0 we know that a maneuver has commenced.

We can see from the residual plot that we have our work cut out for us. We can clearly see the result of the maneuver in the residual plot, but the amount of noise in the signal obscures the start of the maneuver. This is our age old problem of extracting the signal from the noise.

## 14.4 Adjustable Process Noise

The first approach we will consider will use a lower order model and adjust the process noise based on whether a maneuver is occurring or not. When the residual gets “large” (for some reasonable definition of large) we will increase the process noise. This will cause the filter to favor the measurement over the process prediction and the filter will track the signal closely. When the residual is small we will then scale back the process noise.

There are many ways of doing this in the literature, I will consider a couple of choices.

### 14.4.1 Continuous Adjustment

The first method (from Bar-Shalom [1]) normalizes the square of the residual using the following equation:

$$\epsilon = \mathbf{y}^T \mathbf{S}^{-1} \mathbf{y}$$

where  $\mathbf{y}$  is the residual and  $\mathbf{S}$  is the measurement covariance, which has the equation

$$\mathbf{S} = \mathbf{H} \mathbf{P} \mathbf{H}^T + \mathbf{R}$$

If the linear algebra used to compute this confuses you, recall that we can think of matrix inverses as division, so  $\epsilon = \mathbf{y}^T \mathbf{S}^{-1} \mathbf{y}$  can be thought of as computing

$$\epsilon = \frac{\mathbf{y}^2}{\mathbf{S}}$$

Both  $\mathbf{y}$  and  $\mathbf{S}$  are attributes of `filterpy.KalmanFilter` so implementation is of this computation will be straightforward.

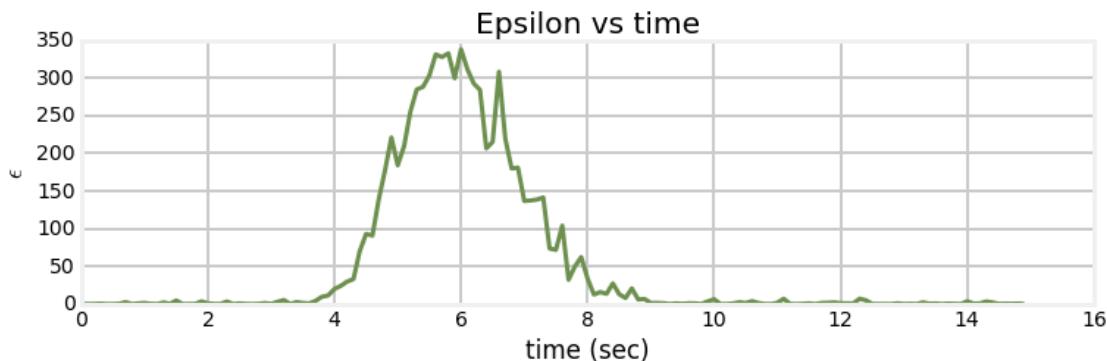
Let's look at a plot of  $\epsilon$  against time.

```
In [15]: from filterpy.common import dot3
        from numpy.linalg import inv

        sensor_std = 0.2
        cvfilter= make_cv_filter(dt, sensor_std)
        _, zs2 = generate_data(120, sensor_std)

        epss = []
        for z in zs2[:, 0]:
            cvfilter.predict()
            cvfilter.update([z])
            y = cvfilter.y
            S = cvfilter.S
            eps = dot3(y.T, inv(cvfilter.S), y)
            epss.append(eps)

        plt.plot(t, epss)
        bp.set_labels(title='Epsilon vs time',
                      x='time (sec)', y='$\epsilon$')
```



This plot should make clear the effect of normalizing the residual. Squaring the residual ensures that the signal is always greater than zero, and normalizing by the measurement covariance scales the signal so that we can distinguish when the residual is markedly changed relative to the measurement noise. The maneuver starts at  $t=3$  seconds, and we can see that  $\epsilon$  starts to increase rapidly not long after that.

We will want to start scaling  $\mathbf{Q}$  up once  $\epsilon$  exceeds some limit, and back down once it again falls below that limit. We multiply  $\mathbf{Q}$  by a scaling factor. Perhaps there is literature on choosing this factor analytically; I derive it experimentally. We can be somewhat more analytical about choosing the limit for  $\epsilon$  (named  $\epsilon_{max}$ ) - generally speaking once the residual is greater than 3 standard deviations or so we can assume the difference is due to a real change and not to noise. However, sensors are rarely truly Gaussian and so a larger number, such as 5-6 standard deviations is used in practice.

I have implemented this algorithm using reasonable values for  $\epsilon_{max}$  and the  $\mathbf{Q}$  scaling factor. To make inspection of the result easier I have limited the plot to the first 10 seconds of simulation.

```
In [16]: from filterpy.common import dot3
        from numpy.linalg import inv

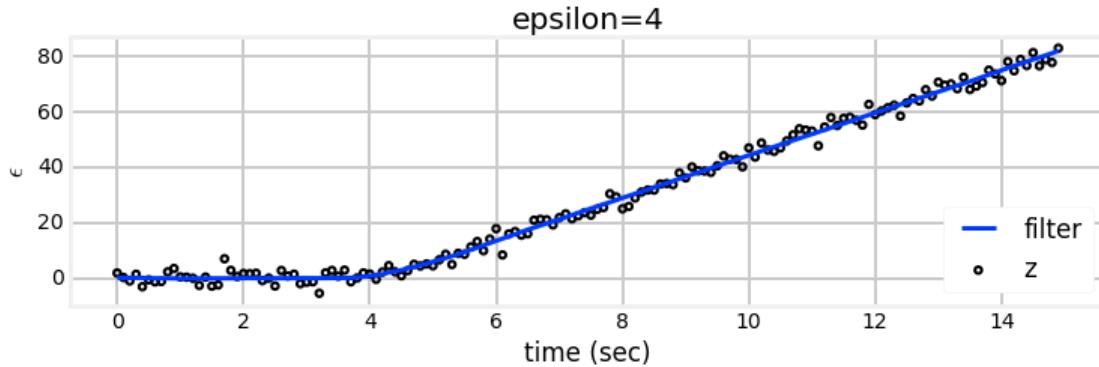
        # reinitialize filter
        sensor_std = 0.2
        cvfilter= make_cv_filter(dt, sensor_std)
        _, zs2 = generate_data(120, sensor_std)

        Q_scale_factor = 1000.
        eps_max = 4.

        epss = []
        xs = []
        count = 0
        for i, z in zip(t, zs2[:, 0]):
            cvfilter.predict()
            cvfilter.update([z])
            y = cvfilter.y
            S = cvfilter.S
            eps = dot3(y.T, inv(cvfilter.S), y)
            epss.append(eps)
            xs.append(cvfilter.x[0])

            if eps > eps_max:
                cvfilter.Q *= Q_scale_factor
                count += 1
            elif count > 0:
                cvfilter.Q /= Q_scale_factor
                count -= 1

        bp.plot_measurements(t, z_xs2, lw=6, label='z')
        bp.plot_filter(t, xs, label='filter')
        plt.legend(loc=4)
        bp.set_labels(title='epsilon=4', x='time (sec)', y='$\epsilon$')
```



The performance of this filter is markedly better than the constant velocity filter. The constant velocity filter took roughly 10 seconds to reacquire the signal after the start of the maneuver. The adaptive filter takes under a second to do the same.

#### 14.4.2 Continuous Adjustment - Standard Deviation Version

Another, very similar method from Zarchan [2] sets the limit based on the standard deviation of the measurement error covariance. Here the equations are:

$$\begin{aligned} std &= \sqrt{\mathbf{H}\mathbf{P}\mathbf{H}^T + \mathbf{R}} \\ &= \sqrt{\mathbf{S}} \end{aligned}$$

If the absolute value of the residual is more than some multiple of the standard deviation computed above we increase the process noise by a fixed amount, recompute Q, and continue.

```
In [17]: from numpy.linalg import inv
        from math import sqrt

def zarchan_adaptive_filter(Q_scale_factor, std_scale,
                           std_title=False,
                           Q_title=False):
    cvfilter = make_cv_filter(dt, std=0.2)
    pos2, zs2 = generate_data(120, std=0.2)
    xs2 = pos2[:,0]
    z_xs2 = zs2[:,0]

    # reinitialize filter
    initialize_filter(cvfilter)
    cvfilter.R = np.eye(1)*0.2

    phi = 0.02
    cvfilter.Q = Q_discrete_white_noise(dim=2, dt=dt, var=phi)
    xs = []
    ys = []
    count = 0
    for i,z in zip(t,z_xs2):
        cvfilter.predict()
```

```

cvfilter.update([z])
y = cvfilter.y
S = cvfilter.S
std = sqrt(S)

xs.append(cvfilter.x)
ys.append(y)

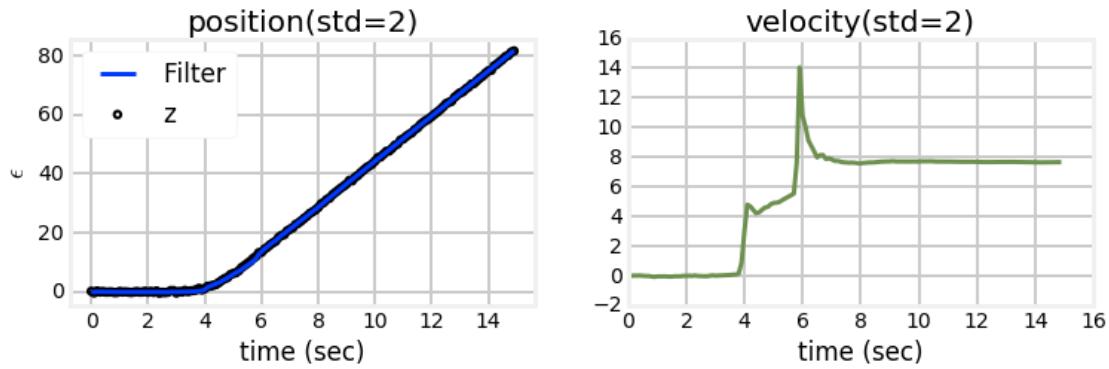
if abs(y[0]) > std_scale*std:
    phi += Q_scale_factor
    cvfilter.Q = Q_discrete_white_noise(2, dt, phi)
    count += 1
elif count > 0:
    phi -= Q_scale_factor
    cvfilter.Q = Q_discrete_white_noise(2, dt, phi)
    count -= 1

xs = np.asarray(xs)
plt.subplot(121)
bp.plot_measurements(t, z_xs2, label='z')
bp.plot_filter(t, xs[:, 0])
bp.set_labels(x='time (sec)', y='$\epsilon$')
plt.legend(loc=2)
if std_title:
    plt.title('position(std={})'.format(std_scale))
elif Q_title:
    plt.title('position(Q scale={})'.format(Q_scale_factor))
else:
    plt.title('position')

plt.subplot(122)
plt.plot(t, xs[:, 1])
plt.xlabel('time (sec)')
if std_title:
    plt.title('velocity(std={})'.format(std_scale))
elif Q_title:
    plt.title('velocity(Q scale={})'.format(Q_scale_factor))
else:
    plt.title('velocity')
plt.show()

zarchan_adaptive_filter(1000, 2, std_title=True)

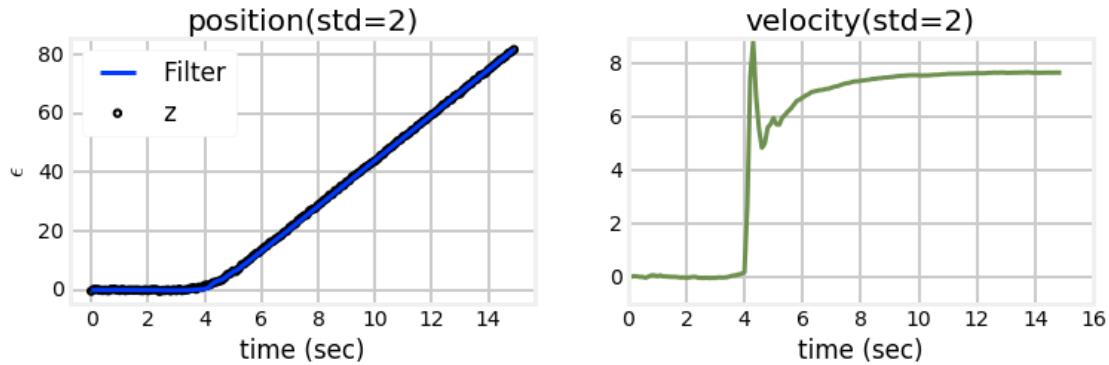
```



So I chose to use 1000 as the scaling factor for the noise, and 2 as the standard deviation limit. Why these numbers? Well, first, let's look at the difference between 2 and 3 standard deviations.

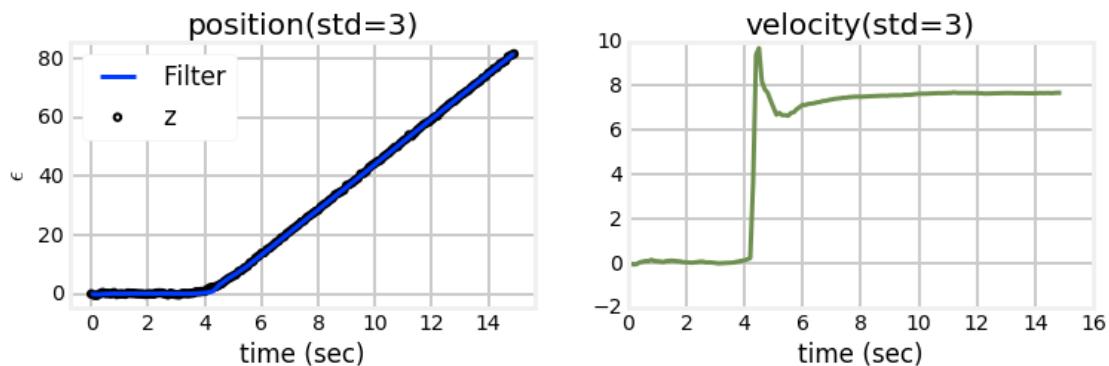
### Two Standard Deviations

```
In [18]: zarchan_adaptive_filter(1000, 2, std_title=True)
```



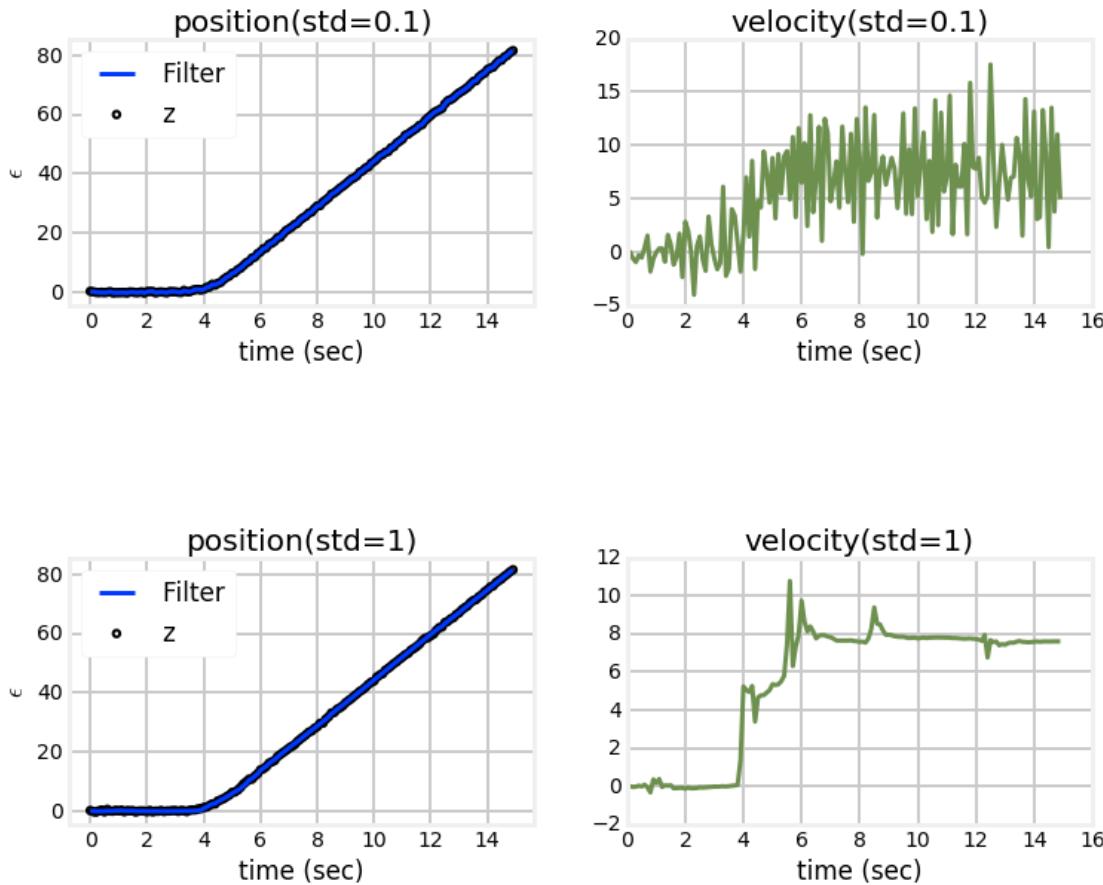
### Three Standard Deviations

```
In [19]: zarchan_adaptive_filter(1000, 3, std_title=True)
```



We can see from the charts that the filter output for the position is very similar regardless of whether we use 2 standard deviations or three. But the computation of the velocity is a different matter. Let's explore this further. First, let's make the standard deviation very small.

```
In [20]: zarchan_adaptive_filter(1000, .1, std_title=True)
zarchan_adaptive_filter(1000, 1, std_title=True)
```



As the standard deviation limit gets smaller the computation of the velocity gets worse. Think about why this is so. If we start varying the filter so that it prefers the measurement over the prediction as soon as the residual deviates even slightly from the prediction we very quickly be giving almost all the weight towards the measurement. With no weight for the prediction we have no information from which to create the hidden variables. So, when the limit is 0.1 std you can see that the velocity is swamped by the noise in the measurement. On the other hand, because we are favoring the measurements so much the position follows the maneuver almost perfectly.

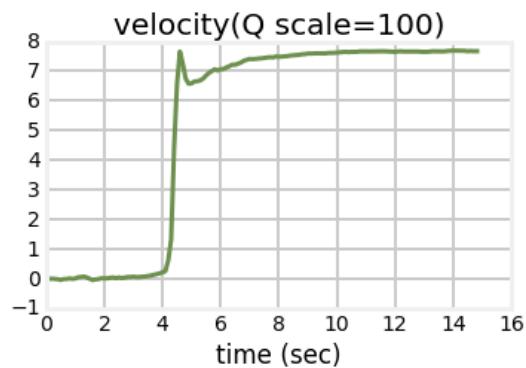
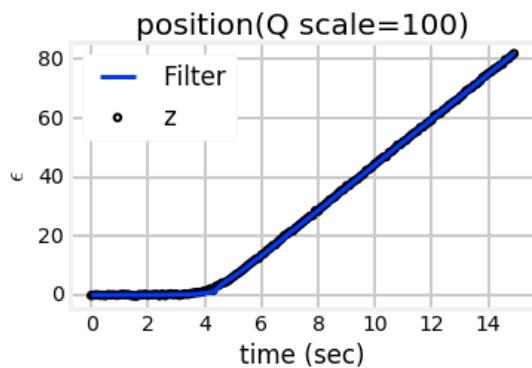
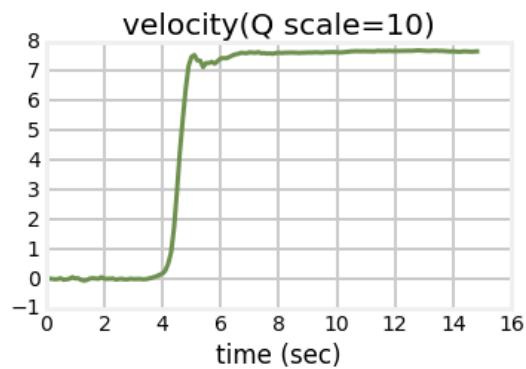
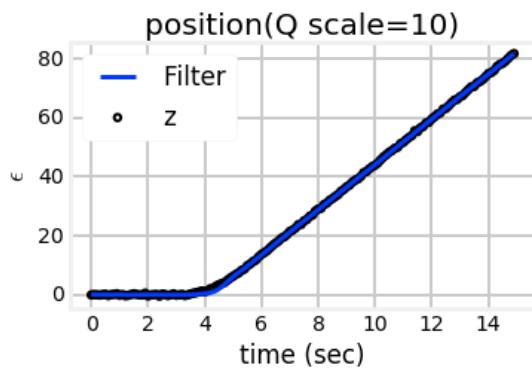
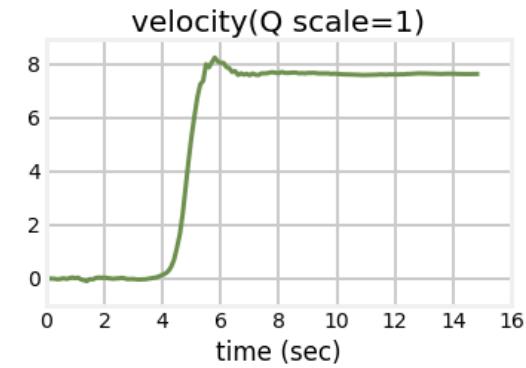
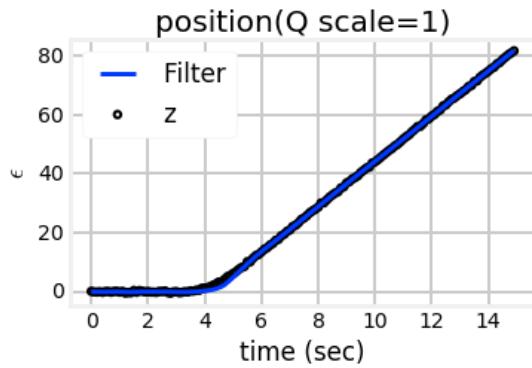
Now let's look at the effect of various increments for the process noise. Here I have held the standard deviation limit to 2 std, and varied the increment from 1 to 10,000.

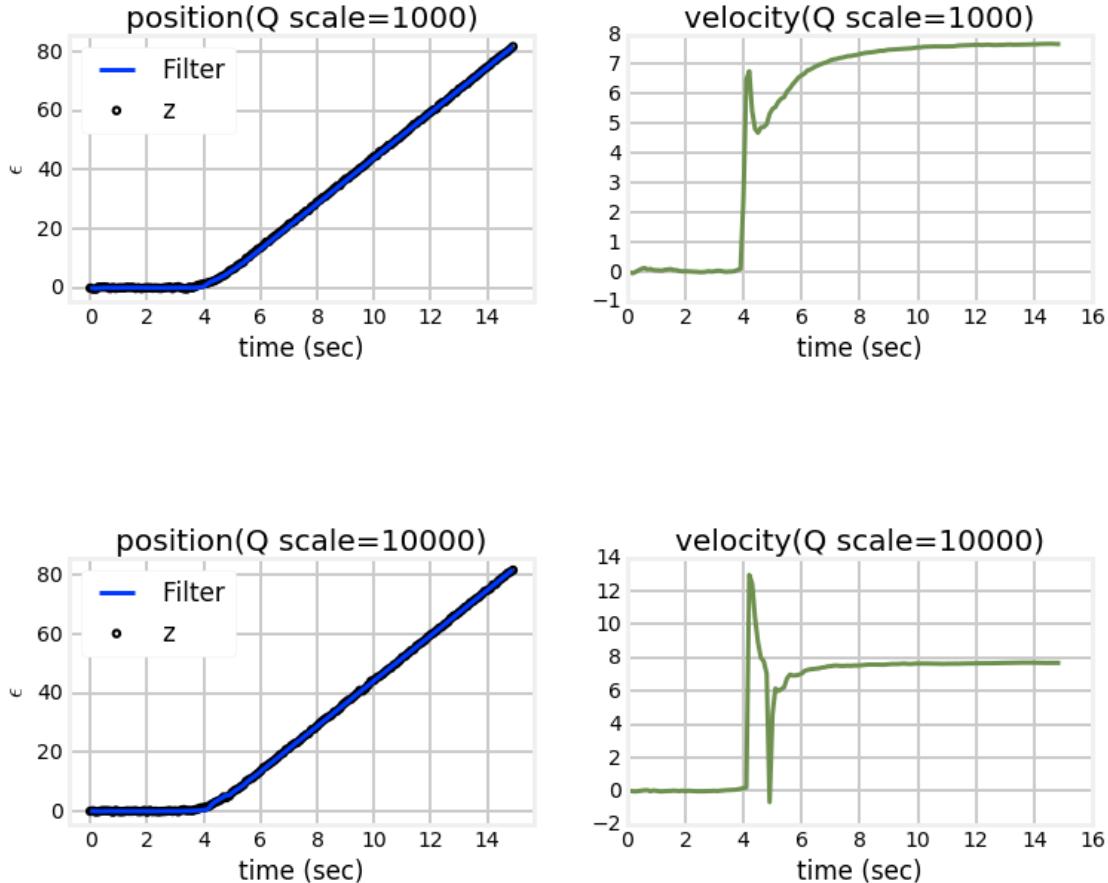
```
In [21]: zarchan_adaptive_filter(1, 2, Q_title=True)
zarchan_adaptive_filter(10, 2, Q_title=True)
```

```

zarchan_adaptive_filter(100, 2, Q_title=True)
zarchan_adaptive_filter(1000, 2, Q_title=True)
zarchan_adaptive_filter(10000, 2, Q_title=True)

```





Here we can see that the position estimate gets marginally better as the increment factor increases, but that the velocity estimate starts to create a large overshoot.

It isn't possible for me to tell you which of these is 'correct'. You will need to test your filter's performance against real and simulated data, and choose the design that best matches the performance you need for each of the state variables.

## 14.5 Fading Memory Filter

Fading memory filters are not normally classified as an adaptive filter since they do not adapt to the input, but they do provide good performance with maneuvering targets. They also have the benefit of having a very simple computational form for first, second, and third order kinematic filters (e.g. the filters we are using in this chapter). This simple form does not require the Riccati equations to compute the gain of the Kalman filter, which drastically reduces the amount of computation. However, there is also a form that works with the standard Kalman filter. I will focus on the latter in this chapter since our focus is more on adaptive filters. Both forms of the fading memory filter are implemented in `FilterPy`.

The Kalman filter is recursive, but it incorporates all of the previous measurements into the current computation of the filter gain. If the target behavior is consistent with the process model than this allows the Kalman filter to find the optimal estimate for every measurement. Consider a ball in flight - we can clearly estimate the position of the ball at time  $t$  better if we take into account all the previous measurement. If we only used some of the measurements we would be less certain about the current position, and thus more

influenced by the noise in the measurement. If this is still not clear, consider the worst case. Suppose we forget all but the last measurement and estimates. We would then have no confidence in the position and trajectory of the ball, and would have little choice but to weight the current measurement heavily. If the measurement is noisy, the estimate is noisy. We see this effect every time a Kalman filter is initialized. The early estimates are noisy, but then they settle down as more measurements are acquired.

However, if the target is maneuvering it is not always behaving like the process model predicts. In this case remembering all of the past measurements and estimates is a liability. We can see this in all of the charts above. The target initiates a turn, and the Kalman filter continues to project movement in a straight line. This is because the filter has built a history of the target's movement, and incorrectly 'feels' confident that the target is moving in a straight line at a given heading and velocity.

The fading memory filter accounts for this problem by giving less weight to older measurements, and greater weight to the more recent measurements.

There are many formulations for the fading memory filter; I use the one provided by Dan Simon in [Optimal State Estimation](#) [3]. I will not go through his derivation, but only provide the results.

The Kalman filter equation for the covariances of the estimation error is

$$\mathbf{P} = \mathbf{F}\mathbf{P}\mathbf{F}^T + \mathbf{Q}$$

We can force the filter to forget past measurements by multiplying a term  $\alpha$

$$\tilde{\mathbf{P}} = \alpha^2 \mathbf{F}\mathbf{P}\mathbf{F}^T + \mathbf{Q}$$

where  $\alpha > 1.0$ . If  $\alpha == 1$  then we get the normal Kalman filter performance.  $\alpha$  is an attribute of `FilterPy` Kalman filter class; its value defaults to 1 so the filter acts like a Kalman filter unless  $\alpha$  is assigned a value other than 1. There is no hard and fast rule for choosing  $\alpha$ , but it is typically very close to 1, such as 1.01. You will need to make many runs with either simulated or real data to determine a value that responds to maneuvers without causing the estimate to become too noisy due to overly weighting the noisy measurement.

Why does this work? If we increase the estimate error covariance the filter becomes more uncertain about it's estimate, hence it gives more weight to the measurement.

One caveat - if we use  $\alpha$  than we are computing  $\tilde{\mathbf{P}}$ , not  $\mathbf{P}$ . In other words, `KalmanFilter.P` is not equal to the covariance of the estimation error, so do not treat it as if it is.

Let's filter our data using the fading memory filter and see the result. I will inject a lot of error into the system so that we can compare various approaches.

```
In [22]: pos2, zs2 = generate_data(70, std=1.2)
xs2 = pos2[:, 0]
z_xs2 = zs2[:, 0]

cvfilter = make_cv_filter(dt, std=1.2)
cvfilter.x.fill(0.)
cvfilter.Q = Q_discrete_white_noise(dim=2, dt=dt, var=0.02)
cvfilter.alpha = 1.00

xs, res = [], []
for z in z_xs2:
    cvfilter.predict()
    cvfilter.update([z])
    xs.append(cvfilter.x[0])
    res.append(cvfilter.y[0])
xs = np.asarray(xs)
```

```

plt.subplot(221)
bp.plot_measurements(t[0:100], z_xs2, label='z')
plt.plot(t[0:100], xs, label='filter')
plt.legend(loc=2)
plt.title('Standard Kalman Filter')

cvfilter.x.fill(0.)
cvfilter.Q = Q_discrete_white_noise(dim=2, dt=dt, var=20.)
cvfilter.alpha = 1.00

xs, res = [], []
for z in z_xs2:
    cvfilter.predict()
    cvfilter.update([z])
    xs.append(cvfilter.x[0])
    res.append(cvfilter.y[0])

xs = np.asarray(xs)

plt.subplot(222)
bp.plot_measurements(t[0:100], z_xs2, label='z')
plt.plot(t[0:100], xs, label='filter')
plt.legend(loc=2)
plt.title('$\\mathbf{Q}=20$')

cvfilter.x.fill(0.)
cvfilter.Q = Q_discrete_white_noise(dim=2, dt=dt, var=0.02)
cvfilter.alpha = 1.02

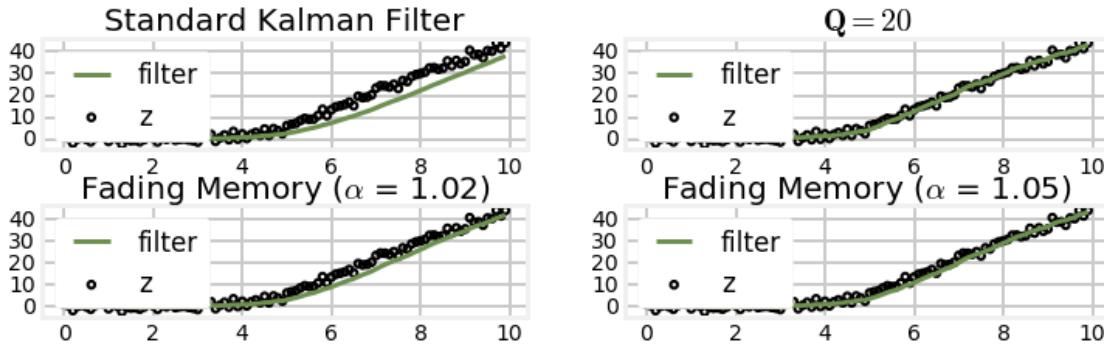
xs, res = [], []
for z in z_xs2:
    cvfilter.predict()
    cvfilter.update([z])
    xs.append(cvfilter.x[0])
    res.append(cvfilter.y[0])
xs = np.asarray(xs)
plt.subplot(223)
bp.plot_measurements(t[0:100], z_xs2, label='z')
plt.plot(t[0:100], xs, label='filter')
plt.legend(loc=2)
plt.title('Fading Memory ($\\alpha = 1.02$)')

cvfilter.x.fill(0.)
cvfilter.Q = Q_discrete_white_noise(dim=2, dt=dt, var=0.02)
cvfilter.alpha = 1.05

xs, res = [], []
for z in z_xs2:
    cvfilter.predict()
    cvfilter.update([z])
    xs.append(cvfilter.x[0])
    res.append(cvfilter.y[0])
xs = np.asarray(xs)
plt.subplot(224)

```

```
bp.plot_measurements(t[0:100], z_xs2, label='z')
plt.plot(t[0:100], xs, label='filter')
plt.legend(loc=2)
plt.title('Fading Memory ($\alpha$ = 1.05)');
```



The first plot shows the performance of the Kalman filter. The filter diverges when the maneuver starts and does not reacquire the signal until about 10 seconds. I then made the filter track the maneuver very quickly by making the process noise large, but this has the cost of making the filter estimate very noisy due to unduly weighting the noisy measurements. I then implemented a fading memory filter with  $\alpha = 1.02$ . The filtered estimate is very smooth, but it does take a few seconds to converge when the target regains steady state behavior. However, the time to do so is considerably smaller than for the Kalman filter, and the amount of lag is much smaller - the estimate for the fading memory is much closer to the actual track than the Kalman filter's track is. Finally, I bumped up  $\alpha$  to 1.05. Here we can see that the filter responds almost instantly to the maneuver, but that the estimate is not as straight during the steady state operation because the filter is forgetting the past measurements.

This is quite good performance for such a small change in code! Note that there is no ‘correct’ choice here. You will need to design your filter based on your needs and the characteristics of the measurement noise, process noise, and maneuvering behavior of the target.

## 14.6 Noise Level Switching

To be written

## 14.7 Variable State Dimension

To be written - vary state order based on whether a maneuver is happening.

## 14.8 Multiple Model Estimation

The example I have been using in this chapter entails a target moving in a steady state, performing a maneuver, and then returning to a steady state. We have been thinking of this as two models - a constant velocity model, and a constant acceleration model. Whenever you can describe the system as obeying one of a finite set of models you can use Multiple Model (MM) Estimation. This means what it sounds like.

We use a bank of multiple filters, each using a different process to describe the system, and either switch between them or blend them based on the dynamics of the tracked object.

As you might imagine this is a broad topic, and there are many ways of designing and implementing MM estimators. But consider a simple approach for the target we have been tracking in this chapter. One idea would be to simultaneously run a constant velocity and a constant acceleration filter, and to switch between their outputs when we detect a maneuver by inspecting the residuals. Even this choice gives us many options. Consider the dynamics of a turning object. For example, an automobile turns on a wheelbase - the front wheels turn, and the car pivots around the rear wheels. This is a nonlinear process, so for best results we would want to use some type of nonlinear filter (EKF, UKF, etc) to model the turns. On the other hand, a linear constant velocity filter would perform fine for the steady state portions of the travel. So our bank of filters might consist of a linear KF and an EKF filter for the turns. However, neither is particularly well suited for modeling behaviors such as accelerating and braking. So a highly performing MM estimator might contain a bank of many filters, each designed to perform best for a certain performance envelope of the tracked object.

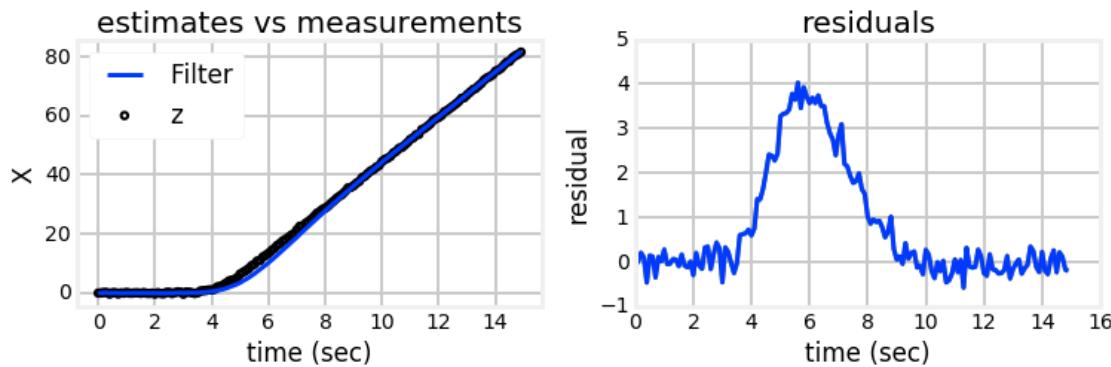
Of course, you do not need to base your filters on the order of the model. You can use different noise models, different adapters in each. For example, in the section above I showed many plots depicting the effects of changing parameters on the estimate of the velocity and position. Perhaps one setting works better for position, and a different setting for velocity. Put both in your bank of filters. You could then take the best estimates for the position from one filter, and the best estimate for the velocity from a different filter.

#### 14.8.1 A Two Filter Adaptive Filter

I trust the idea of switching between filters to get the best performance is clear, but what mathematical foundations should we use to implement it? The problem that we face is trying to detect via noisy measurements when a change in regime should result in a change in model. What aspect of the Kalman filter measures how far the measurement deviates from the prediction? Yes, the residual.

Let's say we have a two state (constant velocity) Kalman filter. As long as the target is not maneuvering the filter will track its behavior closely, and roughly 68% of the measurements should fall within  $1\sigma$ . Furthermore the residual should fluctuate around 0 because as many if the sensor is Gaussian an equal number of measurement should have positive error as have negative errors. If the residual grows and stays beyond predicted bounds then the target must not be performing as predicted by the state model. We saw this earlier in this chart where the residual switched from bouncing around 0 to suddenly jumping and staying above zero once the tracked object began maneuvering.

In [23]: `show_residual_chart()`

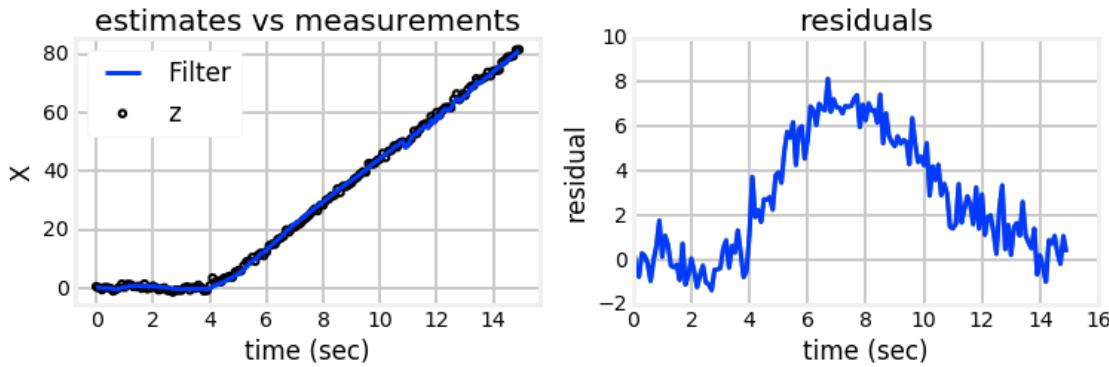


For this problem we saw that the constant velocity filter performed better the constant acceleration filter while the object was in steady state, and the opposite was true when the object is maneuvering. In the chart above that transition occurs at 4 seconds.

So the algorithm is easy. Initialize both a constant velocity and constant acceleration filter and run them together in a predict/update loop. After every update examine the residual of the constant velocity filter. If it falls within theoretical bounds use the estimate from the constant velocity filter as the estimate, otherwise use the estimate from the constant acceleration filter.

```
In [24]: def run_filter_bank(threshold, show_zs=True):
    dt = 0.1
    cvfilter= make_cv_filter(dt, std=0.8)
    cafilter = make_ca_filter(dt, std=0.8)
    pos, zs = generate_data(120, std=0.8)
    z_xs = zs[:, 0]
    t = np.arange(0, len(z_xs) * dt, dt)
    xs, res = [], []
    for z in z_xs:
        cvfilter.predict()
        cafilter.predict()
        cvfilter.update([z])
        cafilter.update([z])
        std = np.sqrt(cvfilter.R[0,0])
        if abs(cvfilter.y[0]) < 2 * std:
            xs.append(cvfilter.x[0])
        else:
            xs.append(cafilter.x[0])
        res.append(cvfilter.y[0])
    xs = np.asarray(xs)
    if show_zs:
        plot_track_and_residuals(t, xs, z_xs, res)
    else:
        plot_track_and_residuals(t, xs, None, res)

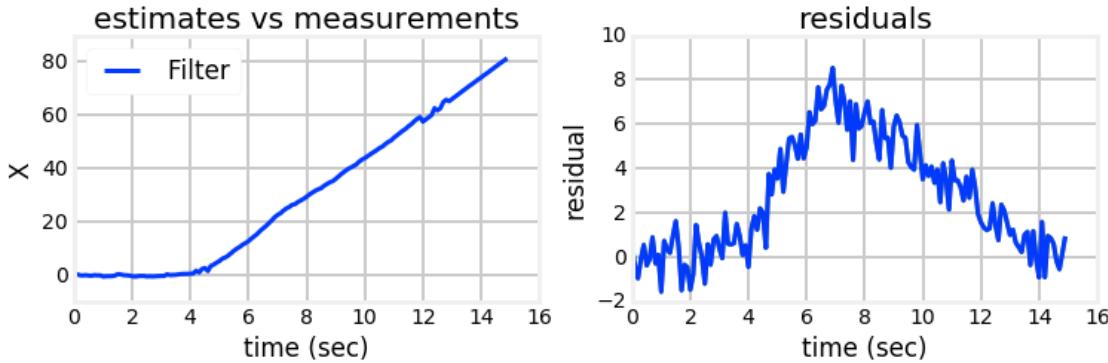
run_filter_bank(threshold=1.4)
```



Here the filter tracks the maneuver closely. While the target is not maneuvering our estimate is nearly noise

free, and then once it does maneuver we quickly detect that and switch to the constant acceleration filter. However, it is not ideal. Here is the filter output plotted alone:

```
In [25]: run_filter_bank(threshold=1.4, show_zs=False)
```



You can see that the estimate jumps when the filter bank switches from one filter to the other. I would not use this algorithm in a production system. The next section gives a state of the art implementation of a filter bank that eliminates this problem.

## 14.9 MMAE

The core idea of using several filters to detect a maneuver is sound, but the estimate is jagged when we abruptly transition between the filters. Choosing one filter over the other flies in the face of this entire book, which uses probability to determine the likelihood of measurements and models. We don't choose either the measurement or prediction, depending on which is more likely, we choose a blend of the two in proportion to their likelihoods. We should do the same here. This approach is called the Multiple Model Adaptive Estimator, or MMAE.

In the **Designing Kalman Filters** chapter we learned the likelihood function

$$\mathcal{L} = \frac{1}{\sqrt{2\pi S}} \exp\left[-\frac{1}{2}\mathbf{y}^\top \mathbf{S}^{-1} \mathbf{y}\right]$$

which tells us how likely a filter is to be performing optimally given the inputs. We can use this to compute the probability that each filter is the best fit to the data. If we have  $N$  filters, we can compute the probability that filter  $i$  is correct in relation to the rest of the filters with

$$p_k^i = \frac{\mathcal{L}_k^i p_{k-1}^i}{\sum_{j=1}^N \mathcal{L}_k^j p_{k-1}^j}$$

That looks messy, but it is straightforward. The numerator is just the likelihood from this time step multiplied by the probability that this filter was correct at the last time frame. We need all of the probabilities for the filter to sum to one, so we normalize by the probabilities for all of the other filters with the term in the denominator.

That is a recursive definition, so we need to assign some initial probability for each filter. In the absence of better information, use  $\frac{1}{N}$  for each. Then we can compute the estimated state as the sum of the state from each filter multiplied by the probability of that filter being correct.

Here is a complete implementation:

```
In [26]: def run_filter_bank():
    dt = 0.1
    cvfilter = make_cv_filter(dt, std=0.2)
    cafilter = make_ca_filter(dt, std=0.2)

    _, zs = generate_data(120, std=0.2)
    z_xs = zs[:, 0]
    #z_xs = [i for i in range(120)]
    t = np.arange(0, len(z_xs) * dt, dt)
    xs, probs = [], []

    pv, pa = 0.8, 0.2
    pvsum, pasum = 0., 0.

    for i, z in enumerate(z_xs):
        cvfilter.predict()
        cafilter.predict()
        cvfilter.update([z])
        cafilter.update([z])

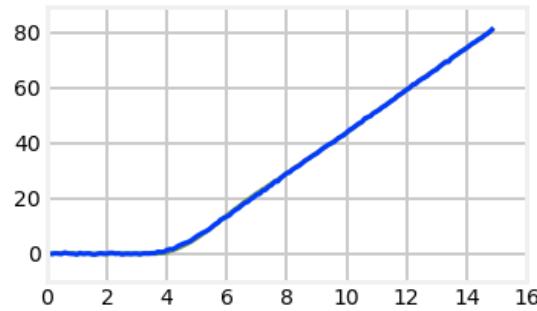
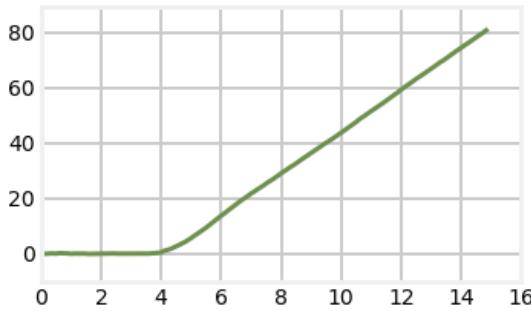
        cv_likelihood = cvfilter.likelihood * pv
        ca_likelihood = cafilter.likelihood * pa

        pv = (cv_likelihood) / (cv_likelihood + ca_likelihood)
        pa = (ca_likelihood) / (cv_likelihood + ca_likelihood)

        x = (pv * cvfilter.x[0]) + (pa*cafilter.x[0])
        xs.append(x)
        probs.append(pv / pa)

    xs = np.asarray(xs)
    plt.subplot(121)
    plt.plot(t, xs)
    plt.subplot(122)
    plt.plot(t, xs)
    plt.plot(t, z_xs)
    return xs, probs

xs, probs = run_filter_bank()
```



I plot the filter's estimates alone on the left so you can see how smooth the result is. On the right I plot both the estimate and the measurements to prove that the filter is tracking the maneuver.

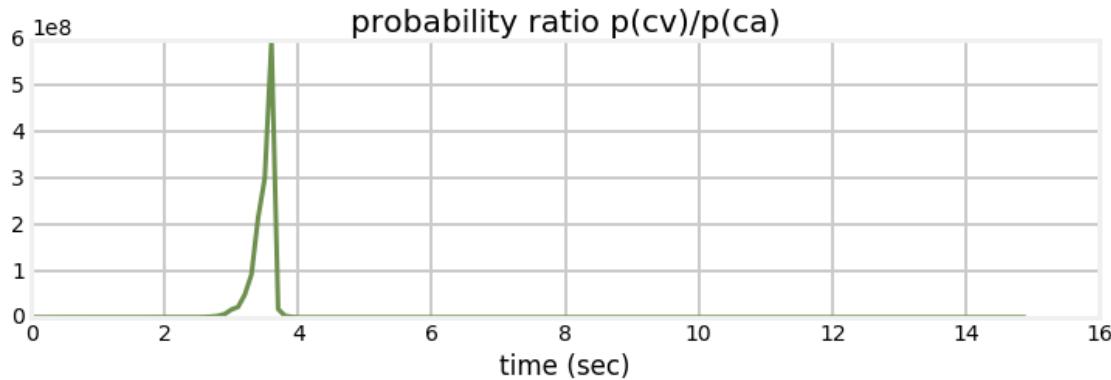
Again I want to emphasize that this is nothing more than the Bayesian algorithm we have been using throughout the book. We have two (or more) measurements or estimate, each with an associated probability. We choose an estimate as a weighted combination of each of those values, where the weights are proportional to the probability of correctness. The computation of the probability at each step is  $\frac{\text{Prob(meas|state)} \times \text{prior}}{\text{normalization}}$ , which is Bayes theorem.

For real world problems you are likely to need more than two filters in your bank. In my job I track objects using computer vision. I track hockey pucks. Pucks slide, they bounce and skitter, they roll, they ricochet, they are picked up and carried, and they are 'dribbled' quickly by the players. I track humans who are athletes, and their capacity for nonlinear behavior is nearly limitless. A two filter bank doesn't get very far in those circumstances. I need to model multiple process models, different assumptions for noise due to the computer vision detection, and so on. But you have the main idea.

#### 14.9.1 Limitations of the MMAE Filter

The MMAE as I have presented it has a significant problem. Look at this chart of the ratio of the probability for the constant velocity vs constant acceleration filter.

```
In [27]: plt.plot(t, probs)
plt.title('probability ratio p(cv)/p(ca)')
plt.xlabel('time (sec)');
```



For the first three seconds, while the tracked object travels in a straight direction, the constant velocity filter becomes much more probable than the constant acceleration filter. Once the maneuver starts the probability quickly changes to favor the constant acceleration model. However, the maneuver is completed by second six. You might expect that the probability for the constant velocity filter would once again become large, but instead it remains at zero.

This happens because of the recursive computation of the probability:

$$p_k = \frac{\mathcal{L}p_{k-1}}{\sum \text{probabilities}}$$

Once the probability becomes very small it can never recover. The result is that the filter bank quickly converges on only the most probable filters. A robust scheme needs to monitor the probability of each filter

and kill off the filters with very low probability and replaced with filters with greater likelihood of performing well. You can subdivide the existing filters into new filters that try to span the characteristics that make them perform well. In the worst case, if a filter has diverged you can reinitialize a filter's state so that it is closer to the current measurements.

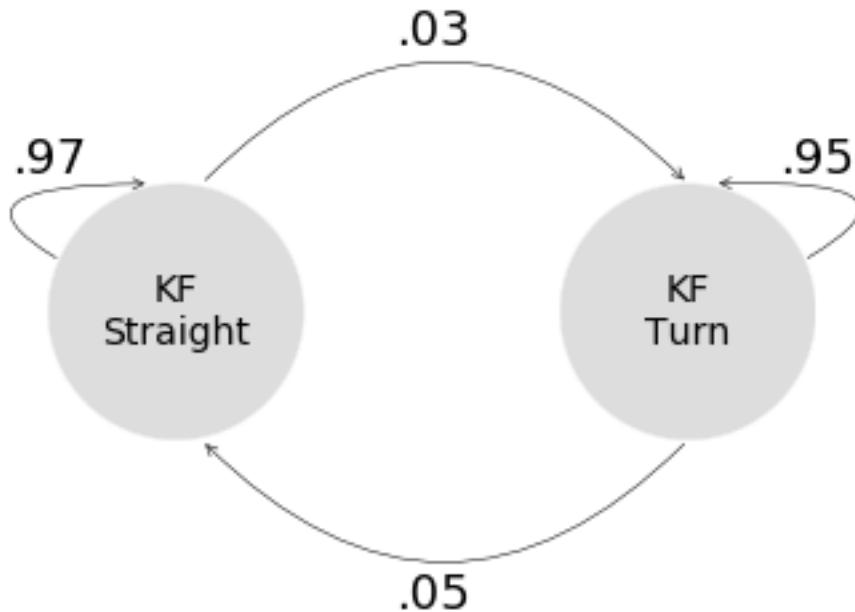
## 14.10 Interacting Multiple Models (IMM)

Let's think about multiple models in another way. The scenario is as before - we wish to track a maneuvering target. We can design a set of Kalman filters which make different modeling assumptions. They can differ in terms of the filter order, or in the amount of noise in the process model. As each new measurement comes in each filter has a probability of being the correct model.

This naive approach leads to combinatorial explosion. At step 1 we generate  $N$  hypotheses, or 1 per filter. At step 2 we generate another  $N$  hypotheses which then need to be combined with the prior  $N$  hypotheses, which yields  $N^2$  hypothesis. Many different schemes have been tried which either cull unlikely hypotheses or merge similar ones, but the algorithms still suffered from computational expense and/or poor performance. I will not cover these in this book, but prominent examples in the literature are the generalized pseudo Bayes (GPB) algorithms.

The Interacting Multiple Models (IMM) algorithm was invented by Blom [4] to solve the combinatorial problem of multiple models. The IMM algorithm uses  $N$  Kalman filters to represent the various possible modes of the system. It models the transition between these modes as a Markov chain, as in this illustration:

```
In [28]: import adaptive_internal
adaptive_internal.plot_markov_chain()
```



This shows an example of two Kalman filters, one that tracks a straight moving target, and a second filter that tracks a turning object. If the current mode of the target is straight, then we predict that there is a 97% chance of the target continuing straight, and a 3% chance of starting a turn. Once the target is turning, we then predict that there is a 95% chance of staying in the turn, and a 5% of returning to a straight path.

The algorithm is not sensitive to the exact numbers, and you will typically use simulation or trials to choose appropriate values. However, these values are quite representative.

We represent Markov chains with a transition probability matrix, which we will call  $p$ . For the Markov chain in the illustration we would write

$$p = \begin{bmatrix} .97 & .03 \\ .05 & .95 \end{bmatrix}$$

This should be clear, but formally  $p_{ij}$  is the probability of transitioning from state  $i$  to state  $j$ .

Next, the algorithm defines a vector of mode probabilities. This is merely a list containing the probability that the system is in any given mode. In the example above we have two modes straight and turn, so we will have a vector of two elements. In the absence of any any additional information we start by assigning equal probability to each mode:

$$\mu = [0.5 \quad 0.5]$$

$\mu$  is typically but not universally used as the symbol for the mode probabilities, so I will use it as well.

We need to update the mode probabilities at each epoch after the measurements are incorporated. We use the likelihood function to compute the likelihood that each Kalman filter, or mode, is the correct one. We can then update the mode probabilities by multiplying them by the likelihood function. In notation, to compute the mode probability:

$$\mu_i = \mu_i \mathcal{L}_i$$

etc.....

## 14.11 References

- [1] Bar-Shalom, Yaakov, Xiao-Rong Li, and Thiagalingam Kirubarajan. Estimation with Applications to Tracking and Navigation. New York: Wiley, p. 424, 2001.
- [2] Zarchan, Paul, and Howard Musoff. Fundamentals of Kalman Filtering: A Practical Approach. Reston, VA: American Institute of Aeronautics and Astronautics, 2000. Print.
- [3] Simon, Dan. Optimal State Estimation: Kalman, H and Nonlinear Approaches. Hoboken, NJ: Wiley-Interscience, p. 208-212, 2006

## Appendix A

# Installation, Python, NumPy, and FilterPy

This book is written in Jupyter Notebook, a browser based interactive Python environment that mixes Python, text, and math. I choose it because of the interactive features - I found Kalman filtering nearly impossible to learn until I started working in an interactive environment. It is difficult to form an intuition of the effect of many of the parameters that you can tune until you can change them rapidly and immediately see the output. An interactive environment also allows you to play ‘what if’ scenarios out. “What if I set  $\mathbf{Q}$  to zero?” It is trivial to find out with Jupyter Notebook.

Another reason I choose it is because I find that a typical textbook leaves many things opaque. For example, there might be a beautiful plot next to some pseudocode. That plot was produced by software, but software that is not available to me as a reader. I want everything that went into producing this book to be available to the reader. How do you plot a covariance ellipse? You won’t know if you read most books. With Jupyter Notebook all you have to do is look at the source code.

Even if you choose to read the book online you will want Python and the SciPy stack installed so that you can write your own Kalman filters. There are many different ways to install these libraries, and I cannot cover them all, but I will cover a few typical scenarios.

### A.1 Installing the SciPy Stack

This book requires IPython, NumPy, SciPy, SymPy, and Matplotlib. The SciPy stack depends on third party Fortran and C code, and is not trivial to install from source code - even the SciPy website strongly urges using a pre-built installation, and I concur with this advice.

I use the Anaconda distribution from Continuum Analytics. This is an excellent package that combines all of the packages listed above, plus many others in one distribution. Installation is very straightforward, and it can be done alongside other Python installation you might already have on your machine. It is free to use, and Continuum Analytics will always ensure that the latest releases of all its packages are available and built correctly for your OS. You may download it from here: <http://continuum.io/downloads>

I strongly recommend using the latest Python 3 version that they provide; I am developing in Python 3.x, and only sporadically test that everything works in Python 2.7. However, it is my long term goal to support 2.7, and if you wish to either not be bothered with 3.x and are willing to be occasionally frustrated with breaking check ins you may stick with Python 2.7.

I am still writing the book, so I do not know exactly what versions of each package is required. I do strongly urge you to use IPython 2.0 or later (this version number has nothing to do with Python 2 vs 3, by the way), as it provides many useful features which I will explain later.

There are other choices for installing the SciPy stack. The SciPy stack provides instructions here: <http://scipy.org/install.html>

## A.2 Installing FilterPy

FilterPy is a Python library that implements all of the filters used in this book, and quite a few not used by this book. Installation is easy using Pip. Issue the following command from the command prompt.

```
pip install filterpy
```

FilterPy is written by me, and the latest development version is always available at [github.com/rlabbe/filterpy](https://github.com/rlabbe/filterpy).

### A.2.1 Manual Install of the SciPy stack

This really isn't easy, and again I advice you to follow the advice of the SciPy website and to use a prepackaged solution. Still, I will give you one example of an install from a fresh operating system. You will have to adapt the instructions to fit your OS and OS version. I will use xubuntu, a linux distribution, because that is what I am most familiar with. I know there are pre-built binaries for Windows (link provided on the SciPy installation web page), and I have no experience with Mac and Python.

I started by doing a fresh install of xubuntu 14.04 on a virtual machine. At that point there is a version of Python 2.7.6 and 3.4 pre-installed. As discussed above you may use either version; I will give the instructions for version 3.4 both because I prefer version 3 and because it can be slightly more tricky because 2.7 is the default Python that runs in xubuntu. Basically the difference is that you have to append a '3' at the end of commands. `python3` instead of `python`, `pip3` instead of `pip`, and so on.

First we will install pip with the following command:

```
sudo apt-get install python3-pip
```

Now you will need to install the various packages from the Ubuntu repositories. Unfortunately they usually do not contain the latest version, so we will also install the development tools necessary to perform an upgrade, which requires compiling various modules.

```
sudo apt-get install python3-numpy python3-scipy python3-matplotlib
sudo apt-get install libblas-dev liblapack-dev gfortran python3-dev
```

Now you can upgrade the packages. This will take a long time as everything needs to be compiled from source. If you get an error I do not know how to help you!

```
sudo pip3 install numpy --upgrade
sudo pip3 install scipy --upgrade
```

Now you get to install SymPy. You can download it from github (replace version number with the most recent version):

```
wget https://github.com/sympy/sympy/releases/download/sympy-0.7.6/sympy-0.7.6.tar.gz
tar -zxvf sympy-0.7.6.tar.gz
```

Now go into the directory you just extracted and run setup.

```
sudo python3 setup.py install
```

If all of this went without a hitch you should have a good install. Try the following. From the command line type `ipython3` to launch ipython, then issue the following commands. Each should run with no exceptions.

```
import numpy as np
import scipy as sp
import sympy as sym
np.__version__
sp.__version__
sym.__version__
```

Now let's make sure plotting works correctly.

```
import matplotlib.pyplot as plt
plt.plot([1, 4, 3])
plt.show()
```

Now you get to fix IPython so Jupyter notebook will run. First, I had to uninstall IPython with

```
sudo pip3 uninstall ipython
```

Then, I reinstalled it with the `[all]` option so that all the required dependencies are installed.

```
sudo pip3 install "ipython3[all]"
```

Now test the installation by typing

```
ipython notebook
```

If successful, it should launch a browser showing you the Jupyter home page.

That was not fun. It actually goes somewhat smoother in Windows, where you can download pre-built binaries for these packages; however, you are dependent on this being done for you in a timely manner. So, please follow the SciPy advice and use a pre-built distribution! I will not be supporting this manual install process.

## A.3 Installing/downloading and running the book

Okay, so now you have the SciPy stack installed, how do you download the book? This is easy as the book is in a github repository. From the command line type the following:

```
git clone https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python.git
```

Alternatively, browse to <https://github.com/rlabbe/Kalman-and-Bayesian-Filters-in-Python> and download it via your browser.

Now, from the command prompt change to the directory that was just created, and then run Jupyter notebook:

```
cd Kalman-and-Bayesian-Filters-in-Python
ipython notebook
```

A browser window should launch showing you all of the chapters in the book. Browse to the first chapter by clicking on it, then open the notebook in that subdirectory by clicking on the link.

## A.4 Using IPython Notebook

A complete tutorial on Jupyter Notebook is beyond the scope of this book. Many are available online. In short, Python code is placed in cells. These are prefaced with text like `In [1]:`, and the code itself is in a boxed area. If you press CTRL-ENTER while focus is inside the box the code will run and the results will be displayed below the box. Like this:

```
In [3]: print(3+7.2)
```

```
10.2
```

If you have this open in Jupyter Notebook now, go ahead and modify that code by changing the expression inside the `print` statement and pressing CTRL+ENTER. The output should be changed to reflect what you typed in the code cell.

## A.5 SymPy

SymPy is a Python package for performing symbolic mathematics. The full scope of its abilities are beyond this book, but it can perform algebra, integrate and differentiate equations, find solutions to differential equations, and much more. For example, we use use to to compute the Jacobian of matrices and the expected value integral computations.

First, a simple example. We will import SymPy, initialize its pretty print functionality (which will print equations using LaTeX). We will then declare a symbol for Numpy to use.

```
In [4]: import sympy
sympy.init_printing(use_latex='mathjax')

phi, x = sympy.symbols('\phi, x')
phi
```

```
Out[4] :
```

$$\phi$$

Notice how we use a latex expression for the symbol `phi`. This is not necessary, but if you do it will render as LaTeX when output. Now let's do some math. What is the derivative of  $\sqrt{\phi}$ ?

```
In [5]: sympy.diff('sqrt(phi)')
```

```
Out[5] :
```

$$\frac{1}{2\sqrt{\phi}}$$

We can factor equations

In [6]: `sympy.factor(phi**3 -phi**2 + phi - 1)`

Out[6]:

$$(\phi - 1) (\phi^2 + 1)$$

and we can expand them.

In [7]: `((phi+1)*(phi-4)).expand()`

Out[7]:

$$\phi^2 - 3\phi - 4$$

You can find the value of an equation by substituting in a value for a variable

In [8]: `w =x**2 -3*x +4  
w.subs(x,4)`

Out[8]:

$$8$$

You can also use strings for equations that use symbols that you have not defined

In [9]: `x = sympy.expand('t+1)*2')  
x`

Out[9]:

$$2t + 2$$

Now let's use SymPy to compute the Jacobian of a matrix. Let us have a function

$$h = \sqrt{x^2 + z^2}$$

for which we want to find the Jacobian with respect to x, y, and z.

In [10]: `x, y, z = sympy.symbols('x y z')  
H = sympy.Matrix([sympy.sqrt(x**2 + z**2)])  
state = sympy.Matrix([x, y, z])  
H.jacobian(state)`

Out[10]:

$$\begin{bmatrix} \frac{x}{\sqrt{x^2+z^2}} & 0 & \frac{z}{\sqrt{x^2+z^2}} \end{bmatrix}$$

Now let's compute the discrete process noise matrix  $\mathbf{Q}_k$  given the continuous process noise matrix

$$\mathbf{Q} = \Phi_s \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

and the equation

$$\mathbf{Q} = \int_0^{\Delta t} \Phi(t) \mathbf{Q} \Phi^T(t) dt$$

where

$$\Phi(t) = \begin{bmatrix} 1 & \Delta t & \Delta t^2/2 \\ 0 & 1 & \Delta t \\ 0 & 0 & 1 \end{bmatrix}$$

```
In [11]: dt = sympy.symbols('Delta{t}')
F_k = sympy.Matrix([[1, dt, dt**2/2],
                    [0, 1, dt],
                    [0, 0, 1]]))
Q = sympy.Matrix([[0,0,0],
                  [0,0,0],
                  [0,0,1]]))

sympy.integrate(F_k*Q*F_k.T,(dt, 0, dt))
```

Out[11]:

$$\begin{bmatrix} \frac{\Delta t^5}{20} & \frac{\Delta t^4}{8} & \frac{\Delta t^3}{6} \\ \frac{\Delta t^4}{20} & \frac{\Delta t^3}{8} & \frac{\Delta t^2}{2} \\ \frac{8}{6} & \frac{3}{2} & \Delta t \end{bmatrix}$$

## Appendix B

# Symbols and Notations

Here is a collection of the notation used by various authors for the linear Kalman filter equations.

### B.1 Labbe

$$\begin{aligned}\bar{\mathbf{x}} &= \mathbf{F}\mathbf{x} + \mathbf{B}\mathbf{u} \\ \bar{\mathbf{P}} &= \mathbf{F}\mathbf{P}\mathbf{F}^T + \mathbf{Q}\end{aligned}$$

$$\begin{aligned}\mathbf{y} &= \mathbf{z} - \mathbf{H}\bar{\mathbf{x}} \\ \mathbf{S} &= \mathbf{H}\bar{\mathbf{P}}\mathbf{H}^T + \mathbf{R} \\ \mathbf{K} &= \bar{\mathbf{P}}\mathbf{H}^T\mathbf{S}^{-1} \\ \mathbf{x} &= \bar{\mathbf{x}} + \mathbf{K}\mathbf{y} \\ \mathbf{P} &= (\mathbf{I} - \mathbf{K}\mathbf{H})\bar{\mathbf{P}}\end{aligned}$$

### B.2 Wikipedia

$$\begin{aligned}\hat{\mathbf{x}}_{k|k-1} &= \mathbf{F}_k \hat{\mathbf{x}}_{k-1|k-1} + \mathbf{B}_k \mathbf{u}_k \\ \mathbf{P}_{k|k-1} &= \mathbf{F}_k \mathbf{P}_{k-1|k-1} \mathbf{F}_k^T + \mathbf{Q}_k \\ \tilde{\mathbf{y}}_k &= \mathbf{z}_k - \mathbf{H}_k \hat{\mathbf{x}}_{k|k-1} \\ \mathbf{S}_k &= \mathbf{H}_k \mathbf{P}_{k|k-1} \mathbf{H}_k^T + \mathbf{R}_k \\ \mathbf{K}_k &= \mathbf{P}_{k|k-1} \mathbf{H}_k^T \mathbf{S}_k^{-1} \\ \hat{\mathbf{x}}_{k|k} &= \hat{\mathbf{x}}_{k|k-1} + \mathbf{K}_k \tilde{\mathbf{y}}_k \\ \mathbf{P}_{k|k} &= (I - \mathbf{K}_k \mathbf{H}_k) \mathbf{P}_{k|k-1}\end{aligned}$$

### B.3 Brookner

$$\begin{aligned} X_{n+1,n}^* &= \Phi X_{n,n}^* \\ X_{n,n}^* &= X_{n,n-1}^* + H_n(Y_n - MX_{n,n-1}^*) \\ H_n &= S_{n,n-1}^* M^\top [R_n + M S_{n,n-1}^* M^\top]^{-1} \\ S_{n,n-1}^* &= \Phi S_{n-1,n-1}^* \Phi^\top + Q_n \\ S_{n-1,n-1}^* &= (I - H_{n-1} M) S_{n-1,n-2}^* \end{aligned}$$

### B.4 Gelb

$$\begin{aligned} \hat{\underline{x}}_k(-) &= \Phi_{k-1} \hat{\underline{x}}_{k-1}(+) \\ \hat{\underline{x}}_k(+) &= \hat{\underline{x}}_k(-) + K_k [Z_k - H_k \hat{\underline{x}}_k(-)] \\ K_k &= P_k(-) H_k^\top [H_k P_k(-) H_k^\top + R_k]^{-1} \\ P_k(+) &= \Phi_{k-1} P_{k-1}(+) \Phi_{k-1}^\top + Q_{k-1} \\ P_k(-) &= (I - K_k H_k) P_k(-) \end{aligned}$$

### B.5 Brown

$$\begin{aligned} \hat{\mathbf{x}}_{k+1}^- &= \phi_k \hat{\mathbf{x}}_k \\ \hat{\mathbf{x}}_k &= \hat{\mathbf{x}}_k^- + \mathbf{K}_k [\mathbf{z}_k - \mathbf{H}_k \hat{x}_k^-] \\ \mathbf{K}_k &= \mathbf{P}_k^- \mathbf{H}_k^\top [\mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^\top + \mathbf{R}_k]^{-1} \\ \mathbf{P}_{k+1}^- &= \phi_k \mathbf{P}_k \phi_k^\top + \mathbf{Q}_k \\ \mathbf{P}_k &= (\mathbf{I} - \mathbf{K}_k \mathbf{H}_k) \mathbf{P}_k^- \end{aligned}$$

### B.6 Zarchan

$$\begin{aligned} \hat{x}_k &= \Phi_k \hat{x}_{k-1} + G_k u_{k-1} + K_k [z_k - H \Phi_k \hat{x}_{k-1} - H G_k u_{k-1}] \\ M_k &= \Phi_k P_{k-1} \phi_k^\top + Q_k \\ K_k &= M_k H^\top [H M_k H^\top + R_k]^{-1} \\ P_k &= (I - K_k H) M_k \end{aligned}$$

### B.7 Bar-Shalom

### B.8 Thrun

## Appendix C

# H Infinity filter

I am still mulling over how to write this chapter. In the meantime, Professor Dan Simon at Cleveant State University has an accessible introduction here:

<http://academic.csuohio.edu/simond/courses/eec641/hinfinity.pdf>

In one sentence the  $H_\infty$  (H infinity) filter is like a Kalman filter, but it is robust in the face of non-Gaussian, non-predictable inputs.

My FilterPy library contains an H-Infinity filter. I've pasted some test code below which implements the filter designed by Simon in the article above. Hope it helps.

```
In [2]: from __future__ import (absolute_import, division, print_function,
                               unicode_literals)

from numpy import array
import matplotlib.pyplot as plt

from filterpy.hinfinity import HInfinityFilter

dt = 0.1
f = HInfinityFilter(2, 1, dim_u=1, gamma=.01)

f.F = array([[1., dt],
             [0., 1.]])  

  
f.H = array([[0., 1.]])
f.G = array([[dt**2 / 2, dt]]).T  

  
f.P = 0.01
f.W = array([[0.0003, 0.005],
             [0.0050, 0.100]])/ 1000 #process noise  

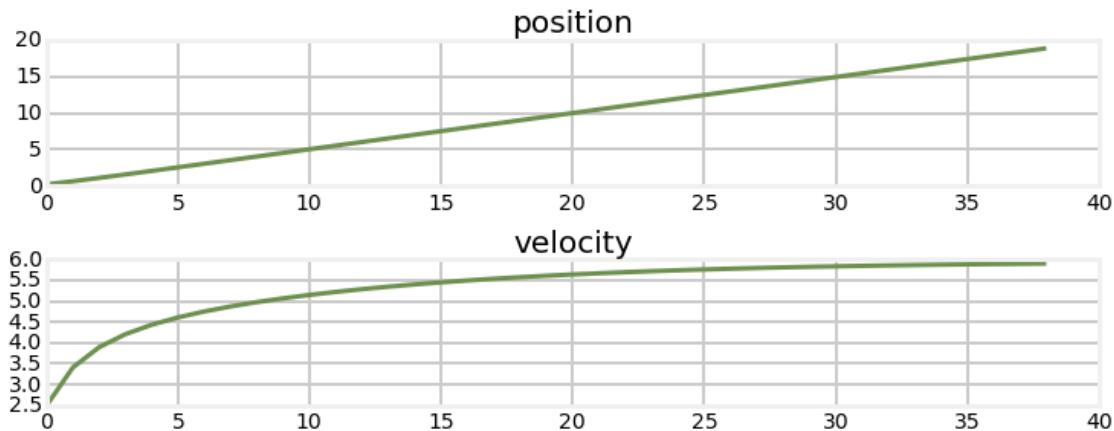
  
f.V = 0.01
f.Q = 0.01
u = 1. #acceleration of 1 f/sec**2  

  
xs = []
vs = []

for i in range(1,40):
```

```
f.update(5)
#print(f.x.T)
xs.append(f.x[0,0])
vs.append(f.x[1,0])
f.predict(u=u)

plt.subplot(211)
plt.plot(xs)
plt.title('position')
plt.subplot(212)
plt.plot(vs)
plt.title('velocity');
```



## Appendix D

# Ensemble Kalman Filters

I am not well versed with Ensemble filters. I have implemented one for this book, and made it work, but I have not used one in real life. Different sources use slightly different forms of these equations. If I implement the equations given in the sources the filter does not work. It is possible that I am doing something wrong. However, in various places on the web I have seen comments by people stating that they do the kinds of things I have done in my filter to make it work. In short, I do not understand this topic well, but choose to present my lack of knowledge rather than to gloss over the subject. I hope to master this topic in the future and to author a more definitive chapter. At the end of the chapter I document my current confusion and questions. In any case if I got confused by the sources perhaps you also will, so documenting my confusion can help you avoid the same.

The ensemble Kalman filter (EnKF) is very similar to the unscented Kalman filter (UKF) of the last chapter. If you recall, the UKF uses a set of deterministically chosen weighted sigma points passed through nonlinear state and measurement functions. After the sigma points are passed through the function, we find the mean and covariance of the points and use this as the filter's new mean and covariance. It is only an approximation of the true value, and thus suboptimal, but in practice the filter is highly accurate. It has the advantage of often producing more accurate estimates than the EKF does, and also does not require you to analytically derive the linearization of the state and measurement equations.

The ensemble Kalman filter works in a similar way, except it uses a [Monte Carlo](#) method to choose a large numbers of sigma points. It came about from the geophysical sciences as an answer for the very large states and systems needed to model things such as the ocean and atmosphere. There is an interesting article on its development in weather modeling in [SIAM News](#) [1]. The filter starts by randomly generating a large number of points distributed about the filter's initial state. This distribution is proportional to the filter's covariance  $\mathbf{P}$ . In other words 68% of the points will be within one standard deviation of the mean, 95% percent within two standard deviations, and so on. Let's look at this in two dimensions. We will use `numpy.random.multivariate_normal()` function to randomly create points from a multivariate normal distribution drawn from the mean  $(5, 3)$  with the covariance

$$\begin{bmatrix} 32 & 15 \\ 15 & 40 \end{bmatrix}$$

I've drawn the covariance ellipse representing two standard deviations to illustrate how the points are distributed.

```
In [2]: import matplotlib.pyplot as plt
        import numpy as np
        from numpy.random import multivariate_normal
```

```

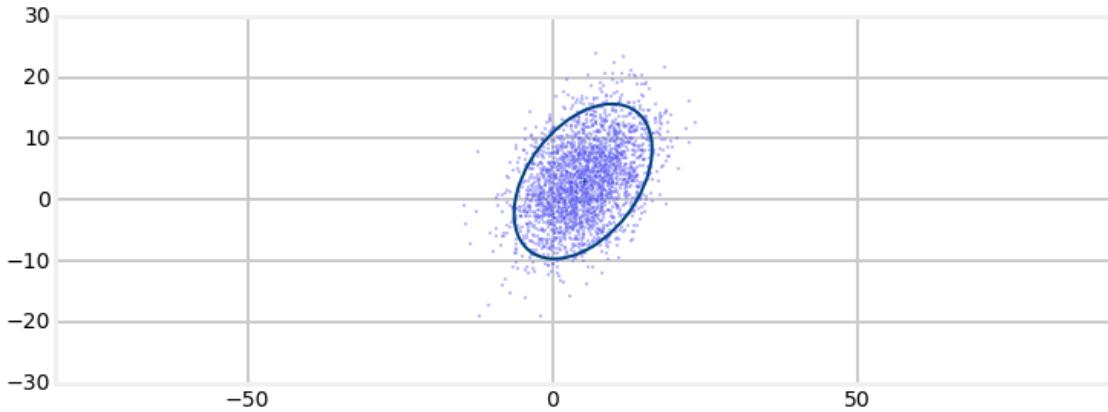
from filterpy.stats import (covariance_ellipse,
                           plot_covariance_ellipse)

mean = (5, 3)
P = np.array([[32, 15],
              [15., 40.]])  
  

x,y = multivariate_normal(mean=mean, cov=P, size=2500).T
plt.scatter(x, y, alpha=0.3, marker='.')
plt.axis('equal')  
  

plot_covariance_ellipse(mean=mean, cov=P,
                        variance=2.**2,
                        facecolor='none')

```



## D.1 The Algorithm

As I already stated, when the filter is initialized a large number of sigma points are drawn from the initial state ( $\mathbf{x}$ ) and covariance ( $\mathbf{P}$ ). From there the algorithm proceeds very similarly to the UKF. During the prediction step the sigma points are passed through the state transition function, and then perturbed by adding a bit of noise to account for the process noise. During the update step the sigma points are translated into measurement space by passing them through the measurement function, they are perturbed by a small amount to account for the measurement noise. The Kalman gain is computed from the

We already mentioned the main difference between the UKF and EnKF - the UKF chooses the sigma points deterministically. There is another difference, implied by the algorithm above. With the UKF we generate new sigma points during each predict step, and after passing the points through the nonlinear function we reconstitute them into a mean and covariance by using the unscented transform. The EnKF keeps propagating the originally created sigma points; we only need to compute a mean and covariance as outputs for the filter!

Let's look at the equations for the filter. As usual, I will leave out the typical subscripts and superscripts; I am expressing an algorithm, not mathematical functions. Here  $N$  is the number of sigma points,  $\chi$  is the set of sigma points.

### D.1.1 Initialize Step

$$\chi \sim \mathcal{N}(\mathbf{x}_0, \mathbf{P}_0)$$

This says to select the sigma points from the filter's initial mean and covariance. In code this might look like

```
N = 1000
sigmas = multivariate_normal(mean=x, cov=P, size=N)
```

### D.1.2 Predict Step

$$\begin{aligned}\chi &= f(\chi, \mathbf{u}) + v_Q \\ \mathbf{x} &= \frac{1}{N} \sum_1^N \chi\end{aligned}$$

That is short and sweet, but perhaps not entirely clear. The first line passes all of the sigma points through a user supplied state transition function and then adds some noise distributed according to the  $\mathbf{Q}$  matrix. In Python we might write

```
for i, s in enumerate(sigmas):
    sigmas[i] = fx(x=s, dt=0.1, u=0.)

sigmas += multivariate_normal(x, Q, N)
```

The second line computes the mean from the sigmas. In Python we will take advantage of `numpy.mean` to do this very concisely and quickly.

```
x = np.mean(sigmas, axis=0)
```

We can now optionally compute the covariance of the mean. The algorithm does not need to compute this value, but it is often useful for analysis. The equation is

$$\mathbf{P} = \frac{1}{N-1} \sum_1^N [\chi - \mathbf{x}^-][\chi - \mathbf{x}^-]^T$$

$\chi - \mathbf{x}^-$  is a one dimensional vector, so we will use `numpy.outer` to compute the  $[\chi - \mathbf{x}^-][\chi - \mathbf{x}^-]^T$  term. In Python we might write

```
P = 0
for s in sigmas:
    P += outer(s-x, s-x)
P = P / (N-1)
```

### D.1.3 Update Step

In the update step we pass the sigma points through the measurement function, compute the mean and covariance of the sigma points, compute the Kalman gain from the covariance, and then update the Kalman state by scaling the residual by the Kalman gain. The equations are

$$\begin{aligned}\boldsymbol{\chi}_h &= h(\boldsymbol{\chi}, u) \\ \mathbf{z}_{mean} &= \frac{1}{N} \sum_1^N \boldsymbol{\chi}_h \\ \mathbf{P}_{zz} &= \frac{1}{N-1} \sum_1^N [\boldsymbol{\chi}_h - \mathbf{z}_{mean}] [\boldsymbol{\chi}_h - \mathbf{z}_{mean}]^\top + \mathbf{R} \\ \mathbf{P}_{xz} &= \frac{1}{N-1} \sum_1^N [\boldsymbol{\chi} - \mathbf{x}^-] [\boldsymbol{\chi}_h - \mathbf{z}_{mean}]^\top \\ \mathbf{K} &= \mathbf{P}_{xz} \mathbf{P}_{zz}^{-1} \\ \boldsymbol{\chi} &= \boldsymbol{\chi} + \mathbf{K} [\mathbf{z} - \boldsymbol{\chi}_h + \mathbf{v}_R]\end{aligned}$$

$$\begin{aligned}\mathbf{x} &= \frac{1}{N} \sum_1^N \boldsymbol{\chi} \\ \mathbf{P} &= \mathbf{P} - \mathbf{K} \mathbf{P}_{zz} \mathbf{K}^\top\end{aligned}$$

This is very similar to the linear KF and the UKF. Let's just go line by line.

The first line,

$$\boldsymbol{\chi}_h = h(\boldsymbol{\chi}, u),$$

just passes the sigma points through the measurement function  $h$ . We name the resulting points  $\boldsymbol{\chi}_h$  to distinguish them from the sigma points. In Python we could write this as

```
sigmas_h = h(sigmas, u)
```

The next line computes the mean of the measurement sigmas.

$$\mathbf{z}_{mean} = \frac{1}{N} \sum_1^N \boldsymbol{\chi}_h$$

In Python we write

```
z_mean = np.mean(sigmas_h, axis=0)
```

Now that we have the mean of the measurement sigmas we can compute the covariance for every measurement sigma point, and the cross variance for the measurement sigma points vs the sigma points. That is expressed by these two equations

$$\begin{aligned}\mathbf{P}_{zz} &= \frac{1}{N-1} \sum_1^N [\boldsymbol{\chi}_h - \mathbf{z}_{mean}] [\boldsymbol{\chi}_h - \mathbf{z}_{mean}]^\top + \mathbf{R} \\ \mathbf{P}_{xz} &= \frac{1}{N-1} \sum_1^N [\boldsymbol{\chi} - \mathbf{x}^-] [\boldsymbol{\chi}_h - \mathbf{z}_{mean}]^\top\end{aligned}$$

We can express this in Python with

```

P_zz = 0
for sigma in sigmas_h:
    s = sigma - z_mean
    P_zz += outer(s, s)
P_zz = P_zz / (N-1) + R

P_xz = 0
for i in range(N):
    P_xz += outer(self.sigmas[i] - self.x, sigmas_h[i] - z_mean)
P_xz /= N-1

```

Computation of the Kalman gain is straightforward  $\mathbf{K} = \mathbf{P}_{xz}\mathbf{P}_{zz}^{-1}$ .

In Python this is

```
python K = np.dot(P_xz, inv(P_zz))
```

Next, we update the sigma points with

$$\chi = \chi + \mathbf{K}[\mathbf{z} - \chi_h + \mathbf{v}_R]$$

Here  $\mathbf{v}_R$  is the perturbation that we add to the sigmas. In Python we can implement this with

```

v_r = multivariate_normal([0]*dim_z, R, N)
for i in range(N):
    sigmas[i] += dot(K, z + v_r[i] - sigmas_h[i])

```

Our final step is recompute the filter's mean and covariance.

```

x = np.mean(sigmas, axis=0)
P = self.P - dot3(K, P_zz, K.T)

```

## D.2 Implementation and Example

I have implemented an EnKF in the `FilterPy` library. It is in many ways a toy. Filtering with a large number of sigma points gives us very slow performance. Furthermore, there are many minor variations on the algorithm in the literature. I wrote this mostly because I was interested in learning a bit about the filter. I have not used it for a real world problem, and I can give no advice on using the filter for the large problems for which it is suited. Therefore I will refine my comments to implementing a very simple filter. I will use it to track an object in one dimension, and compare the output to a linear Kalman filter. This is a filter we have designed many times already in this book, so I will design it with little comment. Our state vector will be

$$\mathbf{x} = \begin{bmatrix} x \\ \dot{x} \end{bmatrix}$$

The state transition function is

$$\mathbf{F} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

and the measurement function is

$$\mathbf{H} = [1 \ 0]$$

The EnKF is designed for nonlinear problems, so instead of using matrices to implement the state transition and measurement functions you will need to supply Python functions. For this problem they can be written as:

```
def hx(x):
    return np.array([x[0]])

def fx(x, dt):
    return np.dot(F, x)
```

One final thing: the EnKF code, like the UKF code, uses a single dimension for **x**, not a two dimensional column matrix as used by the linear kalman filter code.

Without further ado, here is the code.

```
In [3]: from numpy.random import randn
        from filterpy.kalman import EnsembleKalmanFilter as EnKF
        from filterpy.kalman import KalmanFilter
        from filterpy.common import Q_discrete_white_noise
        import book_plots as bp

np.random.seed(1234)

def hx(x):
    return np.array([x[0]])

def fx(x, dt):
    return np.dot(F, x)

F = np.array([[1., 1.],[0., 1.]])

x = np.array([0., 1.])
P = np.eye(2) * 100.
enf = EnKF(x=x, P=P, dim_z=1, dt=1., N=20, hx=hx, fx=fx)

std_noise = 10.
enf.R *= std_noise**2
enf.Q = Q_discrete_white_noise(2, 1., .001)

kf = KalmanFilter(dim_x=2, dim_z=1)
kf.x = np.array([x]).T
kf.F = F.copy()
kf.P = P.copy()
kf.R = enf.R.copy()
kf.Q = enf.Q.copy()
kf.H = np.array([[1., 0.]])

measurements = []
results = []
ps = []
kf_results = []
```

```

zs = []
for t in range (0,100):
    # create measurement = t plus white noise
    z = t + randn()*std_noise
    zs.append(z)

    enf.predict()
    enf.update(np.asarray([z]))

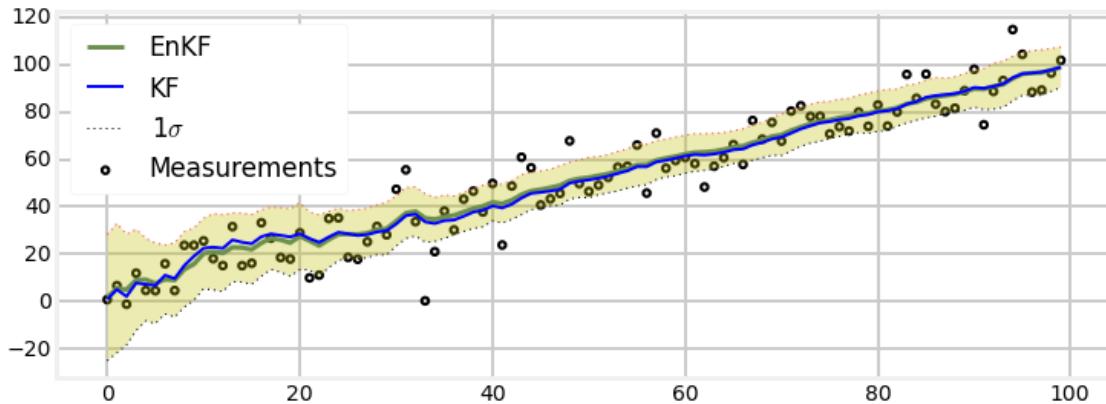
    kf.predict()
    kf.update(np.asarray([[z]]))

    # save data
    results.append (enf.x[0])
    kf_results.append (kf.x[0,0])
    measurements.append(z)
    ps.append(3*(enf.P[0,0]**.5))

results = np.asarray(results)
ps = np.asarray(ps)

plt.plot(results, label='EnKF')
plt.plot(kf_results, label='KF', c='b', lw=2)
bp.plot_measurements(measurements)
plt.plot (results - ps, c='k', linestyle=':', lw=1, label='1$\sigma$')
plt.plot(results + ps, c='k', linestyle=':', lw=1)
plt.fill_between(range(100), results - ps, results + ps, facecolor='y', alpha=.3)
plt.legend(loc='best');

```



It can be a bit difficult to see, but the KF and EnKF start off slightly different, but soon converge to producing nearly the same values. The EnKF is a suboptimal filter, so it will not produce the optimal solution that the KF produces. However, I deliberately chose  $N$  to be quite small (20) to guarantee that the EnKF output is quite suboptimal. If I chose a more reasonable number such as 2000 you would be unable to see the difference between the two filter outputs on this graph.

### D.3 Outstanding Questions

All of this should be considered as my questions, not lingering questions in the literature. However, I am copying equations directly from well known sources in the field, and they do not address the discrepancies.

First, in Brown [2] we have all sums multiplied by  $\frac{1}{N}$ , as in

$$\hat{x} = \frac{1}{N} \sum_{i=1}^N \chi_k^{(i)}$$

The same equation in Crassidis [3] reads (I'll use the same notation as in Brown, although Crassidis' is different)

$$\hat{x} = \frac{1}{N-1} \sum_{i=1}^N \chi_k^{(i)}$$

The same is true in both sources for the sums in the computation for the covariances. Crassidis, in the context of talking about the filter's covariance, states that  $N - 1$  is used to ensure an unbiased estimate. Given the following standard equations for the mean and standard deviation (p.2 of Crassidis), this makes sense for the covariance.

$$\begin{aligned}\mu &= \frac{1}{N} \sum_{i=1}^N [\tilde{z}(t_i) - \hat{z}(t_i)] \\ \sigma^2 &= \frac{1}{N-1} \sum_{i=1}^N \{[\tilde{z}(t_i) - \hat{z}(t_i)] - \mu\}^2\end{aligned}$$

However, I see no justification or reason to use  $N - 1$  to compute the mean. If I use  $N$  in the filter for the mean the filter does not converge and the state essentially follows the measurements without any filtering. However, I do see a reason to use it for the covariance as in Crassidis, in contrast to Brown. Again, I support my decision empirically -  $N - 1$  works in the implementation of the filter,  $N$  does not.

My second question relates to the use of the  $\mathbf{R}$  matrix. In Brown  $\mathbf{R}$  is added to  $\mathbf{P}_{zz}$  whereas it isn't in Crassidis and other sources. I have read on the web notes by other implementers that adding  $\mathbf{R}$  helps the filter, and it certainly seems reasonable and necessary to me, so this is what I do.

My third question relates to the computation of the covariance  $\mathbf{P}$ . Again, we have different equations in Crassidis and Brown. I have chosen the implementation given in Brown as it seems to give me the behavior that I expect (convergence of  $\mathbf{P}$  over time) and it closely compares to the form in the linear KF. In contrast I find the equations in Crassidis do not seem to converge much.

My fourth question relates to the state estimate update. In Brown we have

$$\boldsymbol{\chi} = \boldsymbol{\chi} + \mathbf{K}[\mathbf{z} - \mathbf{z}_{mean} + \mathbf{v}_R]$$

whereas in Crassidis we have

$$\boldsymbol{\chi} = \boldsymbol{\chi} + \mathbf{K}[\mathbf{z} - \boldsymbol{\chi}_h + \mathbf{v}_R]$$

To me the Crassidis equation seems logical, and it produces a filter that performs like the linear KF for linear problems, so that is the formulation that I have chosen.

I am not comfortable saying either book is wrong; it is quite possible that I missed some point that makes each set of equations work. I can say that when I implemented them as written I did not get a filter that worked.

I define “work” as performs essentially the same as the linear KF for linear problems. Between reading implementation notes on the web and reasoning about various issues I have chosen the implementation in this chapter, which does in fact seem to work correctly. I have yet to explore the significant amount of original literature that will likely definitively explain the discrepancies. I would like to leave this here in some form even if I do find an explanation that reconciles the various differences, as if I got confused by these books than probably others will as well.

## D.4 References

- [1] Mackenzie, Dana. Ensemble Kalman Filters Bring Weather Models Up to Date Siam News, Volume 36, Number 8, October 2003. <http://www.siam.org/pdf/news/362.pdf>
- [2] Brown, Robert Grover, and Patrick Y.C. Hwang. Introduction to Random Signals and Applied Kalman Filtering, With MATLAB® excercises and solutions. Wiley, 2012.
- [3] Crassidis, John L., and John L. Junkins. Optimal estimation of dynamic systems. CRC press, 2011.