

PHYSICAL SENSING AND PHYSICS-BASED MACHINE LEARNING FOR
ACTIONABLE INSIGHTS

by

John Waczak

APPROVED BY SUPERVISORY COMMITTEE:

David J. Lary, Chair

Christopher Simmons

David Lumley

Lindsay King

Joseph Izen

Copyright © 2024

John Waczak

All rights reserved

UPDATE REQUIRED!

PHYSICAL SENSING AND PHYSICS-BASED MACHINE LEARNING FOR
ACTIONABLE INSIGHTS

by

JOHN WACZAK, BS, PhD

DISSERTATION

Presented to the Faculty of
The University of Texas at Dallas
in Partial Fulfillment
of the Requirements
for the Degree of

DOCTOR OF PHILOSOPHY IN
PHYSICS

THE UNIVERSITY OF TEXAS AT DALLAS

September 2024

ACKNOWLEDGMENTS

UPDATE REQUIRED! Make sure to include lab colleagues *and* ActivePure colleagues and a nod to Dr. Cooper, OSU teachers, and friends...

August 2024

PHYSICAL SENSING AND PHYSICS-BASED MACHINE LEARNING FOR
ACTIONABLE INSIGHTS

John Waczak, PhD
The University of Texas at Dallas, 2024

Supervising Professor: David J. Lary, Chair

UPDATE REQUIRED!

TABLE OF CONTENTS

ACKNOWLEDGMENTS	v
ABSTRACT	vi
LIST OF FIGURES	xii
LIST OF TABLES	xiii
CHAPTER 1 INTRODUCTION	1
1.1 Dissertation Goals	2
1.2 Global Change	4
1.3 The Role of Sensing	7
1.4 The Role of Computational Modelling	10
1.5 Key Technologies	11
1.5.1 Julia for Scientific Computing	11
1.5.2 Scientific and Physics-based Machine Learning	13
CHAPTER 2 PHYSICAL CONTEXT	16
2.1 Water Quality	16
2.1.1 Properties of Aqueous Solutions	16
2.1.2 Reflectance Spectroscopy	16
2.1.3 Solar Geometry	16
2.2 Air Quality	16
2.2.1 Pollution and Particulate Matter	16
2.2.2 Indoor Air Quality	16
2.3 Physics of Chemical Reactions: Chemical Reaction Kinetics	16
CHAPTER 3 PHYSICAL SENSING	17
3.1 Hyperspectral Imaging	17
3.1.1 Hyperspectral Imaging Sensors	17
3.1.2 Spectral-Spatial Data Collection and Organization	17
3.1.3 Spatial Sampling	18
3.1.4 Spectral Sampling	19
3.1.5 Radiometric Sampling	19

3.1.6	Signal Considerations	20
3.2	Remote Sensing	20
3.2.1	Infrared Sensing Phenomenology	20
3.2.2	Sources of Infrared Radiation	21
3.2.3	Atmospheric Propagation	21
3.2.4	Reflectance and Emissivity Spectra	21
3.3	Coordinated Robotic Teams	22
3.4	In-situ Chemical Sensing in Water	23
3.5	Low Cost Sensors for (Outdoor) Air Quality	23
3.6	The HEART Chamber	23
CHAPTER 4	COMPUTATIONAL TOOLS	24
CHAPTER 5	THEORETICAL TOOLS	25
5.1	Automatic Differentiation	25
5.1.1	Forward Mode AD	25
5.1.2	Reverse Mode AD	25
5.2	Embedding Theorems	25
5.3	Koopman Theory	25
5.4	Bayesian Statistics	25
5.5	Maximum Likelihood Estimation	26
5.6	KL-Divergence	26
5.7	Mathematical Structures	26
5.8	Uncertainty Propagation	26
5.9	Kernelization Methods	26
5.10	Dynamical Systems	26
5.10.1	Chaos???	26
5.11	Endmember Modeling of Reflectance Spectra	26
5.12	Eigen-stuff and the Singular Value Decomposition	27
CHAPTER 6	MACHINE LEARNING METHODS	28
6.1	Exploratory Data Analysis	29

6.1.1	Correlation	29
6.1.2	Mutual Information	29
6.2	Model Training Methodology	29
6.2.1	Data Sampling	29
6.2.2	Model Evaluation	29
6.2.3	Hyperparameter Selection	29
6.2.4	Occam’s Razor and Feature Importance Ranking	30
6.3	Supervised Regression Techniques	30
6.3.1	Neural Networks	30
6.3.2	Gaussian Process Regression	30
6.3.3	Decision Trees	50
6.4	Unsupervised Classification	50
6.4.1	k-Means Clustering	50
6.4.2	Self Organizing Maps	50
6.4.3	Generative Topographic Maps	53
6.5	Dimensionality Redution	53
6.5.1	PCA	53
6.5.2	ICA	53
6.5.3	t-SNE	53
6.6	Model Ensembling	53
6.6.1	Bagging	53
6.6.2	Boosting	53
6.6.3	Stacking	53
6.7	Uncertainty Quantification	53
6.7.1	Quantile Regression	54
6.7.2	Probabilistic Models	54
6.7.3	Conformal Prediction	55
6.8	Scientific Machine Learning	55
6.8.1	Physics-Informed Neural Networks	55

6.8.2	Universal Differential Equations	55
6.8.3	Adjoint Methods for ODEs: Neural Ordinary Differential Equations .	55
6.8.4	Hamiltonian Neural Networks	55
6.9	Generative Methods	55
6.9.1	Auto Encoders	55
6.9.2	Generative Adversarial Networks	55
6.10	Data Assimilation	55
6.10.1	Kalman Filter	58
6.10.2	Extended Kalman Filter	64
6.10.3	continuous-discrete Extended Kalman Filter	65
6.10.4	3d-Var	65
6.10.5	4d-Var	68
CHAPTER 7 AN AUTONOMOUS ROBOTIC TEAM FOR THE RAPID CHARAC-		
TERIZATION OF NOVEL ENVIRONMENTS		70
7.1	Rapid Georectification and Processing of Pushbroom Hyperspectral Imagery	70
7.2	Supervised Regression with Uncertainty Quantification	70
7.3	Methods for Unsupervised Classification of Hyperspectral Scenes in Novel Environments	70
CHAPTER 8 A CHEMICAL DATA ASSIMILATION FRAMEWORK FOR INDOOR		
AIR QUALITY		71
8.1	Characterization of Photolysis	71
8.2	HEART Chamber Sensing System	71
8.3	Chemical Data Assimilation	71
8.4	An Evaluation of Photocatalytic Ionization	71
CHAPTER 9 A DISTRIBUTED NETWORK OF LOW COST AIR QUALITY SEN-		
SORS		72
9.1	MINTS Air Quality Network	72
9.1.1	Central Nodes	72
9.1.2	LoRa Nodes	72
9.1.3	MQTT	72

9.2	Real Time Dashboards via Containerization	72
9.3	Making Data Publicly Accessible via S3 and the Open Storage Network . . .	72
CHAPTER 10 TIME SERIES METHODS FOR AIR QUALITY		73
10.1	Time-Series Methods for Uncertainty Quantification	73
10.1.1	Metrics for Representative Uncertainty of combined Pseudo-Observations	73
10.1.2	Uncertainty Quantification with Temporal Variograms	73
10.2	Physics Informed modeling techniques for Air Quality Data	73
10.2.1	HAVOK + SciML = Nonlinear Modeling with External Forcing . . .	73
10.2.2	Hamiltonian Neural Networks	73
CHAPTER 11 LIMITATIONS AND FUTURE WORK		74
11.1	Robot Team	74
11.1.1	Super Resolution	74
CHAPTER 12 CONCLUSIONS		75
APPENDIX A REPRODUCIBLE RESEARCH TECHNIQUES		76
A.1	Environment Management in Julia	76
A.2	Version Control (git)	76
A.3	CI/CD with github workflows	76
A.4	Literate Programming and Automatic Documentation with Quarto	76
A.5	Containerization with Docker and Docker Compose	76
APPENDIX B OPTIMIZATION METHODS		77
APPENDIX C HIGH PERFORMANCE COMPUTING		78
APPENDIX D CHEMICAL REACTION MECHANISMS		79
D.1	Simple Ion Mechanism	79
D.2	Master Chemical Mechanism	79
REFERENCES		80
BIOGRAPHICAL SKETCH		84
CURRICULUM VITAE		

LIST OF FIGURES

6.1	A standard linear regression fit on some noisy data.	32
6.2	An example of polynomial regression for which we fit a quadratic polynomial to some noisy data.	37
6.3	A Gaussian Process fit to the training data illustrating the predication (means) and a $\pm 2\sigma$ uncertainty interval. We can clearly see how the uncertainty is larger for inputs far away from supplied training data.	47
6.4	Left: the original Gaussian Process obtained with our default choice of kernel parameters. Right: A much better Gaussian Process obtained after performing hyper-parameter optimization on the kernel parameters. Note how the mean function still does an excellent job fitting the data points while <i>also</i> minimizing the uncertainty of the fit.	50

LIST OF TABLES

CHAPTER 1

INTRODUCTION

In this “global” scale introduction, we should introduce the context of global change and the need for improved sensing and modeling to provide actionable insights at a pace that meets societal/human needs. We can then discuss how big data + machine learning can help fill in the gaps where our theoretical knowledge is incomplete, expensive (money and computational) to simulate directly, or plagued by initial condition sensitivity (e.g. Chaos).

Finish with a transition paragraph providing an overview of each of the following chapters
For Water Quality:

- chemical quantification (crude oil, algal blooms, terrorist events, etc...)
- chemical identification (can we identify new constituents in the water?)

Outdoor Air Quality:

- How can we model the uncertainty of low cost sensors for real-time applications. This will have applications on real-time decision making, QA/QC, etc...
- Can we leverage data to effectively model the dynamics of local air quality without the need for complicated microphysics simulations? (We can say something here about the need to go beyond thermodynamic equilibrium)
- Can we learned models provide new insights into real physics? E.g. what do the learned terms of our SciML extended HAVOK model tell us? How can we interpret the forcing function? What do the learned coordinates of the HNN represent? How can we interpret dynamics on the Hamiltonian surface?

Indoor Air Quality:

- Can we leverage cutting edge measurement techniques to effectively model indoor chemical kinetics?
- What is the role of ion chemistry in indoor air quality?
- By combining observed species concentrations with highly detailed kinetics, can we infer the presence and role of species below detectable limits?

For the proposal, let's end the introduction with a timeline (we can use Dr. Lary's Gantt chart).

1.1 Dissertation Goals

The goal of this thesis is advancing physical sensing in service of society to provide actionable insights. This goal is pursued by applying physics informed approaches together with a suite of sensing and computational technologies, implementing the reusable paradigm of software defined sensors, i.e. physical sensing elements wrapped in a software layer. This software layer can serve a variety of purposes such as calibration and the provision of enhanced or derived data products. It is part of a broader effort in the MINTS-AI laboratory at the University of Texas at Dallas. Where MINTS-AI is an acronym, Multi-Scale Multi-Use Integrated Intelligent Interactive Sensing in Service of Society for Actionable Insights.

Comprehensive environmental sensing is a timely and beneficial endeavor for a variety of reasons. The growing awareness of major environmental issues such as climate change, pollution, and habitat loss necessitates effective environmental monitoring and management. Comprehensive environmental sensing can provide real-time data on air and water quality, weather patterns, and other environmental factors, assisting in the identification and resolution of environmental issues. This assists in the development and implementation of policies and strategies aimed at reducing environmental impact and increasing sustainability. Given

that, for instance, air quality can have significant effects on human health, this has particular societal value.

Exposure to air pollution has been linked to a wide range of health effects (Brook et al., 2010; Kelly and Fussell, 2011; Xu et al., 2017), including respiratory and cardiovascular diseases, cancer, and adverse birth outcomes. Further, physical sensing provides valuable data and the basis for international assessments such as the Intergovernmental Panel on Climate Change (IPCC), which seeks to assess the science related to climate change and its impacts on natural and human systems (Houghton et al., 1990, 1996, 2001; Solomon et al., 2007; Parry et al., 2007; Metz et al., 2007; Stocker et al., 2013; Field et al., 2014; Edenhofer et al., 2014; Masson-Delmotte et al., 2018; Friedlingstein et al., 2020; Huang et al., 2017).

Comprehensive sensing of the environment can improve decision-making. The real-time and accurate data provided by environmental sensors can aid in informed decision-making regarding various aspects such as traffic management, industrial regulation, and crop planning. For instance, data on air quality can be used to inform decisions about reducing pollution levels, while data on weather patterns can help farmers to plan their crops and reduce water usage. Comprehensive sensing of the environment can be instrumental in emergency response. Real-time data on weather patterns, air quality, water levels and resources, and seismic activity can help emergency responders to prepare for and respond to natural disasters such as hurricanes, floods, and earthquakes. The quick and accurate information can enable effective and timely response, potentially saving lives and reducing the impact of the disaster.

Many advances in technology have enabled the creation of comprehensive sensing systems that can monitor and analyze data from various sensors and devices in real-time. In this thesis we use a range of technologies including autonomous robotic teams [Dunbabin2012; Rubenstein2014; Chen2017], hyperspectral imaging [Plaza2009; Li2018; Zhu2017], mesh networks utilizing the Internet of Things (IoT) [Gubbi2013; Atzori2010; Al-Fuqaha2015], machine learning (ML) [Goodfellow2016; LeCun2015; Jordan2015], edge

computing, high-performance computing, wearable sensors and modern high-performance dynamic programming languages such as Julia [Bezanson2017] designed for numerical and scientific computing. These technologies have facilitated the collection and processing of large amounts of data from multiple sources, resulting in more accurate and comprehensive environmental monitoring.

1.2 Global Change

Global change refers to the significant and long-term alterations in the Earth’s physical, chemical, and biological systems, resulting from natural and human-induced processes (Edenhofer et al., 2014; Masson-Delmotte et al., 2018; United Nations, 2015). This includes changes in the climate, land use, biodiversity, and biogeochemical cycles, as well as interactions among these systems. Global change can have profound impacts on natural and human systems, including altered weather patterns, sea level rise, increased frequency and severity of extreme events, loss of biodiversity and ecosystem services, and effects on human health and well-being. Understanding and managing global change is a critical challenge facing society today, requiring interdisciplinary approaches and collaboration across sectors and regions.

Global change can have a range of impacts on society, including environmental, social, and economic effects. Some of the aspects of global change that have the biggest impact on society include:

- **Climate Change:** Climate change, driven by human activities such as burning fossil fuels, deforestation, and land-use changes, has impacts on natural systems such as ocean acidification, sea level rise, and changes in precipitation patterns. These impacts can have cascading effects on human systems, including impacts on food security, water availability, and health.

- **Biodiversity Loss:** Global change can lead to the loss of biodiversity, which can have impacts on ecosystem functioning and services, such as pollination, pest control, and carbon storage. These impacts can have indirect effects on human well-being, including impacts on food security, health, and cultural heritage.
- **Land Use Change:** Land use change, such as deforestation, urbanization, and agriculture, can have impacts on natural systems such as soil quality, water availability, and biodiversity. These impacts can have direct and indirect effects on human systems, including impacts on food security, water availability, and cultural heritage.
- **Economic and Social Inequality:** Global change can exacerbate economic and social inequality, with impacts on access to resources, health, and well-being. These impacts can have cascading effects on the ability of societies to adapt and respond to global change.
- **Human Health:** Global change can have significant impacts on human health (World Health Organization, 2018; Costello et al., 2009; Haines et al., 2006), both directly and indirectly, for example:
 - **Heat-related Illness:** As temperatures increase due to global warming, there is an increased risk of heat-related illness, including heat exhaustion and heat stroke, particularly in vulnerable populations such as the elderly, young children, and outdoor workers.
 - **Air Pollution:** Global change can lead to increased air pollution, including from sources such as wildfires and fossil fuel combustion. Exposure to air pollution can increase the risk of respiratory and cardiovascular diseases, including asthma, chronic obstructive pulmonary disease (COPD), and heart disease.

- Vector-borne Diseases: Changes in temperature and precipitation patterns can affect the distribution and abundance of disease vectors such as mosquitoes and ticks, leading to increased risks of vector-borne diseases such as dengue fever, malaria, and Lyme disease.
- Waterborne Diseases: Changes in precipitation patterns and water quality can increase the risk of waterborne diseases, including cholera and other diarrheal diseases.
- Food Security: Global change can affect food production and availability, leading to food shortages and malnutrition, particularly in vulnerable populations such as children and pregnant women.

Effectively addressing these aspects of global change requires interdisciplinary approaches and collaboration across sectors and regions, as well as a commitment to sustainable development and equitable solutions. Adaptation and mitigation are two strategies for addressing global change, which differ in their focus and approach.

Adaptation involves taking measures to adjust and respond to the impacts of global change that are already occurring or are expected to occur in the future. This can include actions such as building sea walls to protect against sea level rise, developing drought-resistant crops, and improving public health infrastructure to address the increased risk of vector-borne diseases. Adaptation strategies aim to reduce the vulnerability of human and natural systems to the impacts of global change and increase their resilience.

Mitigation involves taking measures to reduce the drivers of global change, such as greenhouse gas emissions, land use change, and deforestation. This can include actions such as increasing energy efficiency, shifting to renewable energy sources, and reducing waste and consumption. Mitigation strategies aim to address the root causes of global change and reduce its severity and impact.

Both adaptation and mitigation are important strategies for addressing global change, but they differ in their focus and approach. Adaptation strategies focus on responding to the impacts of global change that are already occurring or are expected to occur in the future, while mitigation strategies focus on reducing the drivers of global change and preventing its impacts from occurring in the first place. A comprehensive approach to global change will require both adaptation and mitigation strategies, as well as efforts to promote sustainable development and equitable solutions.

1.3 The Role of Sensing

Sensing technologies can play a critical role in both adaptation and mitigation efforts by providing data and information that can inform decision-making and improve the effectiveness of strategies (United Nations Environment Programme, 2017; National Research Council, 2010; Centre for Ecology and Hydrology, 2017).

In adaptation efforts, sensing technologies can provide real-time data on environmental conditions such as temperature, precipitation, sea level, air quality, as well as on the status and health of ecosystems and wildlife. This information can be used to inform early warning systems for natural disasters, to track the spread of vector-borne diseases, and to monitor the impacts of climate change on biodiversity and ecosystem services. Sensing technologies can also provide data on the effectiveness of adaptation measures, such as the performance of sea walls and other infrastructure.

In mitigation efforts, sensing technologies can provide data on greenhouse gas emissions and other drivers of global change, as well as on the effectiveness of mitigation measures such as renewable energy and carbon capture and storage. Sensing technologies can also be used to monitor and manage land use changes such as deforestation and urbanization, and to track the impacts of these changes on ecosystems and carbon storage.

Overall, sensing technologies can provide critical data and information for both adaptation and mitigation efforts, helping to improve decision-making and increase the effectiveness of strategies. The integration of sensing technologies with other tools such as modeling and data analysis can also help to identify new strategies and solutions for addressing global change. There are various sensing technologies and approaches used for monitoring the global environment. Here are some of the key ones:

- Remote Sensing: This technology involves using satellites and other airborne platforms to collect data on the Earth's atmosphere, land, and oceans. Remote sensing provides information on environmental parameters such as temperature, humidity, air quality, land use and land cover, and ocean temperature, salinity, and sea level (Thenkabail, 2019; Buyantuyev and Wu, 2017; Gamon et al., 2016; Wang et al., 2017; Pasher et al., 2019). Some examples of remote sensing include:
 - Lidar: This technology uses laser pulses to measure distance and can be used to create detailed three-dimensional maps of the environment. Lidar is commonly used to measure forest canopy height, but can also be used to measure atmospheric conditions such as cloud cover and aerosol concentrations.
 - Imaging Spectroscopy: This technology uses a combination of imaging and spectroscopy to measure the reflectance of different wavelengths of light. Imaging spectroscopy can be used to identify and map different types of vegetation and minerals, and can provide information on the health of plant communities.
 - Unmanned Aerial Vehicles (UAVs): These are remote-controlled or autonomous aircraft that can be equipped with sensors for remote and in-situ environmental monitoring. UAVs can be used for mapping and monitoring of large areas, and can collect high-resolution data on environmental conditions.

- In-Situ Sensors: These sensors are used to collect data directly from the environment at the location of interest. They can measure environmental parameters such as temperature, pressure, and humidity, as well as water quality and soil moisture. In situ sensors are commonly used in marine environments to measure ocean temperature, salinity, and other properties. Some examples of in-situ sensing include:
 - Weather Stations: These are automated weather monitoring systems that collect data on atmospheric conditions such as temperature, humidity, barometric pressure, wind speed and direction, and precipitation. Weather stations can be installed on land or in the ocean to provide continuous monitoring of environmental conditions.
 - Ground-Based Sensors: These sensors are used to monitor the quality of air, water, and soil. They can detect and measure pollutants such as carbon dioxide, nitrogen dioxide, ozone, sulfur dioxide, and particulate matter. Ground-based sensors are installed in various locations such as cities, industrial sites, and rural areas to provide localized environmental monitoring.
 - Acoustic Sensors: These sensors are used to monitor environmental noise levels, including noise from traffic, industrial sources, and natural sources such as wind and waves. Acoustic sensors can provide information on noise levels over time and across different locations.

Overall, these sensing technologies play a critical role in monitoring the global environment and can provide valuable information for environmental research, management, and policy-making.

1.4 The Role of Computational Modelling

Computer modeling can play a valuable role in both understanding and predicting global change (Chen et al., 2019; Hantson et al., 2016; DeLucia et al., 2021; Oleson et al., 2013; Clark et al., 2016). For example:

- **Climate Modeling:** Computer models can be used to simulate the Earth’s climate system and predict future climate conditions. These models can incorporate data on greenhouse gas emissions, land use changes, and other factors to project how the Earth’s climate will change over time.
- **Ecosystem Modeling:** Computer models can be used to simulate how ecosystems will respond to changes in environmental conditions, such as changes in temperature, precipitation, and atmospheric composition. These models can help predict how changes in ecosystems will impact biodiversity, ecosystem services, and human well-being.
- **Carbon Cycle Modeling:** Computer models can be used to simulate the global carbon cycle, which is the exchange of carbon between the Earth’s atmosphere, land, and oceans. These models can help predict how changes in carbon emissions and land use will impact atmospheric carbon dioxide concentrations and global climate.
- **Air Quality Modeling:** Computer models can be used to simulate air quality, including the dispersion of pollutants in the atmosphere. These models can help predict how changes in emissions and atmospheric conditions will impact air quality and human health.
- **Hydrological Modeling:** Computer models can be used to simulate the movement of water through the Earth’s hydrological cycle. These models can help predict how changes in precipitation, land use, and other factors will impact water availability, quality, and distribution.

Overall, computer modeling can provide valuable insights into the complex processes and interactions that drive global change. These insights can inform policy decisions and help guide efforts to mitigate and adapt to the impacts of global change.

1.5 Key Technologies

1.5.1 Julia for Scientific Computing

Julia is designed to combine the ease of use and high-level abstractions of languages like Python with the performance of compiled languages like C++, achieving a unique combination of speed and productivity for numerical and scientific computing. Julia is a high-level, high-performance programming language designed for numerical and scientific computing. It combines the ease of use and readability of Python with the speed and efficiency of Fortran or C. Julia has a wide array of scientific computing functionality, making it a powerful language for numerical analysis, data science, and engineering. It has built-in support for arrays and linear algebra, as well as packages for differential equations, optimization, probabilistic programming, data analysis and visualization, parallel and distributed computing, and machine learning. Julia's combination of performance, expressiveness, and flexibility make it an excellent choice for scientific and engineering applications, allowing for high-level abstractions and rapid prototyping, while still providing low-level control and efficient execution.

Here are some examples of what can be done easily in Julia that may not be as easy or efficient in other widely used scientific computing languages such as Python or Fortran:

- **Multiple dispatch:** Julia has a powerful multiple dispatch system that allows for generic programming and efficient function overloading. This allows for more flexible and expressive code compared to traditional object-oriented programming (OOP) in Python. Multiple dispatch allows a function to behave differently based on the types and/or number of arguments passed to it. In other words, the behavior of a function can be dispatched based on the specific types and/or number of arguments passed to it.

- Just-in-time (JIT) compilation: Julia's JIT compiler translates high-level Julia code into optimized machine code, making Julia programs run nearly as fast as C or Fortran. In contrast, Python code is interpreted, and Fortran requires pre-compilation.
- Distributed computing: Julia has built-in support for distributed computing, making it easy to parallelize and scale up computations across multiple processors or machines. This is not as easy to do in Python or Fortran.
- Units and Error Propagation: The Units package in Julia provides a powerful and flexible framework for handling physical units in computations, useful for error propagation and dimensional analysis, helping to ensure that the results are accurate, consistent, easy to interpret, and dimensionally consistent.
- Built-in unit testing: Julia has a built-in testing framework that makes it easy to write and run unit tests for your code, ensuring that it works correctly.
- ISO characters: Julia supports the use of Greek and other ISO characters in variable and function names, which can make code more readable and expressive, especially in mathematical or scientific contexts.
- Interactive data visualization: Julia has a number of powerful data visualization packages, such as Plots.jl and Makie.jl, that allow for interactive, high-performance data visualization.
- Package management: Julia has a sophisticated package manager that makes it easy to install, manage, and use third-party packages in your code. This is not as easy to do in Fortran, and while Python has a package manager, Julia's package manager is faster and more reliable.

- Inline C/Fortran/Python/R/Matlab code: Julia allows for inline C, Fortran, Python, R or Matlab code, making it easy to use existing libraries and code written in these languages without having to rewrite everything in Julia.

1.5.2 Scientific and Physics-based Machine Learning

Scientific machine learning (SciML) refers to the application of Machine Learning (ML) techniques to scientific problems, where the goal is not only to make predictions but also to gain insights into the underlying physical processes (Raissi et al., 2019; Rackauckas et al., 2020; Carleo et al., 2019). SciML involves the integration of domain-specific knowledge and physical models with data-driven techniques, and it has the potential to revolutionize many areas of science and engineering. In this thesis we explore the use of Physics-based machine learning (PBML) (Raissi and Karniadakis, 2021; Wu and Zhang, 2021) for a variety of applications.

Recent examples include a paper by (Raissi et al., 2019) that introduces a physics-informed neural network (PINN) framework for solving nonlinear partial differential equations, a paper by (Rackauckas et al., 2020) that proposes a universal differential equation (UDE) approach to scientific machine learning, and a review article by (Carleo et al., 2019) that discusses the use of machine learning in various fields of physics, including condensed matter physics, high-energy physics, and quantum physics. PBML has several advantages over purely data-driven approaches, including:

- Improved generalization: PBML models incorporate prior knowledge of the underlying physics, resulting in models that are more interpretable and generalizable. This enables the models to make accurate predictions even with limited training data.
- Incorporation of physical constraints: PBML models can be designed to incorporate physical constraints, such as conservation laws, which can help to ensure physically consistent predictions.

- Improved interpretability: PBML models are more interpretable than purely data-driven models since they are designed to incorporate physical principles. This can enable scientists and engineers to gain deeper insights into the underlying mechanisms of the systems they are studying.
- Reduced data requirements: PBML models require less training data than purely data-driven models since they leverage the physics-based priors, reducing the need for large datasets to train accurate models.
- Better extrapolation: PBML models are better equipped to extrapolate beyond the training data since they incorporate knowledge of the underlying physics, enabling them to make more accurate predictions in new and unseen scenarios.

Overall, PBML has several advantages over purely data-driven approaches, including improved generalization, reduced data requirements, better extrapolation, incorporation of physical constraints, and improved interpretability, making it a valuable tool for scientific and engineering applications.

NOTE: we may want to merge the Introduction with the Physical Context chapter.

CHAPTER 2

PHYSICAL CONTEXT

2.1 Water Quality

2.1.1 Properties of Aqueous Solutions

2.1.2 Reflectance Spectroscopy

2.1.3 Solar Geometry

An explanation of relevant solar angles as well as their determination (i.e. the code I ported to Julia from Matlab script Dr. Lary supplied). We should also comment on the importance of

Use David's figure from his book to illustrate the solar geometry.

2.2 Air Quality

2.2.1 Pollution and Particulate Matter

discuss sources, evolution, global trends, etc...

2.2.2 Indoor Air Quality

2.3 Physics of Chemical Reactions: Chemical Reaction Kinetics

CHAPTER 3

PHYSICAL SENSING

3.1 Hyperspectral Imaging

The following are notes from the Manolakis textbook that I originally kept here. NOTE: we will need to either make new figures or correctly cite these for attribution.

3.1.1 Hyperspectral Imaging Sensors

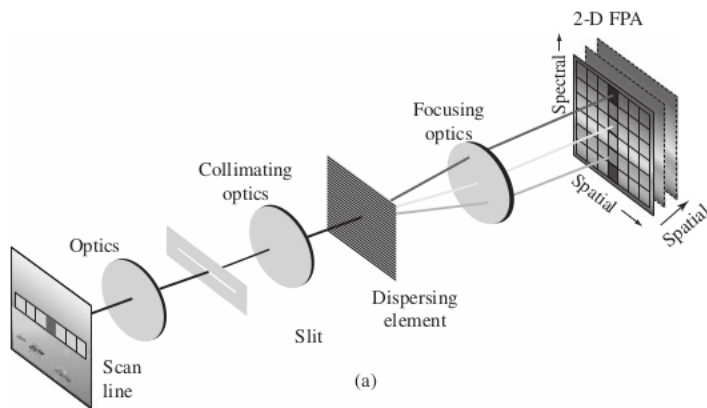
- *Hyperspectral Sensors* aka imaging spectrometers
 - scanning mechanism
 - imaging system
 - spectrometer
- 3 types of resolution
 1. spatial
 2. spectral
 3. radiant
 4. (temporal?)

3.1.2 Spectral-Spatial Data Collection and Organization

Data collected into *datacube* with 2 spatial dimensions, 1 spectral dimension

Different types of rigs:

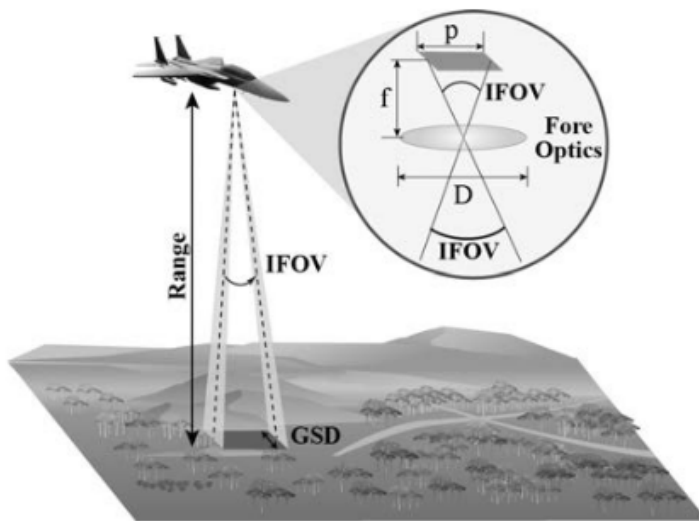
- Pushbroom scanner (ours)
- Staring System
- Fourier Transform Imaging Spectrometer (FTIS)



3.1.3 Spatial Sampling

- ground resolution elements are mapped to picture elements (pixels)
- IFOV: Instantaneous Field of View

- Cross track dimension: the projection of the long axis of the slit (i.e. the axis of the pushbroom sensors)
- Along track dimension: the direction accumulated by traveling
- Ground Sample Distance: physical size of projected pixel element



3.1.4 Spectral Sampling

- Recovery of spectral info is imperfect due to finite sampling
- Spectral Response Function: is the weighing function that describes the wavelengths that are transmitted to a particular spectral sample

3.1.5 Radiometric Sampling

- detector transforms radiant power to electrical signal
- electrical signal converted to number via analog-to-digital converter
- photon detectors

3.1.6 Signal Considerations

Strength of signal is determined by:

- Terrain composition affects amount of radiant energy reflected/emitted from ground resolution element
- Range: Intensity drops off by inverse square law. Further you are away, the worse the signal
- Spectral Bandwidth: output signal of detector element is proportional to spectral bandwidth of the detector
- Instantaneous Field of View: Decreasing IFOV increases spatial resolution but weakens the signal
- Dwell Time: the time required to sweep the IFOV across the ground resolution element, i.e. the time-on-pixel. Longer dwell time \rightarrow more accumulated photons \rightarrow more signal.

3.2 Remote Sensing

mention different types of satellite data platforms (mostly optical), differences in orbits, coverage, etc... Also good to discuss the increasing use of drones for a variety of applications including intelligent agriculture, geophysics, mapping, etc...

3.2.1 Infrared Sensing Phenomenology

Main passive sources of EM radiation for remote sensing are light emitted by the sun and the self-emission via black-body radiation of objects due to their temperature.

3.2.2 Sources of Infrared Radiation

- **spectral radiant exitance** power per unit area emitted by the sun. We can treat this as a black body with temperature $5800K$, maximum emittance at $\lambda = 0.50 \mu m$.
- The Earth is $300K$ with maximum spectral radiant emittance at $\lambda = 9.7 \mu m$. This is known as the **thermal infrared**

3.2.3 Atmospheric Propagation

- Key parameter is the **path length** of atmosphere traveled through before it arrives at the remote sensing system. Main effects are:
 - **Atmospheric Scattering**: diffusion of radiation by particles in the atmosphere
 - **Absorption**
- Useful remote sensing spectral regions are obtained via the **Transmission Spectrum**.
 - **Reflective Range**: $0.35 - 2.5 \mu m$. Dominated by solar illumination
 - **Water Absorption**: $0.2 - 2.5 \mu m$.
- **Atmospheric Windows**: Regions of low atmospheric absorption

3.2.4 Reflectance and Emissivity Spectra

There are three processes that occur when EM radiation meets an interface:

1. **Reflection**: Solar illumination dominates here. Consequently, this part of the spectrum is used to characterize the surface
 - **Specular Reflectors**: Flat surfaces that act like mirrors, i.e. $\theta_i = \theta_r$.
 - **Diffuse (Lambertian) Reflectors**: Rough surfaces that reflect uniformly in all directions.

- **Real Reflectors:** Somewhere between the specular and diffuse.

2. Absorption

3. Transmission

- Fractions vary as a function of λ
- Remote sensing usually cares about *diffuse* reflectors because this is the dominant type of most materials (water being an exception).
- Reflectance of a material is characterized by its **Reflectance Spectrum**, that is, the percent of incident light reflected as a function of wavelength.
 - Dips in reflectance spectrum are called **absorption features**
 - Peaks are called **Reflectance Peaks**
- **Emissivity Spectrum:** The ratio of radiant emittance at a given temperature to the radiant emittance of a black body at the same temperature.

3.3 Coordinated Robotic Teams

- Remote Sensing: data acquisition, processing, and interpretation of images, and related data, obtained from aircraft and satellites that record the interaction between matter and electromagnetic radiation
- Source: the source of electromagnetic radiation, e.g. the sun, black-body radiation, microwave radar, etc...
- Atmospheric Radiation: The EM radiation propagating through the atmosphere. Moderated by various processes including absorption and scattering

- Earth's Surface Interactions: Amount and spectral distribution of radiation emitted/reflected by the earth's surface. This depends on
 - physical properties of the matter
 - wavelength of EM radiation that is sensed

3.4 In-situ Chemical Sensing in Water

Here we can give an overview of the sensors utilized on the boat for the robot team studies... We should include any information about the uncertainties. We can also comment on their use in fresh v.s. salt water environments, e.g. we should discuss fluorometers, sonar, and any other relevant sensors from the boat

3.5 Low Cost Sensors for (Outdoor) Air Quality

The outdoor part here is optional. Here I seek to provide specific details of the classes of low cost sensors used for PM (optical particle counters based on Mie scattering), thermistors (for temperature), humidity sensors, pressure sensors, gas sensors, etc... This can be a rough outline for those that specifically are used in my research projects.

3.6 The HEART Chamber

CHAPTER 4

COMPUTATIONAL TOOLS

I'm not yet sure if I want to keep this chapter or move the content to the appendices

CHAPTER 5

THEORETICAL TOOLS

5.1 Automatic Differentiation

perhaps this can be moved to the Computational tools section instead....

Use Chris Rackauckas and Chris Olah blogs to start with chain rule and justify automatic differentiation. Then describe reverse mode via gradient tapes and the pullback operator.

5.1.1 Forward Mode AD

Dual Numbers

Describe dual numbers because they are cool (and by extension, forward mode AD)

5.1.2 Reverse Mode AD

5.2 Embedding Theorems

Discuss Taken's embedding theorem and other relevant information for the time-series work.

5.3 Koopman Theory

Give an overview of continuous and discrete Koopman operator theory.

5.4 Bayesian Statistics

Derive Baye's rule and other important concepts

5.5 Maximum Likelihood Estimation

5.6 KL-Divergence

5.7 Mathematical Structures

A generic overview of role of mathematical structures in physics, machine learning, and mathematical modeling.

5.8 Uncertainty Propagation

Linear v.s. Nonlinear Techniques, measurements.jl

5.9 Kernelization Methods

5.10 Dynamical Systems

can quote paper on origins of term *phase space*

5.10.1 Chaos???

5.11 Endmember Modeling of Reflectance Spectra

We can discuss previous attempts at modeling measured reflectance spectra via linear combinations of endmember spectra. Similarly, we can do a nice discussion of methods for *discovering* the endmembers from the combined spectral. We can also discuss the limitations of linearity assumptions and outline where that clearly fails, i.e. turbid solutions, etc.

5.12 Eigen-stuff and the Singular Value Decomposition

Begin with a general description of the utility of eigen-stuff type analysis. Then explain how the SVD is the natural extension of this to non-square systems (rectangular matrices). Further expand by illustrating the application to principal component analysis. Make sure to comment on the general utility of decomposing a function/vector/signal into a (infinite) linear combination of (stationary) modes with simple time evolution.

CHAPTER 6

MACHINE LEARNING METHODS

NOTE: Good way to start is with the question: What is a model? We can discuss the differences between mechanistic models (i.e. physics equations), their free parameters (constants of nature) and other types of models such as the non-parametric, nonlinear models used in machine learning. Further, we can discuss how machine learning models often display have the desirable trait of being a **universal approximator**. This will naturally lead to a discussion of function expansions (Taylor, Fourier, other polynomial expansions, etc...) and how they scale poorly (exponential) with increasing dimension. Machine learning models essentially allow us to do the same thing: perform a function expansion given data but in a way that can scale well to arbitrary dimension (features) of data. This allows us to build predictive models without necessarily needing to prescribe *all* of the physics. From another perspective, this framework allows us to incorporate our physics knowledge from well-behaved or linearized domains in order to *fit the residual* behavior with a sophisticated data-driven approach.

discuss role in science in particular (i.e. calibration, modeling, etc...)

discuss pushback against use in science and need for methods which simultaneously provide uncertainty bounds. Also describe types of ML e.g. supervised, unsupervised, generative, etc...

this is a good place for the Chihuahua vs Muffin meme... and use this to motivate incorporating physical knowledge into the machine learning process as a key feature for scientific applications... we have more constraints!

also discuss statistical v.s. deep learning

6.1 Exploratory Data Analysis

6.1.1 Correlation

6.1.2 Mutual Information

6.2 Model Training Methodology

6.2.1 Data Sampling

Dr. Lary's method (from Gaussian Process Code) for representative sampling to reduce data size

- Testing holdout set
- k-fold Cross Validation

We should use this subsection to discuss why k-fold cross validation is ideal where data size/time permits so that you can get a sense of each model's dependence on the input data. Similarly, it is important we end the discussion by pointing out that the final **production** model is to be trained on the *entire* dataset after we have evaluated which model is best to use and the ideal hyperparameters, etc...

6.2.2 Model Evaluation

Here we should discuss how we evaluate models, e.g. scatter plots, quantile-quantile plots, histograms, confusion matrices, various correctness metrics, etc...

6.2.3 Hyperparameter Selection

Discuss grid search vs random search vs Bayesian Optimization, etc...

6.2.4 Occam's Razor and Feature Importance Ranking

I should point out my contributions to various julia packages in MLJ for doing feature importance determination here

Linear Regression

Tree Based Methods

Neural Networks

Perhaps explore ideas from Chris Olah's blog??? From a perspective of basic function composition... as in Rackauckas's blog

Shapley Values

Here we can discuss Shapley Values as a model agnostic method for feature importance rankings with the caveat that their computation can take quite a while...

6.3 Supervised Regression Techniques

6.3.1 Neural Networks

6.3.2 Gaussian Process Regression

Based on my notes from this repo

The following is based on the book **Gaussian Processes for Machine Learning** by Carl Edward Rasmussen and Christopher K. I. Williams (Rasmussen et al., 2006). You can find the free online book here.

To explain/derive the Gaussian Process model for regression, let's first consider a motivating example: **Linear Regression**. We will use this guiding example to derive GRP from a *weight space view*. After this derivation, we will suggest a simpler, but more abstract derivation using a *function space view*. First let's set up some data we can use for training:

The Weight Space View

We consider a dataset \mathcal{D} with n observations

$$\mathcal{D} = \left\{ (\mathbf{x}_i, y_i) \mid i = 1, \dots, n \right\} \quad (6.1)$$

- \mathbf{x}_i is the i^{th} D -dimensional input (feature) vector
- y_i is the i^{th} target

Linear regression is easily understood in terms of *linear algebra*. We therefore collect our dataset \mathcal{D} into a $D \times n$ dimensional **Design Matrix**. Note that we have used a transposed definition (features are rows, records are columns) as Julia is a column-major language (like Matlab and Fortran).

$$X := \begin{pmatrix} \vdots & \vdots & & \vdots \\ \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_n \\ \vdots & \vdots & & \vdots \end{pmatrix} \quad (6.2)$$

and our targets into a target vector

$$\mathbf{y} := (y_1, \dots, y_n) \quad (6.3)$$

so that the full training set becomes

$$\mathcal{D} := (X, \mathbf{y}) \quad (6.4)$$

Standard linear regression is a model of the form

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} \quad (6.5)$$

where \mathbf{w} is the D -dimensional vector of weights. By minimizing the mean-squared-error between our model and targets, one can show that the optimal weights are given by

$$\mathbf{w} = (X X^T)^{-1} X \mathbf{y} \quad (6.6)$$

Note, this can also be easily obtained geometrically by finding the vector with the shortest distance to the hyperplane defined by the column space of X . This corresponds to solving the normal equations.

$$X X^T \mathbf{w} = X \mathbf{y} \quad (6.7)$$

The following demonstrates this procedure on a simple dataset

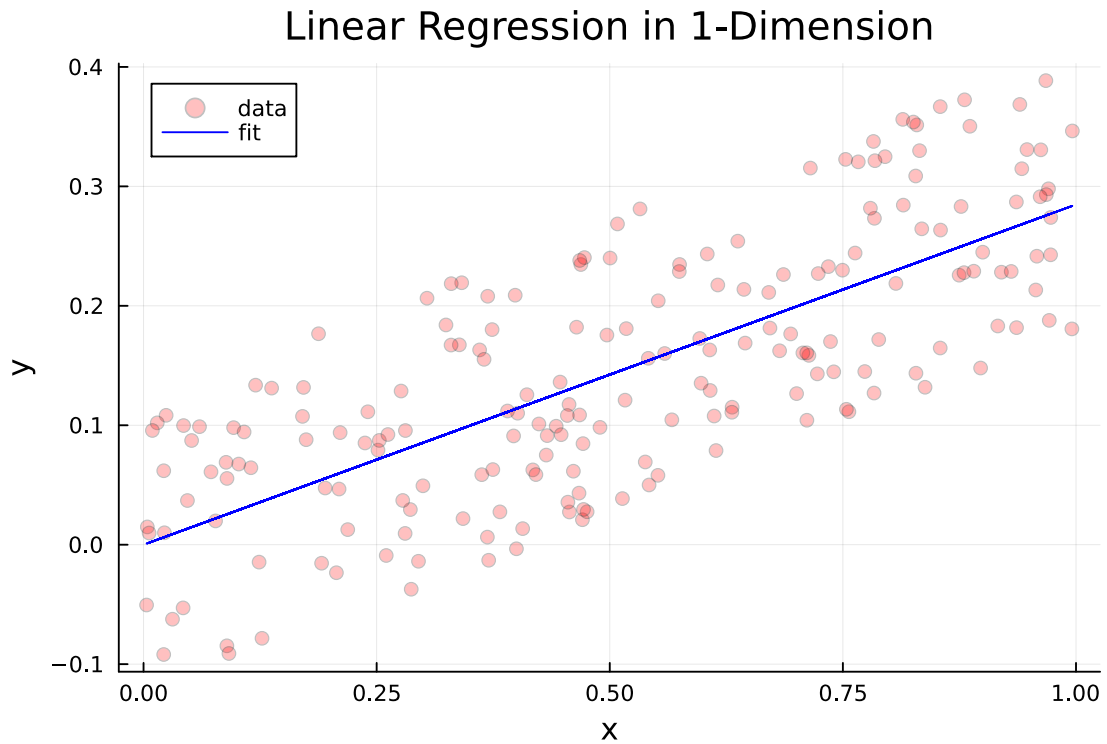


Figure 6.1. A standard linear regression fit on some noisy data.

We can also fit a y-intercept (aka *bias*) by augmenting the design matrix X to contain an extra row with all 1's, i.e.

$$X[D + 1, :] = (1, \dots, 1) \quad (6.8)$$

Standard linear regression assumes that are data \mathcal{D} are perfect but we can clearly see that the above data are noisy. To account for this, we need to make our model *Bayesian* by augmenting it to consider the measurement error. We define

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} \quad (6.9)$$

$$\mathbf{y} = f(\mathbf{x}) + \epsilon \quad (6.10)$$

$$\epsilon \sim \mathcal{N}(0, \sigma_n^2) \quad (6.11)$$

or, in words, our observed values differ from the *truth* by identically, independently, distributed Gaussian noise with mean 0 and variance σ_n^2 . The assumption that the noise is i.i.d. is critical because it allows us to simplify the *likelihood* function by separating out each individual contribution by our datapoints:

$$p(\mathbf{y}|X, \mathbf{w}) := \prod_i^n p(\mathbf{y}_i|\mathbf{x}_i, \mathbf{w}) \quad (6.12)$$

$$= \prod_i^n \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{(\mathbf{y}_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma_n^2}\right) \quad (6.13)$$

$$= \frac{1}{(2\pi\sigma_n^2)^{n/2}} \exp\left(-\frac{1}{2\sigma_n^2}|\mathbf{y} - X^T \mathbf{w}|^2\right) \quad (6.14)$$

$$= \mathcal{N}(X^T \mathbf{w}, \sigma_n^2 I) \quad (6.15)$$

To perform inference with this updated model, we apply Baye's Rule, that is:

$$p(\mathbf{w}|\mathbf{y}, X) = \frac{p(\mathbf{y}|X, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|X)} \quad (6.16)$$

where

- $p(\mathbf{w}|\mathbf{y}, X)$ is the *posterior distribution*
- $p(\mathbf{y}|X, \mathbf{w})$ is the *likelihood*
- $p(\mathbf{w})$ is the *prior distribution*

- $p(\mathbf{y}|X)$ is the *marginal likelihood*, i.e. the normalization constant

It is now that the utility of choosing gaussian distributions for our likelihood and prior becomes clear. We have

$$p(\mathbf{w}|\mathbf{y}, X) \propto \exp\left(-\frac{1}{2\sigma_n^2}(\mathbf{y} - X^T \mathbf{w})^T (\mathbf{y} - X^T \mathbf{w})\right) \exp\left(-\frac{1}{2} \mathbf{w}^T \Sigma_p^{-1} \mathbf{w}\right) \quad (6.17)$$

Taking the log and expanding leads to

$$\log(p(\mathbf{w}|\mathbf{y}, X)) = \frac{1}{2} \left[\frac{1}{\sigma_n^2} \mathbf{y}^T \mathbf{y} - \frac{1}{\sigma_n^2} \mathbf{y}^T X^T \mathbf{w} - \frac{1}{\sigma_n^2} \mathbf{w}^T X \mathbf{y} + \frac{1}{\sigma_n^2} \mathbf{w}^T X X^T \mathbf{w} + \mathbf{w}^T \Sigma_p^{-1} \mathbf{w} \right] \quad (6.18)$$

$$= \frac{1}{2} \left[\mathbf{w}^T \left(\frac{1}{\sigma_n^2} X X^T + \Sigma_p^{-1} \right) \mathbf{w} - \left(\frac{1}{\sigma_n^2} \mathbf{y}^T X^T \right) \mathbf{w} - \mathbf{w}^T \left(\frac{1}{\sigma_n^2} X \mathbf{y} \right) + \mathbf{y}^T \frac{1}{\sigma_n^2} \mathbf{y} \right] \quad (6.19)$$

$$= \mathbf{w}^T A \mathbf{w} - B^T \mathbf{w} - \mathbf{w}^T B + C \quad (6.20)$$

where we have defined

$$A := \frac{1}{\sigma_n^2} X X^T + \Sigma_p^{-1} \quad (6.21)$$

$$B := \frac{1}{\sigma_n^2} X \mathbf{y} \quad (6.22)$$

$$C := \mathbf{y}^T \frac{1}{\sigma_n^2} \mathbf{y} \quad (6.23)$$

Now we can complete the square so that

$$\mathbf{w}^T A \mathbf{w} - B^T \mathbf{w} - \mathbf{w}^T B + C = (\mathbf{w} - \bar{\mathbf{w}})^T A (\mathbf{w} - \bar{\mathbf{w}}) + K \quad (6.24)$$

leading to

$$\bar{\mathbf{w}} = A^{-1} B = \frac{1}{\sigma_n^2} \left(\frac{1}{\sigma_n^2} X X^T + \Sigma_p^{-1} \right)^{-1} X \mathbf{y} \quad (6.25)$$

$$K = C - \bar{\mathbf{w}}^T A \bar{\mathbf{w}} \quad (6.26)$$

Since K does not depend on \mathbf{w} directly, it may be absorbed into the normalization of $p(\mathbf{w}|\mathbf{y}, X)$. Thus we are left with

$$p(\mathbf{w}|\mathbf{y}, X) = \mathcal{N}\left(\bar{\mathbf{w}} = \frac{1}{\sigma_n^2} A^{-1} X \mathbf{y}, \Sigma = A^{-1}\right) \quad (6.27)$$

$$A = \frac{1}{\sigma_n^2} X X^T + \Sigma_p^{-1} \quad (6.28)$$

This result gives us the gaussian distribution over the space of possible parameter vectors \mathbf{w} . To use this distribution to make predictions, consider a newly supplied testpoint \mathbf{x}_* . We want to find

$$p(y_*|\mathbf{x}_*, \mathbf{y}, X) \quad (6.29)$$

We do this by marginalizing over our weight distribution, i.e.

$$p(y_*|\mathbf{x}_*, \mathbf{y}, X) = \int_{\mathbf{w}} p(y_*|\mathbf{x}_*, \mathbf{w}) p(\mathbf{w}|\mathbf{y}, X) d\mathbf{w} \quad (6.30)$$

If we make the further assumption that testing points are i.i.d. gaussian distributed, we see that this integral is the product of two gaussians and therefore is also a gaussian. To find the mean and covariance of the predictive distribution, we check

$$\bar{y}_* = \mathbb{E}[y_*] = \mathbb{E}[\mathbf{x}_*^T \mathbf{w}] = \mathbf{x}_*^T \mathbb{E}[\mathbf{w}] = \mathbf{x}_*^T \bar{\mathbf{w}} \quad (6.31)$$

$$\text{Cov}(y_*) = \mathbb{E}[(y_* - \bar{y}_*)(y_* - \bar{y}_*)^T] \quad (6.32)$$

$$= \mathbb{E}[(\mathbf{x}_*^T \mathbf{w} - \mathbf{x}_*^T \bar{\mathbf{w}})(\mathbf{x}_*^T \mathbf{w} - \mathbf{x}_*^T \bar{\mathbf{w}})^T] \quad (6.33)$$

$$= \mathbb{E}[\mathbf{x}_*^T (\mathbf{w} - \bar{\mathbf{w}})(\mathbf{w} - \bar{\mathbf{w}})^T \mathbf{x}_*] \quad (6.34)$$

$$= \mathbf{x}_*^T \mathbb{E}[(\mathbf{w} - \bar{\mathbf{w}})(\mathbf{w} - \bar{\mathbf{w}})^T] \mathbf{x}_* \quad (6.35)$$

$$= \mathbf{x}_*^T \text{Cov}(\mathbf{w}) \mathbf{x}_* \quad (6.36)$$

$$= \mathbf{x}_*^T A^{-1} \mathbf{x}_* \quad (6.37)$$

so that

$$\boxed{p(y_*|\mathbf{x}_*, \mathbf{y}, X) = \mathcal{N}(\mathbf{x}_*^T \mathbf{w}, \mathbf{x}_*^T A^{-1} \mathbf{x}_*)} \quad (6.38)$$

Let's now take a break from our Bayesian regression discussion and return to the standard linear regression model for a moment. The key drawback of linear models like this is, of course, that they're *linear*!. Considering that many (most) *interesting* relationships are non-linear, how can we extend our simple linear model to enable us to perform complicated non-linear fits?

In the parlance of machine learning, the simple solution is to do feature engineering. If our initial feature vector is

$$\mathbf{x} = (x_1, \dots, x_n) \quad (6.39)$$

we can use our *expertise* to concoct new combinations of these features to produce the augmented vector

$$\tilde{\mathbf{x}} = (x_1, \dots, x_n, x_1^2, \sin(x_2), x_5 x_7 / x_4, \dots) \quad (6.40)$$

As an example, a linear classifier is unable to distinguish points inside a circle from those outside just from the (x, y) coordinates alone. Augmenting the feature vector to include the squared radius $x^2 + y^2$ as a new feature removes this obstacle. This works because the *linear* part of linear regression only refers to the fact that our model takes *linear combinations* of feature variables to produce its output. There is no restriction that the features themselves need to be independent variables! This same idea is what makes methods like SINDy work which use libraries of polynomial combinations of base features.

Constructing new features is often more art than science. To standardize the process, let's abstract the mapping from the original feature vector \mathbf{x} to the augmented vector $\tilde{\mathbf{x}}$. This is accomplished via the projection map $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^N$ where

$$\mathbf{x} \mapsto \tilde{\mathbf{x}} = \phi(\mathbf{x}) \quad (6.41)$$

The result is that our linear model updates to become

$$f(\mathbf{x}) := \phi(\mathbf{x})^T \mathbf{w} \quad (6.42)$$

where the weight vector has gone from D dimensional to N dimensional. Similarly, the normal equations for \mathbf{w} update to become

$$\mathbf{w} = (\Phi\Phi^T)^{-1}\Phi\mathbf{y} \quad (6.43)$$

where $\Phi = \phi(X)$ is the $N \times n$ matrix resulting from applying ϕ columnwise to X .

The following example shows how to use such a mapping to produce a quadratic polynomial fit.

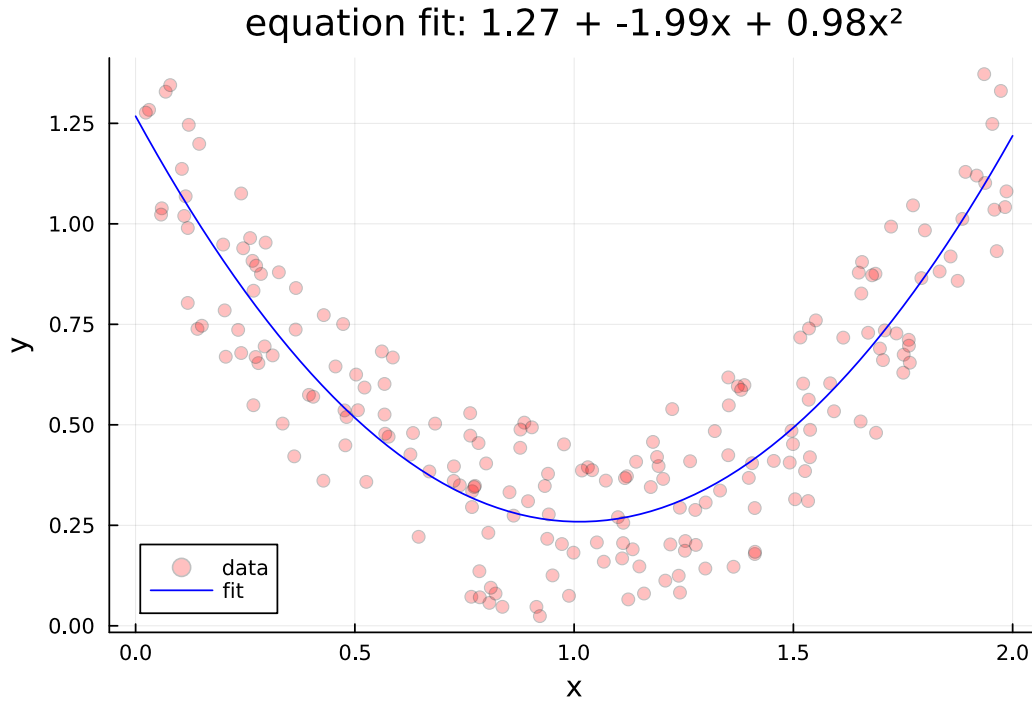


Figure 6.2. An example of polynomial regression for which we fit a quadratic polynomial to some noisy data.

We see from Figure 6.2 that our linear regression model found a great fit for a 2nd order polynomial when supplied with polynomial features. Let's update our Bayesian regression scheme to reflect the use of our feature projection map ϕ . First we define

$$\Phi := \phi(X) \quad (6.44)$$

$$\phi_* := \phi(\mathbf{x}_*) \quad (6.45)$$

Our predictive distribution therefore becomes

$$p(y_*|\mathbf{x}_*, X, \mathbf{y}) = \mathcal{N}\left(\frac{1}{\sigma_n^2}\phi_*^T A^{-1}\Phi\mathbf{y}, \phi_*^T A^{-1}\phi_*\right) \quad (6.46)$$

$$A = \frac{1}{\sigma_n^2}\Phi\Phi^T + \Sigma_p^{-1} \quad (6.47)$$

Great! Now we can do our Bayesian inference with non-linear features given by ϕ . There is a *massive* problem with this approach, however. Our order 2 polynomial map ϕ takes us from a D dimensional feature vector to $(D + 1)!$ many. This means that as we add more features to our feature map, the dimension of the resulting vector will quickly become prohibitively large. Looking at our current model equations, we see that the bottleneck is the matrix inversion of A which requires we invert an $N \times N$ matrix. Our prediction (i.e. the mean) involves multiplication on the right by the n dimensional vector \mathbf{y} . With that in mind, perhaps we can reformulate the above into an equivalent form using at most an $n \times n$ dimensional matrix.

Let $K := \Phi^T \Sigma_p \Phi$. Observe the following:

$$\frac{1}{\sigma_n^2}\Phi(K + \sigma_n^2 I) = \frac{1}{\sigma_n^2}\Phi(\Phi^T \Sigma_p \Phi + \sigma_n^2 I) \quad (6.48)$$

$$= \frac{1}{\sigma_n^2}\Phi\Phi^T \Sigma_p \Phi + \Phi I \quad (6.49)$$

$$= \left(\frac{1}{\sigma_n^2}\Phi\Phi^T\right) \Sigma_p \Phi + (\Phi I \Phi^{-1} \Sigma_p^{-1}) \Sigma_p \Phi \quad (6.50)$$

$$= \left(\frac{1}{\sigma_n^2}\Phi\Phi^T + \Sigma_p^{-1}\right) \Sigma_p \Phi \quad (6.51)$$

$$= A \Sigma_p \Phi \quad (6.52)$$

From there we see that

$$A^{-1} \frac{1}{\sigma_n^2} \Phi (K + \sigma_n^2 I) = \Sigma_p \Phi \quad (6.53)$$

$$\Rightarrow \frac{1}{\sigma_n^2} A^{-1} \Phi = \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \quad (6.54)$$

$$\Rightarrow \frac{1}{\sigma_n^2} \phi_*^T A^{-1} \Phi = \phi_*^T \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \quad (6.55)$$

For the covariance, we utilize the matrix inversion lemma [\(ADD REFERENCE HERE\)](#) which states

$$(Z + UWV^T)^{-1} = Z^{-1} - Z^{-1}U(W^{-1} + V^T Z^{-1}U)^{-1}V^T Z^{-1} \quad (6.56)$$

With the identification

$$Z^{-1} \rightarrow \Sigma_p \quad (6.57)$$

$$W^{-1} \rightarrow \sigma_n^2 I \quad (6.58)$$

$$V \rightarrow \Phi \quad (6.59)$$

$$U \rightarrow \Phi \quad (6.60)$$

we find

$$\Sigma_p - \Sigma_p \Phi (\Sigma_p + \Phi^T \Sigma_p \Phi)^{-1} \Phi^T \Sigma_p = \left(\Sigma_p^{-1} + \Phi \frac{1}{\sigma_n^2} I \Phi^T \right)^{-1} \quad (6.61)$$

$$= \left(\frac{1}{\sigma_n^2} \Phi \Phi^T + \Sigma_p^{-1} \right)^{-1} \quad (6.62)$$

$$= A^{-1} \quad (6.63)$$

Thus, we have the equivalent form for our predictive distribution:

$$\boxed{p(y_* | \mathbf{x}_*, X, \mathbf{y}) = \mathcal{N} \left(\phi_*^T \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \mathbf{y}, \phi_*^T \Sigma_p \phi_* - \phi_*^T \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \Phi^T \Sigma_p \phi_* \right)} \quad (6.64)$$

where the pesky $N \times N$ term has been replaced by the $n \times n$ matrix $\Phi^T \Sigma_p \Phi$.

We now make the *key* observation that the only matrices that appear in the above expression are

$$\Phi^T \Sigma_p \Phi, \quad \phi_*^T \Sigma_p \phi_* \quad (6.65)$$

$$\phi_*^T \Sigma_p \Phi, \quad \Phi^T \Sigma_p \phi_* \quad (6.66)$$

whose matrix elements we can write abstractly as

$$\phi(\mathbf{x})^T \Sigma_p \phi(\mathbf{x}') \quad (6.67)$$

To fit our model, we must determine appropriate values for the symmetric, positive semi-definite covariance matrix Σ_p (and σ_n too, technically). Instead, we observe that this matrix product is a quadratic form which we can think of as representing an inner product on our transformed vectors:

$$K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \quad (6.68)$$

We call the function $k(\mathbf{x}, \mathbf{x}')$ the *kernel function* or the *covariance function*.

All we need to perform the above calculations are the matrix elements of K applied to our data \mathcal{D} and any test points \mathbf{x}_* we wish to apply our model to. In effect, this means we are free to use feature vectors of any dimension, including ∞ . The idea here is that the kernel function represents an inner product over *some* vector space. As it turns out, the RBF kernel corresponds to a an infinite dimensional feature vector.

There are many choices for the kernel function. One of the most popular is the RBF (radial basis function) kernel, also commonly referred to as the *squared exponential kernel*:

$$k_{\text{rbf}}(\mathbf{x}, \mathbf{x}') := \sigma_f^2 \exp\left(-\frac{1}{2\ell^2}|\mathbf{x} - \mathbf{x}'|^2\right) \quad (6.69)$$

where σ_f^2 is the *signal variance* and ℓ denotes the similarity length scale.

For notational convenience, let's define

$$K := k(X, X) \quad (6.70)$$

$$K_{**} := k(X_*, X_*) \quad (6.71)$$

$$K_* := k(X, X_*) \quad (6.72)$$

then, our predictive distribution takes the final, *clean* form

$$\boxed{p(\mathbf{y}_* | X_*, X, \mathbf{y}) = \mathcal{N}\left(K_*^T (K + \sigma_n^2 I)^{-1} \mathbf{y}, K_{**} - K_*^T (K + \sigma_n^2 I)^{-1} K_*\right)} \quad (6.73)$$

This is the *end-result* of Gaussian Process Regression acheived via the *weight-space view*.

The Function-space View

So far our approach has been to generalize the standard linear regression model to allow for fitting over a (possibly infinite) basis of features with consideration for measurement and model uncertainty (our Bayesian priors). In essence, the idea was to fit the distribution of all possible weights conditioned on the available training data, $p(\mathbf{w}|X, \mathbf{y})$. A second *equivalent* approach is to instead consider the distribution of all possible model function $f(\mathbf{x})$. By constructing a Bayesian prior over this space, we constrain the space of possible model functions and learn a *distribution* over all allowed model functions, $p(f|X, \mathbf{y})$. To do so we will need to develop the abstract machinery of distributions over function spaces. When these distributions are Gaussian, the result is called a **Gaussian process**.

By this point, we are very familiar with the Gaussian distribution, a.k.a. the Normal distribution $\mathcal{N}(\mu, \sigma^2)$. This distribution is defined by a mean value μ and a variance σ^2 . It's *big brother* is the **Multivariate Normal Distribution**, $\mathcal{N}(\mu, \Sigma)$, described by a vector of means μ and a covariance matrix Σ . A natural question, then, is can we generalize the concept of the Gaussian distribution from N dimensions to being defined over a continuous field? This leads us naturally to the development of a so-called **Gaussian Process**.

Definition: A *Gaussian Process*, \mathcal{GP} , is a collection of random variables for which any finite subset are described by a joint Gaussian distribution.

To see where this comes from, recall that in our previous derivation, we already made the assumption that all our data points \mathcal{D} are i.i.d. Gaussian distributed. A gaussian process is the natural extension of this and makes the assumption that the continuous set from which are data are sampled are *so Gaussian* that any finite sample will be jointly Gaussian distributed. The term *process* is used to distinguish between finite collections of random variables (distributions) and their continuous counterparts described here.

Because each finite subset of this continuous collection is jointly gaussian, we can completely specify a Gaussian Process with two functions: the mean function $m(\mathbf{x})$ and the covariance function $k(\mathbf{x}, \mathbf{x}')$. To denote this, we typically write

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (6.74)$$

To see this in action, recall our Bayesian regression model

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w} \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p) \quad (6.75)$$

where we have set the prior on \mathbf{w} to have zero mean. The mean function is given by the expectation value of our model:

$$\mathbb{E}[f(\mathbf{x})] = \phi(\mathbf{x})^T \mathbb{E}[\mathbf{w}] = 0 \quad (6.76)$$

and the covariance function is given by

$$\mathbb{E}[f(\mathbf{x})f(\mathbf{x}')] = \phi(\mathbf{x})^T \mathbb{E}[\mathbf{w}\mathbf{w}^T] \phi(\mathbf{x}') = \phi(\mathbf{x})^T \Sigma_p \phi(\mathbf{x}') \quad (6.77)$$

To repeat the point, the key feature of Gaussian processes is that finite subsets are jointly Gaussian distributed. Thus we can split our data into the testpoints $\mathcal{D} = (X, \mathbf{y})$ and testpoints X_* and treat each collection as joint distributions with the following priors:

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix} \right) \quad (6.78)$$

where $\mathbf{f} := f(X)$ and $\mathbf{f}_* = f(X_*)$.

To obtain our predictive distribution, $p(\mathbf{f}_* | X_*, X, \mathbf{y})$, we *condition the joint prior distribution* on the observations. To see how this works, consider a general joint gaussian distribution given by

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right) \quad (6.79)$$

define the centered values $\tilde{x} := x - \mu_x$ and $\tilde{y} := y - \mu_y$. Define the intermediate variable

$$z := \tilde{x} - A\tilde{y} \quad (6.80)$$

Note that since we've subtracted out the mean we have $\mathbb{E}[\tilde{x}] = \mathbb{E}[\tilde{y}] = \mathbb{E}[z] = 0$. Let's now find A .

$$\mathbb{E}[z\tilde{y}^T] = \mathbb{E}[(\tilde{x} - A\tilde{y})\tilde{y}^T] \quad (6.81)$$

$$= \mathbb{E}[\tilde{x}\tilde{y}^T - A\tilde{y}\tilde{y}^T] \quad (6.82)$$

$$= \mathbb{E}[\tilde{x}\tilde{y}^T] - \mathbb{E}[A\tilde{y}\tilde{y}^T] \quad (6.83)$$

$$= \Sigma_{xy} - A\mathbb{E}[\tilde{y}\tilde{y}^T] \quad (6.84)$$

$$= \Sigma_{xy} - A\Sigma_{yy} \quad (6.85)$$

Therefore if we choose A so that z and \tilde{y} are independent and uncorrelated, then $\Sigma_{zy} = \mathbb{E}[z\tilde{y}^T] = 0$. Using this assumption, we find

$$0 = \mathbb{E}[z\tilde{y}^T] = \Sigma_{xy} - A\Sigma_{yy} \Rightarrow \boxed{A = \Sigma_{xy}\Sigma_{yy}^{-1}} \quad (6.86)$$

If we now condition \tilde{x} on \tilde{y} (i.e. look at \tilde{x} when \tilde{y} is constant), we find

$$\mathbb{E}[\tilde{x}|\tilde{y}] = A\tilde{y} + \mathbb{E}[z] \quad (6.87)$$

$$= A\tilde{y} + 0 \quad (6.88)$$

$$= \Sigma_{xy}\Sigma_{yy}^{-1}\tilde{y} \quad (6.89)$$

$$(6.90)$$

By manipulating this expression, we can now derive $\mathbb{E}[x|y]$ as follows:

$$\mathbb{E}[x|\tilde{y}] = \mathbb{E}[\tilde{x}|\tilde{y}] + \mu_x \quad (6.91)$$

$$= \mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}\tilde{y} \quad (6.92)$$

$$(6.93)$$

$$\boxed{\mathbb{E}[x|y] = \mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y)} \quad (6.94)$$

Similarly for the covariance, we have

$$\text{Cov}(x|y) = \text{Cov}(\tilde{x} + \mu_x|\tilde{y}) \quad (6.95)$$

$$= \text{Cov}(\tilde{x} + \mu_x|\tilde{y} + \mu_y) \quad (6.96)$$

$$= \text{Cov}(\tilde{x}|(\tilde{y} + \mu_y)) \quad (6.97)$$

$$= \text{Cov}(\tilde{x}|\tilde{y}) \quad (6.98)$$

$$= \text{Cov}((z + A\tilde{y})|\tilde{y}) \quad (6.99)$$

$$= \text{Cov}(z) + A\text{Cov}(\tilde{y}) \quad (6.100)$$

$$= \text{Cov}(z) + 0 \quad (6.101)$$

$$= \mathbb{E}[zz^T] \quad (6.102)$$

$$= \mathbb{E}[(\tilde{x} - A\tilde{y})(\tilde{x} - A\tilde{y})^T] \quad (6.103)$$

$$= \mathbb{E}[\tilde{x}\tilde{x}^T - A\tilde{y}\tilde{x}^T - x(A\tilde{y})^T + A\tilde{y}\tilde{y}^T A^T] \quad (6.104)$$

$$= \Sigma_{xx} - A\Sigma_{yx} - \Sigma_{xy}A^T + A\Sigma_{yy}A^T \quad (6.105)$$

$$= \Sigma_{xx} - (\Sigma_{xy}\Sigma_{yy}^{-1})\Sigma_{yx} - \Sigma_{xy}(\Sigma_{yy}^{-1})^T\Sigma_{xy}^T + \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_y(\Sigma_y^{-1})^T\Sigma_{xy}^T \quad (6.106)$$

$$= \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{xy}^T - \Sigma_{xy}(\Sigma_{yy}^{-1})^T\Sigma_{xy}^T + \Sigma_{xy}(\Sigma_{yy}^{-1})^T\Sigma_{xy}^T \quad (6.107)$$

$$= \Sigma_{xx} - \Sigma_{xy}[\Sigma_{yy}^{-1} - (\Sigma_{yy}^{-1})^T + (\Sigma_{yy}^{-1})^T]\Sigma_{xy}^T \quad (6.108)$$

$$= \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx} \quad (6.109)$$

$$\boxed{\text{Cov}(x|y) = \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx}} \quad (6.110)$$

Armed with this identity for joint Guassian distributions, we are ready to derive the predictive distribution for Gaussian Process Regression. We find:

$$p(\mathbf{f}_*|X_*, X, \mathbf{y}) = \mathcal{N}\left(K_*^T K^{-1}\mathbf{f}, K_{**} - K_*^T K^{-1}K_*\right). \quad (6.111)$$

To account for noisy observations, we can augment our correlation function to include a noise offset. The joint distribution then becomes:

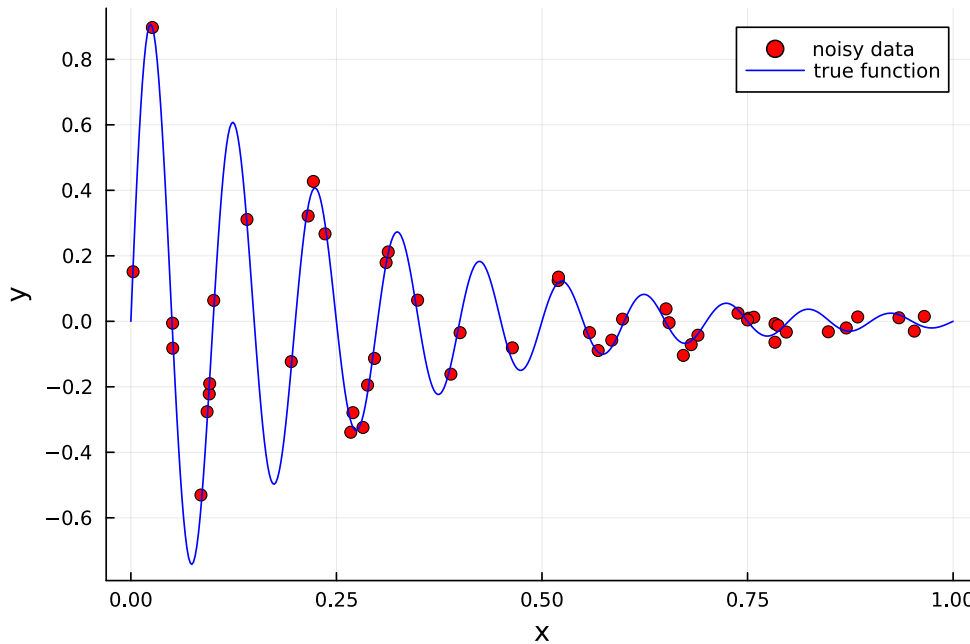
$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix} \right) \quad (6.112)$$

which leads to the predictive distribution

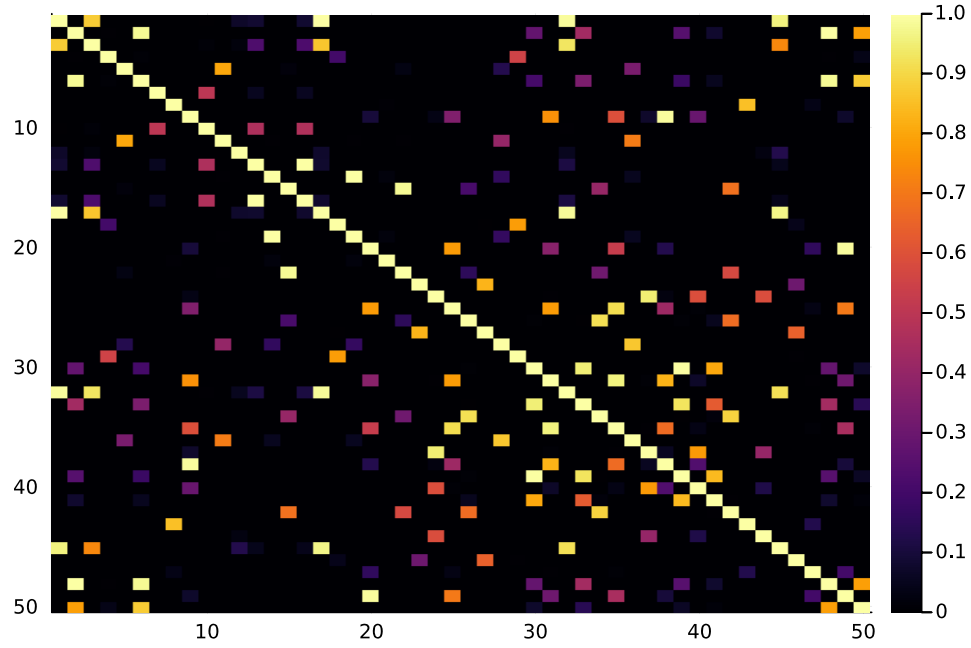
$$p(\mathbf{f}_* | X_*, X, \mathbf{y}) = \mathcal{N} \left(K_*^T [K + \sigma_n^2 I]^{-1} \mathbf{f}, K_{**} - K_*^T [K + \sigma_n^2 I]^{-1} K_* \right) \quad (6.113)$$

Doing it in Julia

To see how we can use this in practice, let's first construct some sample data which we seek to model as a Gaussian Process



The excellent package `KernelFunctions.jl` provides a clean interface to create various kernel functions and apply them to data to create our K -matrices. Due to the fact that kernel functions obey composition laws, we can easily build up complicated Kernels from basic pieces via function composition with `∘`



Unsurprisingly, there is a lot of activation on the diagonal as for a single datapoint \mathbf{x} , we have

$$k(\mathbf{x}, \mathbf{x}) = \exp\left(-\frac{0}{2\ell^2}\right) = 1.0 \quad (6.114)$$

The package `AbstractGPs.jl` provides an excellent way to define Gaussian Processes by supplying mean and kernel functions. We can then sample from our GPs with a simple interface designed to extend the basic functions from `Statistics.jl`. From an `AbstractGP` we can construct a `FiniteGP` by *indexing* into our datasets. The procedure can be summarized as follows

1. Build a kernel function $k(\cdot, \cdot)$ via composition using `KernelFunctions.jl`
2. Construct an a Gaussian Process $f \sim \mathcal{GP}$ abstractly using `AbstractGPs.jl`
3. Construct a finite representation of our GP, f_x , over our training data
4. Construct a posterior Gaussian Process from f_x and our training targets \mathbf{y} .
5. Construct a finite representation of the posterior GP applied to our prediction data

6. Sample this final distribution to obtain a prediction via `mean()` and variances via `var()`. Alternatively, we can obtain a multivariate normal distribution for each point by calling `marginals()`.

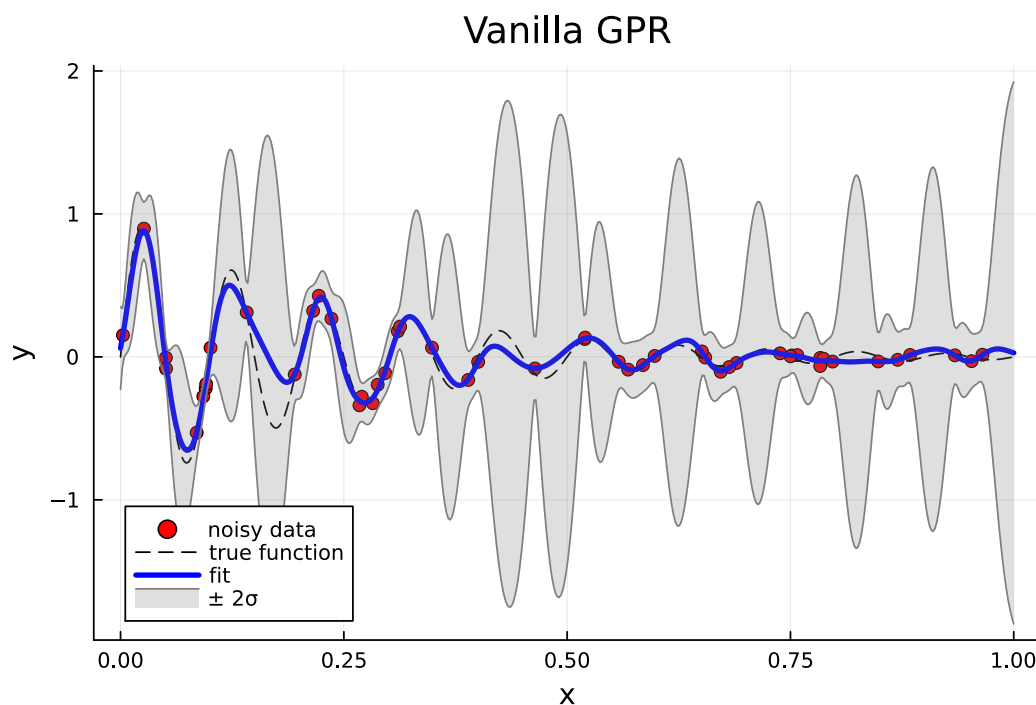


Figure 6.3. A Gaussian Process fit to the training data illustrating the prediction (means) and a $\pm 2\sigma$ uncertainty interval. We can clearly see how the uncertainty is larger for inputs far away from supplied training data.

Fitting the Kernel Hyperparameters

At this point, it is easy to think we are finished; we *already* fit the Gaussian process. However, we were forced to choose values for both ℓ and σ^2 kernel parameters. How can we optimally select the ideal hyperparameters for our Gaussian Process? This question leads us into the realm of Bayesian Model Selection. Rather than focusing specifically on our Gaussian Process model, let's take a step back and think about the process of model selection from a Bayesian perspective.

There are several levels of parameters in machine learning. At the lowest level, we have the model weights \mathbf{w} . Above that, we have model hyperparameters, θ . At the top we have model structure \mathcal{H} . In our Bayesian framework, we can consider prior distributions defined at each of these levels which codify credences for how much we trust a particular model, hyperparameter, etc. At the bottom, we have

$$p(\mathbf{w}|X, \mathbf{y}, \theta, \mathcal{H}_i) = \frac{p(\mathbf{y}|X, \mathbf{w}, \theta, \mathcal{H}_i)p(\mathbf{w}|\theta, \mathcal{H}_i)}{p(\mathbf{y}|X, \theta, \mathcal{H}_i)} \quad (6.115)$$

If this looks confusing, consider Bayes rule for 3 events R, H, S . We have:

$$P(R|H, S) = \frac{P(R, H, S)}{P(H, S)} \quad (6.116)$$

$$= \frac{P(H|R, S)P(R, S)}{P(H, S)} \quad (6.117)$$

$$= \frac{P(H|R, S)P(R|S)P(S)}{P(H|S)P(S)} \quad (6.118)$$

$$= \frac{P(H|R, S)P(R|S)}{P(H|S)} \quad (6.119)$$

To get the result, just think of θ and \mathcal{H}_i as a single *event* and translate the above to distribution functions. See [stack exchange link](#).

The prior $p(\mathbf{w}|\theta, \mathcal{H}_i)$ encodes any knowledge we have about the parameters prior to seeing the data. The denominator is the *marginal likelihood* and is given by

$$p(\mathbf{y}|X, \theta, \mathcal{H}_i) = \int d\mathbf{w} p(\mathbf{y}|X, \mathbf{w}, \theta, \mathcal{H}_i)p(\mathbf{w}|\theta, \mathcal{H}_i) \quad (6.120)$$

The next level up is to express the distribution of hyper-parameters θ :

$$p(\theta|X, \mathbf{y}, \mathcal{H}_i) = \frac{p(\mathbf{y}|X, \theta, \mathcal{H}_i)p(\theta|\mathcal{H}_i)}{p(\mathbf{y}|X, \mathcal{H}_i)}. \quad (6.121)$$

Here $p(\theta|\mathcal{H}_i)$ is called the *hyper-prior*. Similarly, the normalization constant is given by

$$p(\mathbf{y}|X, \mathcal{H}_i) = \int d\theta p(\mathbf{y}|X, \theta, \mathcal{H}_i)p(\theta|\mathcal{H}_i). \quad (6.122)$$

Finally, at the top level we have the set of possible model structures $\{\mathcal{H}_i\}$. This leads to

$$p(\mathcal{H}_i|X, \mathbf{y}) = \frac{p(\mathbf{y}|X, \mathcal{H}_i)p(\mathcal{H}_i)}{p(\mathbf{y}|X)} \quad (6.123)$$

with normlization constant

$$p(\mathbf{y}|X) = \sum_i p(\mathbf{y}|X, \mathcal{H}_i)p(\mathcal{H}_i). \quad (6.124)$$

Depending on the model details, these integrals may be intractable without approximations or Monte Carlo methods. Since we rarely have sufficient knowledge to form a hyperparameter prior, one often attempts to maximize the marginal likelihood $p(\mathbf{y}|X, \theta, \mathcal{H}_i)$ with respect to the hyperparameters θ instead. This is known as Type II Maximum Likelihood Estimation. (Refer to prior section on this in the Theoretical Techniques chapter). In the case of Gaussian Process Regression, we are once again saved by the fact that every piece has a convenient functional form resulting in analytically tractible integrals for the marginal likelihood function. We find

$$\ln p(\mathbf{y}|X, \theta) = -\frac{1}{2}\mathbf{y}^T(K_f + \sigma_n^2 I)^{-1}\mathbf{y} - \frac{1}{2}\ln|K_f + \sigma_n^2 I| - \frac{n}{2}\ln(2\pi) \quad (6.125)$$

where we have employed the natural logarithm to remove the pesky exponentials and arrive at a very convenient form. Note that because the logarithm is monotonically increasing, the maximum of the log-marginal-likelihood will be the same as the marginal-likelihood itself. Provided this functional form, we can use our favorite optimization routine to obtain the kernel hyperparameters which maximize the likelihood of obtaining our data as illustrated in the following comparison figure.

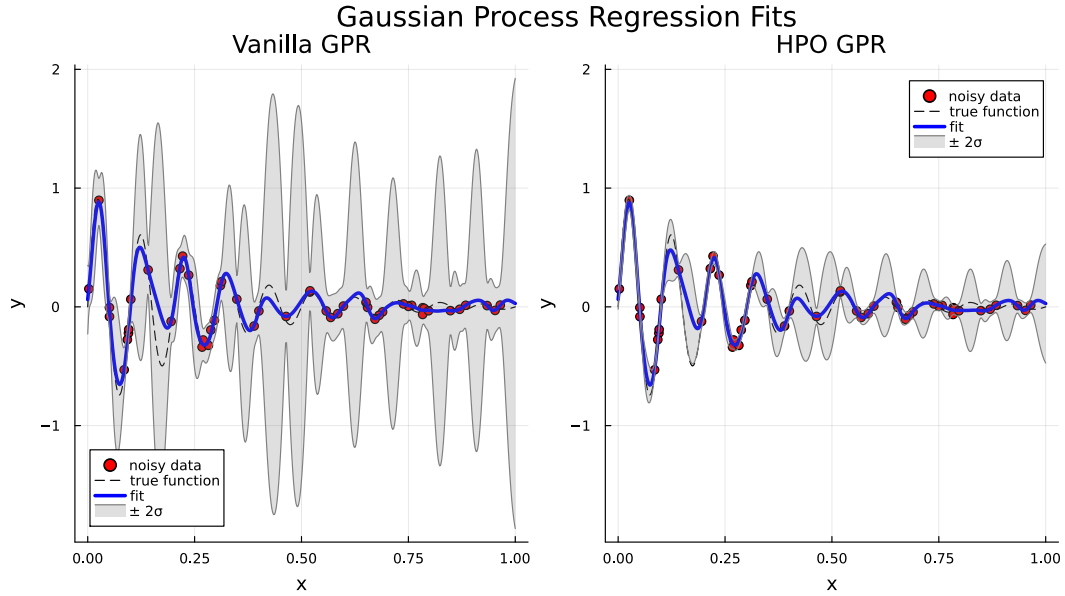


Figure 6.4. Left: the original Gaussian Process obtained with our default choice of kernel parameters. Right: A much better Gaussian Process obtained after performing hyperparameter optimization on the kernel parameters. Note how the mean function still does an excellent job fitting the data points while *also* minimizing the uncertainty of the fit.

6.3.3 Decision Trees

6.4 Unsupervised Classification

6.4.1 k-Means Clustering

6.4.2 Self Organizing Maps

Self organizing maps (SOMs) are an unsupervised machine learning technique developed by Kohonen (Kohonen, 1982) based on the simple biological principle that **neurons near each other fire together**. This observation that the topological **closeness** of similar computational units is a critical feature of intelligent systems leads to a natural reinterpretation of the familiar perceptron model into a new form amenable for a variety of clustering and dimensionality reduction tasks. In particular, the SOM enables a rapid unsupervised classification of multidimensional data into a (typically) one or two dimensional **simplicial complex**, the discrete realization of a topological manifold, whose vertices correspond to representative

points in the original data space \mathcal{D} . While a tad esoteric compared to other popular unsupervised methods like KMeans clustering or DBSCAN, the SOM distinguishes itself with the added benefit that it's training procedure guarantees nodes (i.e. classes) which are topologically close in the SOM simplex share similar weights. This additional structure makes the SOM particularly attractive when an interpretation of the discovered clusters *as well as* the relationships between them is desired.

We should add some more context (and references) for how the SOM has been used here. Perhaps we can ask David about his diatom identification from the Saharan dust sources.

The original treatment of the SOM by Kohonen was made in terms of processing units, sensory signals, and relaying networks [Kohonen-1982], however, in the modern era of deep learning, a more palatable derivation can be obtained by re-interpreting the weights of a simple perceptron model to provide the foundation for a clustering approach. As described in [Goodfellow-2016], a perceptron is a function of the form

$$\mathbf{y} = \sigma(W\mathbf{x}) \quad (6.126)$$

where $W \in \mathbb{R}^{n \times m}$ is a matrix of weights which transform the input $\mathbf{x} \in \mathbb{R}^m$ into \mathbb{R}^n and σ is a nonlinear *activation function* applied element-wise to the outputs of the matrix multiplication (indicated by the \cdot syntax). If we instead think of the weight matrix as an ordered collection of vectors $\{\mathbf{w}_i\}_{i=1}^n$, then this formula can be further decomposed into

$$\mathbf{y} = \sum_{i=1}^n \sigma(\mathbf{w}_i^T \mathbf{x}) = \sum_{i=1}^n \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) \quad (6.127)$$

The function of the perceptron is now clear: given an input vector \mathbf{x} and a collection of n -many weight vectors \mathbf{w}_i , compute the inner product of \mathbf{x} with each weight vector, apply the nonlinear activation function σ , and concatenate the results. If we allow ourselves to imagine the weight vectors \mathbf{w}_i as members of the same space as the inputs \mathbf{x} , a reasonable

question to ask is: *how similar is the input \mathbf{x} to each \mathbf{w}_i ?. Further, the application of the inner product $\langle \cdot, \cdot \rangle$ suggests we may answer this question in terms of the distance

$$\langle \mathbf{w}_i - \mathbf{x}, \mathbf{w}_i - \mathbf{x} \rangle = d(\mathbf{w}_i, \mathbf{x})^2. \quad (6.128)$$

In other words, given a set of weight vectors \mathbf{w}_i which we may now think of as the clusters for our unsupervised model, we can measure the similarity between a given datum \mathbf{x}_j and each cluster by computing the distance

$$d_{ij} = d(\mathbf{w}_i, \mathbf{x}_j) \quad (6.129)$$

The Training Process

Common SOM topologies

- Square
- Cylindrical
- Toroidal
- Spherical

A simple example: partitioning color spaces

Drawbacks of the SOM model

6.4.3 Generative Topographic Maps

6.5 Dimensionality Reduction

6.5.1 PCA

6.5.2 ICA

6.5.3 t-SNE

6.6 Model Ensembling

MLJ and SciKitLearn documentation sites will have good references for this, I think.

6.6.1 Bagging

e.g. Random Forests

6.6.2 Boosting

e.g. XGBoost

6.6.3 Stacking

e.g. Super Learners. Use the example of model stacking from MLJ documentation to describe our approach.

6.7 Uncertainty Quantification

The focus of this section can be on the concept of uncertainty quantification. For many physics theories that are nicely linearized, uncertainty analysis can be easily accomplished

at the level of first order sensitivities. That is, we can look at the Jacobian of our model to infer the behavior of small deviations about initial conditions. This approach does not easily extend to more complicated domains where the nonlinear effects dominate. Further, we also often want to establish ways to think about the fundamental instrument uncertainty for a measuring device. This can require meticulous calibrations which often assume a linear or polynomial fit... We can do better. Why not let the data tell us what the measurement uncertainty really is?

A good motivating example for the discussion of instrumental uncertainty is the use of a thermistor to measure temperature. One must assume a reasonable range of temperatures to establish the linear relationship between temperature and resistivity that is used to determine the temperature. However, the material characteristics of the thermistor that introduce nonlinearities at extreme temperatures don't necessarily mean we should have to throw out measurements that do not fall within this well-behaved range. Rather, we can perform a more sophisticated *calibration* to learn a model mapping resistivity to temperature that can account for these effects.

This is the bread-and-butter of the MINTS sensing efforts. Often low-cost sensing solutions provide decent measurements within a limited domain. With quality data from superior (but often prohibitively expensive) reference instruments, we can improve the default calibration to improve the reliability of data (by reducing uncertainty) and extend its domain of usefulness.

6.7.1 Quantile Regression

6.7.2 Probabilistic Models

e.g. discuss models like the GPR which fundamentally produce *distributions*. Similarly for models like Neural Networks that can produce multiple outputs, we can train a model to predict both μ and σ to make the model probabilistic.

6.7.3 Conformal Prediction

6.8 Scientific Machine Learning

6.8.1 Physics-Informed Neural Networks

6.8.2 Universal Differential Equations

6.8.3 Adjoint Methods for ODEs: Neural Ordinary Differential Equations

6.8.4 Hamiltonian Neural Networks

6.9 Generative Methods

6.9.1 Auto Encoders

6.9.2 Generative Adversarial Networks

6.10 Data Assimilation

NOTE: We should add further derivations from a Bayesian framework to fit in with what we've done for the Gaussian process regression and generative topographic maps. We should also add some more generic background to the topic of Data Assimilation, it's first uses in Meteorology,

The proper application of scientific models to make real-world predictions requires that we commit ourselves to a full accounting of all possible sources of uncertainty when reporting results. Further, the explosion of *big data* across scientific fields provides a plethora of observational data that our models are typically unequipped to incorporate when making predictions. The field of *Data Assimilation* addresses this problem by providing a family of techniques engineered to combine model output together with observational data whilst enabling a complete accounting of the sources of uncertainty. For chaotic systems in particular, data assimilation enables integration on long time scales that would be impossible via models alone. In this overview, we will follow the examples from (Ahmed et al., 2020).

Data assimilation can be cleanly developed in the framework of discrete dynamical systems. Since, at the end of the day, all of our scientific models must be discretized to be evaluated numerically, this is a reasonable course of action. Our goal is to find the best prediction for the system state vector u that combines our model predictions, also known as *forecasts*, with observational data. Model predictions are described via the discrete update equation:

$$u_{k+1} = \mathcal{M}(u_k; \theta) \quad (6.130)$$

For ODE systems, \mathcal{M} represents the time integration scheme for a models like

$$\frac{du}{dt} = f(u, t; \theta), \quad (6.131)$$

in other words, a particular choice of ODE integration scheme (Runge Kutta for example) used to evolve the state vector[✓] from time u_k to u_{k+1} .

To measure the performance of our assimilation scheme, we denote the *true* value of the state vector as $u^{(t)}$. The output of our model is denoted $u^{(b)}$ (b superscript for *background*). The discrepancy between the true value and our forecast is denoted $\xi^{(b)} = u^{(t)} - u^{(b)}$ characterizing the extent to which our model prediction is imperfect. The observations of our system are denoted by $w_k = w(t_k)$. These observations do not necessarily need to be components of the state vector u , but rather, are related to it via the *observation function*, $h : u_k \mapsto w_k$. For example, one may attempt to predict surface sea temperatures using data assimilation with data from satellite observations. The function h would then be the Stefan-Boltzmann relating the measured spectra to temperature. However, real world data is *also* imperfect, which we must take into account let we over constrain our model with poor quality data. We write

$$w_k = h(u_k) + \xi_k^{(m)} \quad (6.132)$$

where $\xi_k^{(m)}$ denotes this measurement error.

Given our model predictions $u_k^{(b)}$ and observations w_k , we seek to obtain the *optimal* or best-possible prediction called the **analysis**, $u^{(a)}$. Despite our care, this analysis will still be imperfect, so we further define the analysis error as

$$\xi^{(a)} = u^{(t)} - u^{(a)} \quad (6.133)$$

In summary, we have defined the following relevant quantities (at time t_k):

$$u_k^{(t)} \in \mathbb{R}^n \quad \text{the true state vector} \quad (6.134)$$

$$u_k^{(b)} \in \mathbb{R}^n \quad \text{the model forecast} \quad (6.135)$$

$$u_k^{(a)} \in \mathbb{R}^n \quad \text{the analysis} \quad (6.136)$$

$$w_k \in \mathbb{R}^m \quad \text{the observation vector} \quad (6.137)$$

$$\xi^{(b)} \in \mathbb{R}^n \quad \text{the model forecast error} \quad (6.138)$$

$$\xi^{(m)} \in \mathbb{R}^m \quad \text{the observation noise vector} \quad (6.139)$$

$$\xi^{(a)} \in \mathbb{R}^n \quad \text{the analysis error} \quad (6.140)$$

$$\xi^{(p)} \in \mathbb{R}^n \quad \text{the process noise if we used our model on the true state} \quad (6.141)$$

$$\mathcal{M} : \mathbb{R}^n \rightarrow \mathbb{R}^n \quad \text{the discrete model evolution function} \quad (6.142)$$

$$f : \mathbb{R}^n \rightarrow \mathbb{R}^n \quad \text{differential equation model (the right hand side)} \quad (6.143)$$

$$h : \mathbb{R}^n \rightarrow \mathbb{R}^m \quad \text{observation function} \quad (6.144)$$

To move forward, we now establish the following (prior) assumptions about the distribution of errors in each component in order to make it possible to derive a closed form for the final analysis. We require

$$\mathbb{E}[\xi_k^{(b)}] = 0 \quad \mathbb{E}[\xi_k^{(b)}(\xi_j^{(b)})^T] = 0 \text{ for } k \neq j \quad (6.145)$$

$$\mathbb{E}[\xi_k^{(m)}] = 0 \quad \mathbb{E}[\xi_k^{(m)}(\xi_j^{(m)})^T] = 0 \text{ for } k \neq j \quad (6.146)$$

$$\mathbb{E}[\xi_k^{(b)}(u_0)^T] = 0 \quad \mathbb{E}[\xi_k^{(m)}(u_0)^T] = 0 \quad (6.147)$$

$$\mathbb{E}[\xi_k^{(b)}\xi_j^{(m)}] = 0 \quad (6.148)$$

or in words, the errors in our model and observation are unbiased (e.g. mean zero) and uncorrelated.

We also define the error covariance matrices

$$Q_k := \mathbb{E}[\xi_k^{(p)}(\xi_k^{(p)})^T] \quad (6.149)$$

$$R_k := \mathbb{E}[\xi_k^{(m)}(\xi_k^{(m)})^T] \quad (6.150)$$

$$B_k := \mathbb{E}[\xi_k^{(b)}(\xi_k^{(b)})^T] \quad (6.151)$$

which we will use in our consideration of the final error of our analysis.

6.10.1 Kalman Filter

Given some model for the error covariance matrices Q_k and R_k , we would like a method that propagates *both* our model *and* the errors forward. This way we may guarantee that the accuracy of our analysis doesn't come at the cost of higher uncertainty.

The original implementation of the Kalman filter was for strictly linear systems, in particular as a filter for signal processing. We will first develop the analysis for this simplified case and then will generalize to the *Extended Kalman Filter* (EKF) that can handle fully nonlinear situations.

In the linear case, our system may be written as

$$u_{k+1}^{(t)} = M_k u_k^{(t)} + \xi_{k+1}^{(p)} \quad (6.152)$$

$$w_k = H_k u_k^{(t)} + \xi_k^{(m)} \quad (6.153)$$

where M_k and H_k are now matrices defining the linear problem.

The goal of the Kalman filter is to derive the analysis $u^{(a)}$ which minimizes the trace of the analysis error covariance matrix (i.e. sum of squared errors):

$$\text{Tr}(P_k) := \mathbb{E}[(u_k^{(t)} - u_k^{(a)})^T (u_k^{(t)} - u_k^{(a)})] \quad (6.154)$$

Finding the analysis consists of two steps: the forecast step and the assimilation step.

Beginning with the forecasting step, assume we have the analysis at time t_k denoted $u_k^{(a)}$.

Then the forecast for time t_{k+1} is

$$u_{k+1}^{(b)} = M_k u_k^{(a)}, \quad (6.155)$$

or in words, we obtain the next prediction by evolving the previous analysis forward. The background error is therefore

$$\xi_{k+1}^{(b)} = u_{k+1}^{(t)} - u_{k+1}^{(b)} \quad (6.156)$$

$$= M_k u_k^{(t)} + \xi_{k+1}^{(p)} - M_k u_k^{(a)} \quad (6.157)$$

$$= M_k \left(u_k^{(t)} - u_k^{(a)} \right) + \xi_{k+1}^{(p)} \quad (6.158)$$

$$= M_k \xi_k^{(a)} + \xi_{k+1}^{(p)} \quad (6.159)$$

We may now evaluate the covariance matrix of our background estimate as:

$$B_{k+1} = \mathbb{E}[\xi_{k+1}^{(b)} (\xi_{k+1}^{(b)})^T] \quad (6.160)$$

$$= \mathbb{E} \left[\left(M_k \xi_k^{(a)} + \xi_{k+1}^{(p)} \right) \left(M_k \xi_k^{(a)} + \xi_{k+1}^{(p)} \right)^T \right] \quad (6.161)$$

$$(6.162)$$

If we presume that $\mathbb{E}[\xi_k^{(b)} (\xi_{k+1}^{(p)})^T] = 0$, then the cross terms vanish and we are left with

$$\boxed{B_{k+1} = M_k P_k M_k^T + Q_{k+1}} \quad (6.163)$$

Thus we now have the background (i.e forecast) estimate of the state at t_{k+1} and its covariance matrix. Given a measurement w_{k+1} at the same time with covariance matrix R_{k+1} , we may now perform the assimilation step where we fuse the two sources of information to obtain $u_{k+1}^{(a)}$ and P_{k+1} .

Let's suppose that the analysis has the form

$$u_{k+1}^{(a)} = \nu + K_{k+1} w_{k+1} \quad (6.164)$$

for some vector $\nu \in \mathbb{R}^n$ and matrix $K_{k+1} \in \mathbb{R}^{m \times n}$. In a perfect world, we would have

$$\mathbb{E}[u_k^{(t)} - u_k^{(a)}] = 0. \text{ Therefore,}$$

$$0 = \mathbb{E}[u_k^{(t)} - u_k^{(a)}] \quad (6.165)$$

$$= \mathbb{E}[(u_k^{(b)} + \xi_k^{(b)}) - (\nu + K_k w_k)] \quad (6.166)$$

$$= \mathbb{E}[(u_k^{(b)} + \xi_k^{(b)}) - (\nu + K_k H_k u_k^{(t)} + K_k \xi_k^{(m)})] \quad (6.167)$$

$$= \mathbb{E}[u_k^{(b)}] + \mathbb{E}[\xi_k^{(b)}] - \mathbb{E}[\nu] - K_k H_k \mathbb{E}[u_k^{(t)}] - K_k \mathbb{E}[\xi_k^{(m)}] \quad (6.168)$$

$$= u_k^{(b)} + 0 - \nu - K_k H_k u_k^{(b)} - 0 \quad (6.169)$$

$$= u_k^{(b)} - \nu - K_k H_k u_k^{(b)} \quad (6.170)$$

$$\Rightarrow \nu = u_k^{(b)} - K_k H_k u_k^{(b)} \quad (6.171)$$

which we now substitute to obtain

$$\boxed{u_k^{(a)} = u_k^{(b)} + K_k(w_k - H_k u_k^{(b)})} \quad (6.172)$$

Now that we know the form for the analysis we may derive the optimal matrix K_k by optimization of P_k . We have

$$\xi_k^{(a)} = u_k^{(t)} - u_k^{(a)} \quad (6.173)$$

$$= M_{k-1}u_{k-1}^{(t)} + \xi_k^{(p)} - u_k^{(b)} - K_k \left(w_k - H_k u_k^{(b)} \right) \quad (6.174)$$

$$= M_{k-1}u_{k-1}^{(t)} + \xi_k^{(p)} - M_{k-1}u_{k-1}^{(a)} - K_k \left(H_k u_k^{(t)} + \xi_k^{(m)} - H_k u_k^{(b)} \right) \quad (6.175)$$

$$= M_{k-1}u_{k-1}^{(t)} + \xi_k^{(p)} - M_{k-1}u_{k-1}^{(a)} - K_k H_k u_k^{(t)} - K_k \xi_k^{(m)} + K_k H_k u_k^{(b)} \quad (6.176)$$

$$= M_{k-1}u_{k-1}^{(t)} + \xi_k^{(p)} - M_{k-1}u_{k-1}^{(a)} - K_k H_k u_k^{(t)} - K_k \xi_k^{(m)} + K_k H_k u_k^{(b)} \quad (6.177)$$

$$= \left\{ M_{k-1}(\xi_{k-1}^{(a)} + u_{k-1}^{(a)}) + \xi_k^{(p)} - M_{k-1}u_{k-1}^{(a)} \right\} - K_k H_k u_k^{(t)} - K_k \xi_k^{(m)} + K_k H_k u_k^{(b)} \quad (6.178)$$

$$= \left\{ M_{k-1}\xi_{k-1}^{(a)} + \xi_k^{(p)} \right\} - K_k H_k u_k^{(t)} + K_k H_k u_k^{(b)} - K_k \xi_k^{(m)} \quad (6.179)$$

$$= M_{k-1}\xi_{k-1}^{(a)} + \xi_k^{(p)} - K_k H_k (M_{k-1}u_{k-1}^{(t)} + \xi_k^{(b)}) + K_k H_k u_k^{(b)} - K_k \xi_k^{(m)} \quad (6.180)$$

$$= M_{k-1}\xi_{k-1}^{(a)} + \xi_k^{(p)} - K_k H_k M_{k-1}(\xi_{k-1}^{(a)} + u_{k-1}^{(a)}) - K_k H_k \xi_k^{(b)} + K_k H_k u_k^{(b)} - K_k \xi_k^{(m)} \quad (6.181)$$

$$= M_{k-1}\xi_{k-1}^{(a)} + \xi_k^{(p)} - K_k H_k M_{k-1}(\xi_{k-1}^{(a)} + u_{k-1}^{(a)}) - K_k H_k \xi_k^{(b)} + K_k H_k M_{k-1}u_{k-1}^{(a)} - K_k \xi_k^{(m)} \quad (6.182)$$

$$= M_{k-1}\xi_{k-1}^{(a)} + \xi_k^{(p)} - K_k H_k M_{k-1}\xi_{k-1}^{(a)} - K_k H_k \xi_k^{(b)} - K_k \xi_k^{(m)} \quad (6.183)$$

$$= (I - K_k H_k)(M_{k-1}\xi_{k-1}^{(a)} - \xi_k^{(p)}) - K_k \xi_k^{(m)} \quad (6.184)$$

$$(6.185)$$

and therefore the covariance matrix is

$$P_k = \mathbb{E}[\xi_k^{(a)}(\xi_k^{(a)})^T] \quad (6.186)$$

$$= \mathbb{E} \left[\left((I - K_k H_k)(M_{k-1} \xi_{k-1}^{(a)} - \xi_k^p) - K_k \xi_k^{(m)} \right) \left((I - K_k H_k)(M_{k-1} \xi_{k-1}^{(a)} - \xi_k^p) - K_k \xi_k^{(m)} \right)^T \right] \quad (6.187)$$

$$\begin{aligned} &= (I - K_k H_k) M_{k-1} \mathbb{E}[\xi_{k-1}^{(a)}(\xi_{k-1}^{(a)})^T] M_{k-1}^T (I - K_k H_k)^T \\ &\quad + (I - K_k H_k) \mathbb{E}[\xi_k^{(p)}(\xi_k^{(p)})^T] (I - K_k H_k)^T \\ &\quad - K_k \mathbb{E}[\xi_k^{(m)}(\xi_k^{(m)})^T] K_k^T \end{aligned} \quad (6.188)$$

$$= (I - K_k H_k) B_k (I - K_k H_k)^T - K_k R_k K_k^T \quad (6.189)$$

Now all that remains is to derive the specific form for K_k . The Kalman filter is defined as that K_k which minimizes the sum of squared analysis errors, i.e. the trace of the analysis error covariance matrix. Therefore, the following identities from matrix calculus will be useful:

$$\nabla_A \text{tr}(AB) = B^T \quad (6.190)$$

$$\nabla_A \text{tr}(BA^T) = B \quad (6.191)$$

$$\nabla_A \text{tr}(ABA^T) = AB^T + AB \quad (6.192)$$

$$(6.193)$$

from which we obtain

$$0 = \nabla_{K_k} \text{tr}(P_k) \quad (6.194)$$

$$= \nabla_{K_k} \left\{ B_k - B_k H_k^T K_k^T - K_k H_k B_k + K_k H_k B_k H_k^T B_k^T - K_k R_k K_k^T \right\} \quad (6.195)$$

$$= -B_k H_k^T - (H_k B_k)^T + K_k [H_k B_k H_k^T + (H_k B_k H_k^T)^T - R_k + R_k^T] \quad (6.196)$$

$$= -2B_k H_k^T + 2K_k (H_k B_k H_k^T - R_k) \quad (6.197)$$

$$\Rightarrow K_k = B_k H_k^T [H_k B_k H_k^T - R_k]^{-1} \quad (6.198)$$

we now substitute this result to obtain a simplified form for P_k :

$$P_k = (I - K_k H_k) B_k (I - K_k H_k)^T + K_k R_k K_k^T \quad (6.199)$$

$$= (I - K_k H_k) B_k - (I - K_k H_k) B_k (K_k H_k)^T + K_k R_k K_k^T \quad (6.200)$$

$$= (I - K_k H_k) B_k - \left\{ (I - K_k H_k) B_k (K_k H_k)^T + K_k R_k K_k^T \right\} \quad (6.201)$$

$$= (I - K_k H_k) B_k - \left\{ (I - K_k H_k) B_k H_k^T K_k^T + K_k R_k K_k^T \right\} \quad (6.202)$$

$$= (I - K_k H_k) B_k - \left\{ (I - K_k H_k) B_k H_k^T + K_k R_k \right\} K_k^T \quad (6.203)$$

$$= (I - K_k H_k) B_k - \left\{ B_k H_k^T - K_k (H_k B_k H_k^T + R_k) \right\} K_k^T \quad (6.204)$$

$$= (I - K_k H_k) B_k - \left\{ B_k H_k^T - B_k H_k^T \right\} K_k^T \quad (6.205)$$

$$= (I - K_k H_k) B_k \quad (6.206)$$

where we have used the fact that covariance matrices are symmetric.

Now that we have all of the pieces, let's summarize the procedure:

1. **Initialization:** We must set the system to some initial condition. This means we must define u_0^a and P_0 . We must also come up with a model for the process noise covariance Q_k and measurement error covariance R_k .
2. **Forecast Step:** Starting with the analysis state vector $u_{k-1}^{(a)}$ and covariance P_{k-1} , we apply the model time evolution operator M_{k-1} to obtain the predicted state $u_k^{(b)}$ and covariance B_k at the next time step.

$$u_k^{(b)} = M_{k-1} u_{k-1}^{(a)} \quad (6.207)$$

$$B_k = M_{k-1} P_{k-1} M_{k-1}^T + Q_k \quad (6.208)$$

3. **Assimilation Step:** Combine our observations w_k and their associated uncertainties contained in R_k together with the background prediction $u_k^{(b)}$ and its error covariance

matrix B_k to obtain the analysis state $u_k^{(a)}$ and error covariance P_k .

$$K_k = B_k H_k^T \left[H_k B_k H_k^T - R_k \right]^{-1} \quad (6.209)$$

$$u_k^{(a)} = u_k^{(b)} + K_k (w_k - H_k u_k^{(b)}) \quad (6.210)$$

$$P_k = (I - K_k H_k) B_k \quad (6.211)$$

6.10.2 Extended Kalman Filter

Given the nonlinear nature of many scientific models, it is desirable to extend out notion of the *Kalman Filter* to be able to handle nonlinear models $f(\cdot)$ (and by extension, their update function $\mathcal{M}(\cdot)$), and nonlinear observation functions $h(\cdot)$. This can be accomplished so long as these functions are sufficiently smooth (C^1 to be precise) so as to admit valid Taylor approximations to first order. That is,

$$\mathcal{M}(u_k) \approx \mathcal{M}(u_k^{(a)}) + D_M(u_k^{(a)}) \xi_k^{(a)} \quad h(u_k) \approx h(u_k^{(b)}) + D_h(u_k^{(b)}) \xi_k^{(b)} \quad (6.212)$$

$$D_M := \left[\frac{\partial \mathcal{M}_i}{\partial u_j} \right] \quad D_h := \left[\frac{\partial h_i}{\partial u_j} \right] \quad (6.213)$$

where \mathcal{M}_i and h_i denote the i th component functions of \mathcal{M} and h respectively. For the so-called *Extended Kalman Filter* (EKF), approach. We assume that the analysis errors can be reasonably approximated to evolve linearly such that we can use the above *linearized* approximations to our model update and observation functions to reasonably model our uncertainties.

We therefore replace the (previously) linear functions M_k and H_k with their nonlinear Jacobian derived counterparts to obtain the following procedure

1. **Initialization:** To begin we must choose values for $u_0^{(a)}$ and P_0 . We must also provide models for Q_k and R_k .
2. **Forecast Step:**

$$u_k^{(b)} = \mathcal{M}(u_{k-1}^{(a)}) \quad (6.214)$$

$$B_k = D_M(u_{k-1}^{(a)}) P_{k-1} D_M^T(u_{k-1}^{(a)}) + Q_k \quad (6.215)$$

3. Assimilation Step:

$$K_k = B_k D_M^T(u_k^{(b)}) \left[D_h(u_k^{(b)}) B_k D_h^T(u_k^{(b)}) + R_k \right]^{-1} \quad (6.216)$$

$$u_k^{(a)} = u_k^{(b)} + K_k (w_k - h(u_k^{(b)})) \quad (6.217)$$

$$P_k = \left(I - K_k D_h(u_k^{(b)}) \right) B_k \quad (6.218)$$

Note that the most challenging piece in this implementation is to stably compute the model update Jacobian D_M , as this is not just the Jacobian of our continuous model $f(\cdot)$. For example, if we are modeling the state of the system via a set of ODEs, then D_M is the Jacobian of the output of a single integrator step with respect to the initial starting state. This can pose significant computational challenges where adaptive solvers are desired as we must be able to back-propagate out derivative information for all evaluations of the function $f(\cdot)$ taken during the step.

6.10.3 continuous-discrete Extended Kalman Filter

UPDATE REQUIRED! SEE NOTEBOOK

6.10.4 3d-Var

For the *Kalman Filter* and the *EKF*, we derived an optimal way to combine observation with simulation so as to minimize the trace of the analysis error covariance matrix, P_k . An alternative approach is to recast the problem as a pure optimization problem where rather than finding a filter K_k that will add an innovation to $u_k^{(b)}$ to obtain the analysis $u_k^{(a)}$, we obtain the analysis by directly optimizing a cost function

$$J(u) = \frac{1}{2} (w - h(u))^T R^{-1} (w - h(u)) + \frac{1}{2} (u - u^{(b)})^T B^{-1} \frac{1}{2} (u - u^{(b)}) \quad (6.219)$$

which we can justify as coming from the joint probability distribution (assuming Gaussian errors)

$$\mathcal{P}(u|w) = C \exp \left(-\frac{1}{2} (u - u^{(b)})^T B^{-1} \frac{1}{2} (u - u^{(b)}) \right) \cdot \exp \left(-\frac{1}{2} (w - h(u))^T R^{-1} (w - h(u)) \right) \quad (6.220)$$

with model error covariance B and measurement error covariance R as before. This is clearly a *very strong assumption*.

To optimize $J(u)$, we begin by taking it's gradient.

$$\nabla_u J(u) = -D_h^T R^{-1} (w - h(u)) + B^{-1} (u - u^{(b)}) \quad (6.221)$$

Thus, finding the analysis $u^{(a)}$ amounts to solving the system

$$a - D_h^T R^{-1} (w - h(u^{(a)})) + B^{-1} (u^{(a)} - u^{(b)}) = 0 \quad (6.222)$$

As for Kalman filtering, let's begin with the assumption that our model and observation function are linear. Suppose that we have $h(u) = Hu$ so that $D_h(u) = H$. It follows that

$$D_h^T R^{-1} (w - Hu^{(a)}) = B^{-1} (u^{(a)} - u^{(b)}) \quad (6.223)$$

$$D_h^T R^{-1} w - D_h^T R^{-1} H u^{(a)} = B^{-1} u^{(a)} - B^{-1} u^{(b)} \quad (6.224)$$

$$(D_h^T R^{-1} H + B^{-1}) u^{(a)} = D_h^T R^{-1} + B^{-1} u^{(b)} \quad (6.225)$$

$$(H^T R^{-1} H + B^{-1}) u^{(a)} = H^T R^{-1} + B^{-1} u^{(b)} \quad (6.226)$$

Thus we see that the analysis is given by

$$u^{(a)} = u^{(b)} + B H^T (R + H B^T H)^{-1} (w - H u^{(b)}) \quad (6.227)$$

which agrees with what we found for the Linear Kalman Filter.

Moving on to the non-linear case, we can expand h about an initial guess $u^{(c)}$ which we will later choose to be $u^{(b)}$ for convenience.

$$h(u^{(a)}) \approx h(u^{(c)}) + D_h(u^{(c)}) \Delta u \quad (6.228)$$

Using this, we have

$$D_h^T(u^{(a)})R^{-1}(w - h(u^{(a)})) = B^{-1}(u^{(a)} - u^{(b)}) \quad (6.229)$$

$$D_h^T(u^{(a)})R^{-1}(w - h(u^{(c)}) - D_h(u^{(c)})\Delta u) \approx B^{-1}(u^{(c)} + \Delta u - u^{(b)}) \quad (6.230)$$

$$D_h^T(u^{(c)})R^{-1}(w - h(u^{(c)}) - D_h(u^{(c)})\Delta u) \approx B^{-1}(u^{(c)} + \Delta u - u^{(b)}) \quad (6.231)$$

which we now solve for the update Δu to obtain the linear system

$$(B^{-1} + D_h^T(u^{(c)})R^{-1}D_h(u^{(c)})) \Delta u = B^{-1}(u^{(b)} - u^{(c)}) + D_h^T(u^{(c)})R^{-1}(w - h(u^{(c)})) \quad (6.232)$$

Thus we have the following prescription

1. To begin, take $u^{(c)} == u^{(b)}$.
2. Solve the system

$$(B^{-1} + D_h^T(u^{(c)})R^{-1}D_h(u^{(c)})) \Delta u = B^{-1}(u^{(b)} - u^{(c)}) + D_h^T(u^{(c)})R^{-1}(w - h(u^{(c)})) \quad (6.233)$$

to obtain Δu

3. Update your guess using your favorite optimization algorithm. For example, in steepest descent, choose a learning rate η and set

$$u_{\text{new}}^{(c)} = u_{\text{prev}}^{(c)} - \eta \Delta u \quad (6.234)$$

4. Repeat the procedure until $|u_{\text{new}}^{(c)} - u_{\text{prev}}^{(c)}|$ converges to a desired tolerance.

In both the linear and nonlinear case, it should be noted that we have not added time indices to our state vectors. This is an indication that the 3d-var procedure is performed *at every time where you have observation data*.

6.10.5 4d-Var

The *3d-Var* algorithm attempts to optimize a cost function point-by-point to obtain the ideal analysis at each time where we have observation data. This can become computationally expensive as we require model evaluations *and* an optimization routine for every observation point. An alternative approach is to simultaneously optimize across all observations in order to obtain the ideal *initial condition* which achieves the best model fit. In other words, this amounts to parameter fitting for our model where we think of the initial conditions as the parameters. The difference is two-fold: We extend the usual ℓ_2 function into a quadratic form utilizing the observation error covariance matrix R_k . We also typically only have a small number of observables related to the state vector (potentially non-linearly) and therefore must transform the state vector u_k into the observation space via h before computing the loss against our observations w_k . This results in a loss function of the form

$$\begin{aligned} J(u_0) &= \frac{1}{2} \left(u_0 - u_0^{(b)} \right)^T B^{-1} \left(u_0 - u_0^{(b)} \right) + \frac{1}{2} \sum_k (w_k - h(u_k))^T R_k^{-1} (w_k - h(u_k)) \quad (6.235) \\ &= J_b(u_0) + J_m(u_0). \end{aligned}$$

Note that the first term is usefull if we already have an initial guess $u_0^{(b)}$ for the inital condition in mind. If we do not have one, we may ommit this term.

As before, we now want to optimize this cost function. To do so, we first observe that

$$u_k = \mathcal{M}^{(k)}(u_0; \theta) \quad (6.236)$$

It is easy to obtain the gradient of J_b so we shall focus on the second term. We find that

$$\nabla_{u_0} J_m = \nabla_{u_0} \left\{ \sum_k \frac{1}{2} (w_k - h(u_k))^T R_k^{-1} (w_k - h(u_k)) \right\} \quad (6.237)$$

$$= - \sum_k \left[\frac{\partial}{\partial u_0} h(\mathcal{M}^{(k-1)}(u_0)) \right]^T R_k^{-1} (w_k - h(u_k)) \quad (6.238)$$

$$= - \sum_k [D_h(u_k) D_M(u_{k-1}) D_M(u_{k-2}) \cdots D_M(u_0)]^T R_k^{-1} (w_k - h(u_k)) \quad (6.239)$$

$$= - \sum_k [D_M^T(u_0) D_M^T(u_1) \cdots D_M^T(u_{k-1}) D_h^T(u_k)] R_k^{-1} (w_k - h(u_k)) \quad (6.240)$$

With this in hand, the procedure is nearly identical to 3d-var:

1. Integrate your model forward to obtain $\{u_k\}$
2. Evaluate each of the $D_M^T(u_{k-1:0})$ and $D_h(u_k)$.
3. Using these values, compute $\nabla J_m(u)$
4. Set $u_0^{(new)} = u_0^{(prev)} - \eta \nabla J(u_0^{(prev)})$
5. Stop when $|u_0^{(new)} - u_0^{(prev)}|$ converges to your desired tolerance.

You can of course substitute another optimization scheme after step 3.

There are a few further observations that we can make about this procedure. The first, is that should we feel our model $f(\cdot)$ is not sufficient to fit the entire time series without the periodic corrections of the Kalman filter (perhaps there is some missing physics we haven't captured in our model), then it may be ideal to restrict the optimization to a subset of the full time series of observations. The second observation is that the above formulation has assumed the adjoint D_M^T can be computed easily. We can take advantage of the methods from the *Adjoint Sensativity Analysis* of ODEs to take advantage of modern methods for automatic differentiation of ODE solutions. Finally, it should be noted that the $4d - var$ algorithm yields the optimal starting conditions assuming our model is perfect. At the sacrifice of smoothness, it is often desirable to first obtain the ideal initial condition with $4d - var$ and *then* employ a filtering method like the EKF to obtain the final analysis.

CHAPTER 7
AN AUTONOMOUS ROBOTIC TEAM FOR THE RAPID
CHARACTERIZATION OF NOVEL ENVIRONMENTS

- 7.1 Rapid Georectification and Processing of Pushbroom Hyperspectral Imagery**
- 7.2 Supervised Regression with Uncertainty Quantification**
- 7.3 Methods for Unsupervised Classification of Hyperspectral Scenes in Novel Environments**

CHAPTER 8

A CHEMICAL DATA ASSIMILATION FRAMEWORK FOR INDOOR AIR QUALITY

- 8.1 Characterization of Photolysis**
- 8.2 HEART Chamber Sensing System**
- 8.3 Chemical Data Assimilation**
- 8.4 An Evaluation of Photocatalytic Ionization**

CHAPTER 9

A DISTRIBUTED NETWORK OF LOW COST AIR QUALITY SENSORS

9.1 MINTS Air Quality Network

9.1.1 Central Nodes

9.1.2 LoRa Nodes

9.1.3 MQTT

9.2 Real Time Dashboards via Containerization

9.3 Making Data Publicly Accessible via S3 and the Open Storage Network

discuss Rclone and using OSN to make data highly available

CHAPTER 10

TIME SERIES METHODS FOR AIR QUALITY

10.1 Time-Series Methods for Uncertainty Quantification

10.1.1 Metrics for Representative Uncertainty of combined Pseudo-Observations

10.1.2 Uncertainty Quantification with Temporal Variograms

10.2 Physics Informed modeling techniques for Air Quality Data

10.2.1 HAVOK + SciML = Nonlinear Modeling with External Forcing

10.2.2 Hamiltonian Neural Networks

CHAPTER 11

LIMITATIONS AND FUTURE WORK

11.1 Robot Team

11.1.1 Super Resolution

Propose methods for spatial/temporal/spectral super-resolution. In particular, comment on upcoming satellite deployments like EnMAP and others which will provide hyperspectral imagery.

CHAPTER 12
CONCLUSIONS

UPDATE REQUIRED!!!

APPENDIX A

REPRODUCIBLE RESEARCH TECHNIQUES

UPDATE REQUIRED!

A.1 Environment Management in Julia

A.2 Version Control (git)

A.3 CI/CD with github workflows

discuss automated tests as well as automatic document generation here

A.4 Literate Programming and Automatic Documentation with Quarto

A.5 Containerization with Docker and Docker Compose

Discuss NodeRed, InfluxDB, Grafana, etc...

APPENDIX B

OPTIMIZATION METHODS

Describe standard gradient descent and its utility for machine learning. Expand to briefly describe the extensions used in our work:

- Gradient Descent
- Gradient Descent with Momentum
- ADAM
- BFGS
- LBFGS
- The method we used for the variogram method that works specifically for quadratic loss functions...

We should also comment on which particular methods are best (and when)

APPENDIX C

HIGH PERFORMANCE COMPUTING

Provide an overview of relevant concepts in high performance computing i.e.

- slurm
- multi-threading
- parallelization (distributed computing)
- Memory management (i.e. preallocating data containers and writing functions that mutate, not allocate)

APPENDIX D

CHEMICAL REACTION MECHANISMS

UPDATE REQUIRED!

D.1 Simple Ion Mechanism

Here we can include the automatically generated documentation for the Ion Mechanism and any tohers

D.2 Master Chemical Mechanism

REFERENCES

- Ahmed, S. E., S. Pawar, and O. San (2020). Pyda: A hands-on introduction to dynamical data assimilation with python. *Fluids* 5(4), 225.
- Brook, R. D., S. Rajagopalan, C. A. Pope III, J. R. Brook, A. Bhatnagar, A. V. Diez-Roux, F. Holguin, Y. Hong, R. V. Luepker, M. A. Mittleman, et al. (2010). Particulate matter air pollution and cardiovascular disease. *Circulation* 121(21), 2331–2378.
- Buyantuyev, A. and J. Wu (2017). Remote sensing applications for land cover and land-use transformations in semiarid and arid environments. *Journal of Arid Environments* 140, 1–5.
- Carleo, G., K. Choo, J. Hofmann, E. Huang, C. Hughes, M. Hush, R. Iten, J. McClean, C. Miles, J. Preskill, et al. (2019). Machine learning and the physical sciences. *Reviews of Modern Physics* 91(4), 045002.
- Centre for Ecology and Hydrology (2017). Ecological sensing: a revolution in biodiversity monitoring.
- Chen, J., X. Shi, X. Li, M. Wang, W. Shen, and Y. Liu (2019). A review of air quality modeling: From gas-phase to particulate matter. *Advances in Atmospheric Sciences* 36(10), 921–947.
- Clark, M. P., W. N. Adger, S. Dessai, M. Goulden, D. W. Cash, and R. a. N. M. a. A. Stern, Nicholas and Gonzalez (2016). Urbanization, climate change and economic growth: Challenges and opportunities for policy makers. *Science of the Total Environment* 557-558, 279–291.
- Costello, A., M. Abbas, A. Allen, S. Ball, S. Bell, R. Bellamy, S. Friel, N. Groce, A. Johnson, M. Kett, M. Lee, C. Levy, M. Maslin, D. McCoy, B. McGuire, H. Montgomery, D. Napier, C. Pagel, J. Patel, J. A. P. de Oliveira, N. Redclift, H. Rees, D. Rogger, J. Scott, J. Stephenson, J. Twigg, J. Wolff, and C. Patterson (2009). Managing the health effects of climate change: Lancet and university college london institute for global health commission. *The Lancet* 373(9676), 1693–1733.
- DeLucia, E. H., N. Gomez-Casanovas, S. P. Long, M. A. Mayes, R. A. Montgomery, W. J. Parton, W. J. Sacks, J. P. Schimel, J. Verfaillie, and W. L. Silver (2021). The missing soil n: detecting processes driving soil nitrogen storage in complex ecosystems. *Journal of Ecology* 109(2), 447–459.
- Edenhofer, O., R. Pichs-Madruga, Y. Sokona, E. Farahani, S. Kadner, K. Seyboth, A. Adler, I. Baum, S. Brunner, P. Eickemeier, et al. (2014). *Climate change 2014: Mitigation of climate change*. Cambridge University Press.

- Field, C., V. Barros, D. Dokken, K. Mach, M. Mastrandrea, T. Bilir, M. Chatterjee, K. Ebi, Y. Estrada, R. Genova, et al. (2014). *Climate change 2014: Impacts, adaptation, and vulnerability. Part A: Global and sectoral aspects*. Cambridge University Press.
- Friedlingstein, P., M. W. Jones, M. O’Sullivan, R. M. Andrew, J. Hauck, G. P. Peters, W. Peters, J. Pongratz, S. Sitch, C. Le Quéré, et al. (2020). Global carbon budget 2020. *Earth System Science Data* 12(4), 3269–3340.
- Gamon, J. A., K. F. Huemmrich, R. S. Stone, and C. E. Tweedie (2016). Spatial and temporal variation in primary productivity (ndvi) of coastal alaskan tundra: Decreased vegetation growth following earlier snowmelt. *Remote Sensing of Environment* 175, 233–242.
- Haines, A., R. S. Kovats, D. Campbell-Lendrum, and C. Corvalán (2006). Climate change and human health: Impacts, vulnerability and public health. *Public Health* 120(7), 585–596.
- Hantson, S., A. Arneth, S. P. Harrison, D. I. Kelley, I. C. Prentice, S. S. Rabin, S. Archibald, F. Mouillot, S. R. Arnold, P. Artaxo, et al. (2016). The status and challenge of global fire modelling. *Biogeosciences* 13(11), 3359–3375.
- Houghton, J., Y. Ding, D. Griggs, M. Noguer, P. van der Linden, X. Dai, K. Maskell, and C. Johnson (2001). *Climate change 2001: The scientific basis*. Cambridge University Press.
- Houghton, J., G. Jenkins, and J. Ephraums (1990). *Climate change: The IPCC scientific assessment*. Cambridge University Press.
- Houghton, J., L. Meira Filho, B. Callander, N. Harris, A. Kattenberg, and K. Maskell (1996). *Climate change 1995: The science of climate change*. Cambridge University Press.
- Huang, J., L. Yu, J. Guo, X. Guo, W. Wang, C. Liu, and D. Ji (2017). Assessment of global surface energy budget datasets using flux tower observations. *Journal of Geophysical Research: Atmospheres* 122(14), 7452–7475.
- Kelly, F. J. and J. C. Fussell (2011). Air pollution and public health: emerging hazards and improved understanding of risk. *Environmental Geochemistry and Health* 33(4), 363–373.
- Kohonen, T. (1982). Self-organized formation of topologically correct feature maps. *Biological cybernetics* 43(1), 59–69.
- Masson-Delmotte, V., P. Zhai, H.-O. Pörtner, D. Roberts, J. Skea, P. Shukla, A. Pirani, W. Moufouma-Okia, C. Péan, R. Pidcock, et al. (2018). *Global warming of 1.5°C. An IPCC special report on the impacts of global warming of 1.5°C above pre-industrial levels and related global greenhouse gas emission pathways, in the context of strengthening the global response to the threat of climate change, sustainable development, and efforts to eradicate poverty*. Intergovernmental Panel on Climate Change.

- Metz, B., O. Davidson, P. Bosch, R. Dave, and L. Meyer (2007). *Climate change 2007: Mitigation of climate change*. Cambridge University Press.
- National Research Council (2010). Verifying greenhouse gas emissions: Methods to support international climate agreements.
- Oleson, K., D. Lawrence, G. Bonan, and M. Flanner (2013). Interactions between land use change and carbon cycle feedbacks. *Global Biogeochemical Cycles* 27(4), 972–983.
- Parry, M., O. Canziani, J. Palutikof, P. van der Linden, and C. Hanson (2007). *Climate change 2007: Impacts, adaptation and vulnerability*. Cambridge University Press.
- Pasher, J., B.-J. Park, J. Théau, F. Pimont, and S. Goetz (2019). Remote sensing of wetlands: An overview and practical guide. *Wetlands Ecology and Management* 27(2), 129–147.
- Rackauckas, C., D. Kelly, Q. Nie, J. Li, C. Warner, M. Dhairya, J. Fang, E. Zhou, R. Supekar, S. Sandhu, et al. (2020). Universal differential equations for scientific machine learning. *arXiv preprint arXiv:2012.09345*.
- Raissi, M. and G. E. Karniadakis (2021). Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics* 378, 686–707.
- Raissi, M., P. Perdikaris, and G. E. Karniadakis (2019). Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics* 378, 686–707.
- Rasmussen, C. E., C. K. Williams, et al. (2006). *Gaussian processes for machine learning*, Volume 1. Springer.
- Solomon, S., D. Qin, M. Manning, Z. Chen, M. Marquis, K. Averyt, M. Tignor, and H. Miller Jr (2007). *Climate change 2007: The physical science basis*. Cambridge University Press.
- Stocker, T., D. Qin, G.-K. Plattner, M. Tignor, S. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex, and P. Midgley (2013). *Climate change 2013: The physical science basis*. Cambridge University Press.
- Thenkabail, P. S. (2019). Remote sensing of global croplands for food security. *Remote Sensing* 11(10), 1261.
- United Nations (2015). Transforming our world: the 2030 agenda for sustainable development. *UN General Assembly*.
- United Nations Environment Programme (2017). Adaptation gap report 2017.

- Wang, D., D. Xie, Y. Xie, and C. Chen (2017). Remote sensing applications for urban water resources: A review. *Remote Sensing* 9(8), 829.
- World Health Organization (2018). Climate change and health.
- Wu, J. and X. Zhang (2021). A review on physics-informed machine learning: Basic principles, recent developments and future directions. *Physics Reports* 903, 1–45.
- Xu, X., F. Deng, X. Guo, P. Lv, H. Zhong, Y. Hao, G. Hu, J. Huang, Y. Guo, Y. Liu, et al. (2017). Association between particulate matter air pollution and hospital admissions in patients with chronic obstructive pulmonary disease in beijing, china. *Science of the Total Environment* 579, 1616–1621.

BIOGRAPHICAL SKETCH

UPDATE REQUIRED!!!

CURRICULUM VITAE

UPDATE REQUIRED!