

# Notes on statistics

Justin L. Ripley  
ripley@illinois.edu

September 16, 2023

# Contents

<b>1</b>	<b>Overview of parametric Bayesian statistics</b>	<b>2</b>
1.1	Definitions . . . . .	2
1.2	Parameter estimation . . . . .	3
1.3	Model selection/Hypothesis testing . . . . .	4
<b>2</b>	<b>Model comparison</b>	<b>5</b>
2.1	Bayes factor . . . . .	5
2.2	Nested models and the Savage-Dickey ratio . . . . .	6
2.3	Occam factor . . . . .	6
<b>3</b>	<b>Forecasting</b>	<b>8</b>
3.1	Injection analysis . . . . .	8
3.2	Fisher forecasting . . . . .	9
<b>4</b>	<b>Times series analysis</b>	<b>11</b>
4.1	Basic definitions . . . . .	11
4.2	Correlation and covariance . . . . .	11
4.3	Stationary and weak-sense stationary stochastic processes . . . . .	12
4.4	Wiener-Khinchin theorem . . . . .	13
4.5	Gaussian white noise . . . . .	14

4.6	Likelihood function for a series of measurements with colored stationary noise	14
4.7	Matched filter theorem . . . . .	17
<b>5</b>	<b>Numerical integration</b>	<b>19</b>
5.1	Monte Carlo integration . . . . .	21
5.2	Markov chain Monte Carlo (MCMC) . . . . .	22
5.3	Nested sampling . . . . .	24
<b>6</b>	<b>Numerical optimization</b>	<b>28</b>
6.1	Convex functions . . . . .	29
6.2	Gradient descent . . . . .	29
6.3	Newton's method . . . . .	30
<b>A</b>	<b>Probability theory</b>	<b>32</b>
A.1	Conditional probability and Bayes theorem . . . . .	32
A.2	Change of variables . . . . .	32
A.3	Expectation and covariance . . . . .	33
A.4	Characteristic/moment generating function . . . . .	34
A.5	Central limit theorem . . . . .	35
A.6	Fisher information and the Bernstein–von Mises theorem . . . . .	36
A.7	Fisher information and the Cramér–Rao bound . . . . .	38
<b>B</b>	<b>Fourier and other transforms</b>	<b>40</b>
B.1	Brief review of complex analysis . . . . .	40
B.2	The Fourier transform . . . . .	40
B.3	The Laplace transform . . . . .	41
B.4	The Hilbert transform . . . . .	41

<b>C Stationary phase approximation</b>	<b>42</b>
C.1 Stationary phase approximation . . . . .	42

## Abstract

I briefly review some of the basic notions of parametric Bayesian statistical inference that have come up in my research. These notes are not self contained; I assume some familiarity with probability theory and statistics, on the level of the first few chapters of [Was10]. I may sometimes implicitly assume some knowledge of differential geometry and partial differential equations. I try to cite sources whenever possible (whenever I can remember the source I learned something from), although the purpose of these notes are to serve more as a statistics “cheat sheet” than a formal review. If you think I am missing a reference please let me know. These notes are a work in progress, and they likely contain errors. Please let me know if you find any, or if you find any section unclear!

The notation is: vectors/tensors are in **boldfont**. Indices are denoted with lower case latin letters, e.g. the  $i^{th}$  component of the vector  $\boldsymbol{v}$  is  $(\boldsymbol{v})_i = v_i$ . We typically do not use boldfont when we explicitly write down indices. Repeated indices are summed over. Capital  $P$  always represents a probability distribution,  $\boldsymbol{x}$  always represents an instantiation of measured data,  $\boldsymbol{\theta}$  always represents model parameters. More generally, model parameters are represented by greek letters, while data is represented by latin letters. Random variables are always capitalized. Partial derivatives are denoted by  $\partial$ , and covariant derivatives by  $\nabla_i$  (for our purposes, you can usually replace covariant derivatives with partial derivatives).

I thank Rohit Chandramouli for helpful conversations, and a lecture on model selection that inspired the creation of these notes, and Simone Mezzasoma for helpful comments that have led to a clearer presentation.

Copyright 2023 Justin Ripley. You may copy and distribute this document provided that you make no changes to it.

# Chapter 1

## Overview of parametric Bayesian statistics

### 1.1 Definitions

We use lower case latin letters to index vector/tensor components. We use a lower case latin letter in parenthesis to index a particular vector/tensor. We also bold font vectors. Repeated indices are summed (Einstein summation notation). We denote models with capital Latin letters, model parameters with lower case greek letters, and data with lower case latin letters. Notice that we use latin indices to index both model parameters and data with lower case latin indices, even though in general model parameters and data will live in different dimensional vector spaces. We will drop the instantiation index (the latin index in parenthesis) unless otherwise needed.

Here we focus on **parametric Bayesian statistics**. By parametric, we mean that we have explicit functional models for the probability distributions of parameters, and by Bayesian, we mean we mean that we are interested in the probability distribution of those parameters (and/or models), given the observed data.

Bayes theorem gives us

$$P(\boldsymbol{\theta}|\mathbf{x}, M) = \frac{P(\mathbf{x}|\boldsymbol{\theta}, M) P(\boldsymbol{\theta}, M)}{P(\mathbf{x}, M)}, \quad (1.1)$$

Here  $P(\mathbf{x}|\boldsymbol{\theta}, M)$  is a statistical model  $M$  that reflects our beliefs about the data  $\mathbf{x}$  given the values of the parameters  $\boldsymbol{\theta}$  of a model  $M$ . The **posterior**  $P(\boldsymbol{\theta}|\mathbf{x}, M)$  is a probability distribution for the model parameters  $\boldsymbol{\theta}$  given  $\mathbf{x}$ . The **likelihood function** is  $P(\mathbf{x}|\boldsymbol{\theta}, M)$ , and is denoted by  $\mathcal{L}(\boldsymbol{\theta}, M)$ . The **prior distribution**  $P(\boldsymbol{\theta}, M)$  quantifies our certainty of the model parameters  $\boldsymbol{\theta}$  before we see the current data, and is often denote by  $\pi(\boldsymbol{\theta}, M)$ . The **evidence** [Ski06] (or marginal distribution of  $\mathbf{x}$  [Was10])  $P(\mathbf{x}, M)$  essentially acts as a normalizing constant, as  $P(\boldsymbol{\theta}|\mathbf{x}, M)$  must sum (integrate) to one. The evidence is often

denoted by  $\mathcal{Z}(\mathbf{x}, M)$ . If there are  $N$  independent observations of the data  $\mathbf{x}$ , the likelihood is

$$\mathcal{L}(\boldsymbol{\theta}, M) = \prod_{n=1}^N P(\mathbf{x}_{(n)}|\boldsymbol{\theta}, M). \quad (1.2)$$

We can write the evidence as the integral (or sum) over the model parameter values

$$\mathcal{Z}(\mathbf{x}, M) = \int d\boldsymbol{\theta} \mathcal{L}(\boldsymbol{\theta}, M) \pi(\boldsymbol{\theta}). \quad (1.3)$$

Much of applied Bayesian statistics centers around finding efficient ways to evaluate the likelihood and evidence, given an assumed model  $P(\mathbf{x}|\boldsymbol{\theta}, M)$  and prior  $P(\boldsymbol{\theta}, M)$ .

## 1.2 Parameter estimation

Assume you have one fixed model  $M$ . You can find the distribution of the parameters for the model, given a set of observed data, using Bayes theorem. Rewriting (5.1), we have

$$P(\boldsymbol{\theta}|\mathbf{x}, M) = \frac{\mathcal{L}(\boldsymbol{\theta}, M) \pi(\boldsymbol{\theta})}{\mathcal{Z}(\mathbf{x}, M)}. \quad (1.4)$$

The posterior probability distribution  $P(\boldsymbol{\theta}|\mathbf{x}, M)$  for most problems is complicated and cannot be written in closed form. Determining the posterior can usually only be accomplished numerically. Additionally it can be computationally expensive to compute the posterior distribution, especially if there are many parameters in the model ( $\boldsymbol{\theta}$  has many components).

This being said, it is straightforward to compute the relative probability of two different values of parameters  $\boldsymbol{\theta}_{(n)}$  and  $\boldsymbol{\theta}_{(m)}$ . We have

$$\frac{P(\boldsymbol{\theta}_{(n)}|\mathbf{x}, M)}{P(\boldsymbol{\theta}_{(m)}|\mathbf{x}, M)} = \frac{\mathcal{L}(\boldsymbol{\theta}_{(n)}, M) \pi(\boldsymbol{\theta}_{(n)}, M)}{\mathcal{L}(\boldsymbol{\theta}_{(m)}, M) \pi(\boldsymbol{\theta}_{(m)}, M)}. \quad (1.5)$$

We can write this in terms of the **likelihood ratio**

$$\lambda(\boldsymbol{\theta}_{(n)}, \boldsymbol{\theta}_{(m)}) \equiv \frac{\mathcal{L}(\boldsymbol{\theta}_{(n)}, M)}{\mathcal{L}(\boldsymbol{\theta}_{(m)}, M)}, \quad (1.6)$$

and the **prior odds**

$$R(\boldsymbol{\theta}_{(n)}, \boldsymbol{\theta}_{(m)}) \equiv \frac{\pi(\boldsymbol{\theta}_{(n)}, M)}{\pi(\boldsymbol{\theta}_{(m)}, M)}. \quad (1.7)$$

We discuss computational methods later, but we note that the value of  $\boldsymbol{\theta}$  that maximizes  $\mathcal{L}(\boldsymbol{\theta})$  is the **maximum likelihood estimator (MLE)**, and the value of  $\boldsymbol{\theta}$  that maximizes  $\mathcal{L}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta})$  is the **maximum a posteriori probability estimator (MAP)**. Note that the MLE and MAP do not give us any knowledge of the variance of those parameters—that requires knowledge of the full posterior probability distribution.

### 1.3 Model selection/Hypothesis testing

Second, you could have a collection of models  $M_{(1)}, \dots, M_{(N)}$ . Given a set of observations, you may be interested in the relative ability of each model to explain the data. Using Bayes' theorem, we have

$$P(M_{(n)}|\mathbf{x}) = \frac{P(\mathbf{x}|M_{(n)}) P(M_{(n)})}{P(\mathbf{x})}. \quad (1.8)$$

This means that

$$\frac{P(M_{(n)}|\mathbf{x})}{P(M_{(m)}|\mathbf{x})} = \frac{P(\mathbf{x}|M_{(n)})}{P(\mathbf{x}|M_{(m)})} \frac{P(M_{(n)})}{P(M_{(m)})}. \quad (1.9)$$

Notice that we have essentially marginalized over the parameters of the models. That is, we have

$$P(\mathbf{x}|M_{(n)}) = \int d\theta P(\mathbf{x}|\boldsymbol{\theta}, M_{(n)}) P(\boldsymbol{\theta}) = \mathcal{Z}(\mathbf{x}, M_{(n)}). \quad (1.10)$$

We see that the odds ratio for two models is given by the ratio of the evidence for each model multiplied by the prior odds for each model.

$$\frac{P(M_{(n)}|\mathbf{x})}{P(M_{(m)}|\mathbf{x})} = \frac{\mathcal{Z}(\mathbf{x}, M_{(n)})}{\mathcal{Z}(\mathbf{x}, M_{(m)})} \frac{P(M_{(n)})}{P(M_{(m)})}. \quad (1.11)$$

The odds ratio of the evidence is called the **Bayes factor**

$$B(M_{(n)}, M_{(m)}) \equiv \frac{\mathcal{Z}(\mathbf{x}, M_{(n)})}{\mathcal{Z}(\mathbf{x}, M_{(m)})}. \quad (1.12)$$



# Chapter 2

## Model comparison

### 2.1 Bayes factor

To briefly review, in **parameter estimation**, one finds the best fit parameters from the data given a model  $h(\theta)$ . What “best fit” means depends on the test statistic being used. Here we are concerned **Model selection**, which concerns finding which model better fits the data. In order to find the better fitting model, we compute the **Bayes factor**, which is the ratio of the evidence for each model

$$B_{2,1} \equiv \frac{P(d|H_2)}{P(d|H_1)}. \quad (2.1)$$

As a basic rule of thumb, if  $B_{2,1} \sim 1$ , then neither hypothesis is preferred compared to the other. If  $B_{2,1} \ll 1$ , then model 1 is preferred, while if  $B_{2,1} \gg 1$ , then model 2 is preferred. There are several subtleties to this interpretation, which we discuss more below.

If a model  $H$  has parameters  $\theta$ , we can compute the likelihood by marginalizing over the model's parameters for the likelihood (c.f. (1.3))

$$P(d|H) = \int d\theta P(d|\theta, H) P(\theta, H). \quad (2.2)$$

Doing this integral is typically challenging, since the dimension of the parameter space is very large, and the likelihood  $P(d|\theta, H)$  can be complicated (its functional form can only be guessed at in general). There are various approximations for how to compute this integral (analytically and numerically).

## 2.2 Nested models and the Savage-Dickey ratio

We consider a method to compute the Bayes factor for nested models. Consider a model  $M_1$  which is nested in a model  $M_2$ . The model  $M_2$  has one more parameter than  $M_1$  (generalizing to more parameters is straightforward). We call the extra parameter  $\lambda$ . We call the rest of the parameters  $\boldsymbol{\theta}$  nuisance parameters, as they do not distinguish the two models. In this setup we have that

$$P(d|\boldsymbol{\theta}, M_1) = P(d|\boldsymbol{\theta}, \lambda = \lambda_0, M_2), \quad (2.3)$$

where  $\lambda_0$  is a constant.

The evidence of  $M_1$  is

$$\begin{aligned} P(d|M_1) &= P(d|\lambda = \lambda_0, M_2) \\ &= \frac{P(\lambda = \lambda_0|d, M_2) P(d|M_2)}{P(\lambda = \lambda_0|M_2)}. \end{aligned} \quad (2.4)$$

We then see that the Bayes factor is

$$\begin{aligned} B_{2,1} &= \frac{P(d|H_2)}{P(d|H_1)} \\ &= \frac{P(\lambda = \lambda_0|M_2)}{P(\lambda = \lambda_0|d, M_2)}. \end{aligned} \quad (2.5)$$

This is the **Savage-Dickey ratio** [DL70]. We can also write this as

$$B_{2,1} = \left( \frac{\text{prior}}{\text{posterior}} \right)_{\lambda=\lambda_0}. \quad (2.6)$$

The advantage of this method is that you only need to compute the evidence of the model  $M_2$ , instead of computing the evidence of both  $M_2$  and the nested model  $M_1$ . Also, you do not need to divide two noisy numbers (the evidence of model 1 and model 2), you only need to divide a known number (the prior) by one noisy number (the evidence of model 2).

## 2.3 Occam factor

For more discussion see for example [Mac03]. We consider another measure of the power of a model to explain a given data set. The **Occam factor** is defined to be

$$O \equiv \frac{\text{posterior volume}}{\text{prior volume}} \sim \frac{\sigma_{\boldsymbol{\theta}|d}}{\sigma_{\boldsymbol{\theta}}}. \quad (2.7)$$

By volume, we mean the integral over parameter space of the probability distribution. Here  $\sigma_{\boldsymbol{\theta}}$  is some measure of the variance of the prior probability distribution, and  $\sigma_{\boldsymbol{\theta}|d}$  is variance

of the posterior probability distribution. The Occam factor measures how much the data shrinks the probability distribution as compared to its prior distribution. If the Occam factor is  $\sim 1$ , the data doesn't constrain the model well, since the variance parameters of the model do not shrink. We can interpret this as saying that the model does not explain the observed data well either. We can write

$$\text{evidence} \sim \text{max likelihood} \times \text{occam factor}. \quad (2.8)$$

From this, we see that the Occam factor accounts for the fact that models with more parameters can fit data better, and should be penalized for having more parameters.

For example, consider two hypothesis:  $H_1$  and  $H_2$ . Say they are nested:  $H_2(\boldsymbol{\theta}) \sim H_1(\boldsymbol{\theta}, \lambda)$ . If  $\lambda$  is unconstrained,  $O \sim 1$ , and if  $\lambda$  is well-constrained,  $O \ll 1$ .

# Chapter 3

## Forecasting

When it is hard or expensive to collect data, it can be useful to predict (or forecast) how well parameters of a given model could be measured with simulated data. Forecasting can inform whether a more in-depth analysis of a model on real data is worth doing—that is whether or not real data could place any meaningful measurement of the parameters of a model. Here we review several semi-analytic methods for forecasting.

### 3.1 Injection analysis

We consider a model  $P(\mathbf{x}|\boldsymbol{\theta})$  with prior  $P(\boldsymbol{\theta})$ . For whatever reason, we do not have any data  $\mathbf{x}$ . For example, it may be expensive to collect data, so we do not want to collect it until we have some confidence that we could meaningfully measure parameters in the model  $P(\mathbf{x}|\boldsymbol{\theta})$ . An injection analysis involves determining the distribution of  $P(\boldsymbol{\theta}_i|\mathbf{x}_0)$  where  $\mathbf{x}_0$  is fake (generated) data set that we hope represents a characteristic realization of the data we may measure. In other words, we have “injected” the data  $\mathbf{x}_0$  into our model. If we can meaningfully measure/determine the parameters  $\boldsymbol{\theta}$  given  $\mathbf{x}_0$ , it may be worth collecting real data/observations.

A reasonable choice for  $\mathbf{x}_0$  is to choose  $\mathbf{x}_0$  to maximize  $P(\mathbf{x}|\boldsymbol{\theta}_0)$ , where  $\boldsymbol{\theta}_0$  are the values of parameters that you expect to hope to measure. In equations we choose  $\mathbf{x}_0$  to satisfy

$$\forall \mathbf{x} \ P(\mathbf{x}|\boldsymbol{\theta}_0) \leq P(\mathbf{x}_0|\boldsymbol{\theta}_0) \quad (3.1)$$

Sometimes it is worth adding a realization of noise,  $\mathbf{n}$ , to  $\mathbf{x}_0$ ; we call this

$$\mathbf{x}_{0,n} = \mathbf{x}_0 + \mathbf{n} \quad (3.2)$$

For example, the components of  $\mathbf{n}$  may be drawn from a Gaussian with zero mean and unit diagonal covariance matrix (although the choice of  $\mathbf{n}$  will depend on your understanding nature of the experiment/observation). It is common to  $\mathbf{n} = 0$ , which can be considered the

“best” possible situation for recovering parameters. We then inject  $\mathbf{x}_{0,n}$  into the likelihood, and sample on  $\boldsymbol{\theta}$ , that is we consider

$$P(\boldsymbol{\theta}|\mathbf{x}_{0,n}) = \frac{P(\mathbf{x}_{0,n}|\boldsymbol{\theta}) P(\boldsymbol{\theta})}{P(\mathbf{x}_{0,n})}. \quad (3.3)$$

The probability distribution  $P(\boldsymbol{\theta}|\mathbf{x}_{0,n})$  gives us an understanding of how well we could measure  $\boldsymbol{\theta}_0$ , given (near, if  $n \neq 0$ ) optimal data. Moreover, it can give us an idea of how the different components of  $\boldsymbol{\theta}$  may be correlated with one another.

To determine  $P(\boldsymbol{\theta}|\mathbf{x})$ , we either need to directly sample  $\boldsymbol{\theta}$  from (3.3), or use an approximation method.

## 3.2 Fisher forecasting

Injection analysis with a multivariate normal approximation for the likelihood is called **Fisher forecasting**. Let  $\hat{\boldsymbol{\theta}}$  be the maximum likelihood estimator for  $\boldsymbol{\theta}$ . Under appropriate regularity conditions, in the limit of a large number of observations,  $P(\mathbf{x}|\hat{\boldsymbol{\theta}})$  tends towards the normal distribution (in parameter space), with mean  $\hat{\boldsymbol{\theta}}$  and covariance matrix

$$\Sigma_{ij}^{(F)} = \frac{1}{N} F_{ij}^{-1}(\hat{\boldsymbol{\theta}}). \quad (3.4)$$

Here  $N$  is the number of observations. This is known as the **Bernstein–von Mises theorem**. We provide a proof of this result in Appendix A.

Estimating the maximum likelihood estimator  $\hat{\boldsymbol{\theta}}$  is a challenging task of its own, that we do not explore further here. Using the Bernstein-von Mises theorem, the posterior probability distribution near  $\hat{\boldsymbol{\theta}}$ , in the limit of a large number of observations, is approximately

$$P(\boldsymbol{\theta}|\mathbf{x}) = \frac{\pi(\boldsymbol{\theta})}{\mathcal{Z}(\mathbf{x}, M)} \exp \left[ -\frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^i F_{ij} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^j \right]. \quad (3.5)$$

If we assume a Gaussian prior on the parameters  $\boldsymbol{\theta}$ , then the posterior is a multivariate Gaussian with an inverse covariance matrix given by

$$\Sigma_{ij}^{-1} = N F_{ij} + \frac{1}{\sigma_i^2} \delta_{ij}. \quad (3.6)$$

That is, the posterior probability distribution within this approximation is

$$P(\boldsymbol{\theta}|\mathbf{x}) \approx \frac{1}{\sqrt{(2\pi)^k \det \Sigma}} \exp \left[ -\frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^i \Sigma_{ij}^{-1} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^j \right], \quad (3.7)$$

where  $k$  is the dimensionality of  $\boldsymbol{\theta}$ , that is the number of parameters.

To perform a **Fisher forecast** for a given model  $P(\mathbf{x}|\boldsymbol{\theta})$ , we pick a set of parameters  $\boldsymbol{\theta}_0$ , and then compute the Fisher information matrix (A.22). That is, we assume that  $\boldsymbol{\theta}_0$  are the “true” model parameters, and also are the maximum likelihood estimators. We then “inject” those parameters into the likelihood, which we approximate as a multivariate Gaussian with inverse covariance matrix given by (3.6). This analysis can be useful to determine the strength of correlation between different the different components of  $\theta^i$  (through the off-diagonal terms in  $\Sigma_{ij}$ ). The diagonal of the covariance matrix additionally gives us the  $1-\sigma$  error bars of the parameters. If we could make  $N$  measurements of the same data, each element in the covariance matrix would decrease  $1/N$ , as follows from (3.4). We see that the Fisher matrix can also give us a rough estimate of the number of observations  $N$  that are needed to make a  $n - \sigma$  observation of a parameter  $\theta^i$ .

Fisher forecasting is sometimes said to provide an optimal estimate of the variance of the parameters in a given measurement. This statement is justified by the **Cramér-Rao bound**, which states that the covariance matrix of an unbiased estimator for  $\boldsymbol{\theta}$ ,  $\boldsymbol{\Theta}$ , (that is  $\mathbb{E}[\boldsymbol{\Theta}(\mathbf{x})] = \boldsymbol{\theta}$ ) is bounded from below by the inverse of the Fisher information matrix

$$\Sigma_{ij} \Big|_{\boldsymbol{\theta}=\boldsymbol{\mu}_{\boldsymbol{\theta}}} \geq F_{ij}^{-1}(\boldsymbol{\theta}). \quad (3.8)$$

This bound should be interpreted with caution though, as (A.35) only holds for unbiased estimators to the parameters  $\boldsymbol{\theta}$ . Consider a general estimator  $\boldsymbol{\Theta}(\mathbf{x})$ , and denote its expectation by

$$\mathbb{E}[\boldsymbol{\Theta}(\mathbf{x})] = \boldsymbol{\psi}(\boldsymbol{\theta}). \quad (3.9)$$

The Cramér-Rao bound states that

$$\nabla_{\theta_m} \psi_i \nabla_{\theta_n} \psi_j F_{mn}^{-1}(\boldsymbol{\theta}). \quad (3.10)$$

If  $\boldsymbol{\Theta}$  is an unbiased estimator ( $\boldsymbol{\psi} = \boldsymbol{\theta}$ ), then (A.34) reduces to (A.35). We outline a proof of (A.34) in Appendix A.

# Chapter 4

## Times series analysis

We consider the problem of determining the signal from a data timestream. Calling the data  $x(t)$ , we then want to find a signal  $s(t)$  given noise  $n(t)$ , where

$$x(t) = s(t) + n(t). \quad (4.1)$$

We model the noise  $n(t)$  can be modeled as a **stochastic process** (which implies that  $x(t)$  is a stochastic process) We assume  $x, s, n$  are all real functions. Our main goal is to derive the likelihood function for a time series of the form (4.1) when  $n$  takes the form of colored stationary noise, and to derive the matched filtering theorem.

### 4.1 Basic definitions

A function  $n(t)$  is a **stochastic process**, if  $n(t)$  is a random variable at each time  $t$  that is described by some probability distribution. This probability distribution may depend on  $t$ , and the previous history of values of  $x$ , for example. If  $t$  is a discrete variable, and the probability distribution for  $n(t_i)$  depends on  $n(t_{i-1})$ , then  $n(t_i)$  is a **Markov chain**. If the probability distribution for  $n(t)$  is independent of  $t$ , then  $n(t)$  is a **stationary process**.

### 4.2 Correlation and covariance

We denote the mean and variance of a time series  $x(t)$  with  $\mu_x$  and  $\sigma_x$ , respectively. We define the **covariance** between two stochastic processes  $x_1$  and  $x_2$  at times  $t_1$  and  $t_2$  to be

$$C_{x_1, x_2}(t_1, t_2) \equiv \mathbb{E} \left[ (x_1(t_1) - \mu_{x_1(t_1)}) (x_2^*(t_2) - \mu_{x_2(t_2)}^*) \right]. \quad (4.2)$$

We define the **autocovariance** for a stochastic process to be

$$K_x(t_1, t_2) \equiv C_{xx}(t_1, t_2). \quad (4.3)$$

We define the **autocorrelation** for a stochastic process to be

$$R_x(t_1, t_2) \equiv \mathbb{E}[x(t_1)x^*(t_2)] = K_x + \mu_x^2. \quad (4.4)$$

We define the **energy** of a time series  $s$  to be

$$E_x \equiv \int_{-\infty}^{\infty} dt |x(t)|^2 = \int_{-\infty}^{\infty} df |\tilde{x}(f)|^2. \quad (4.5)$$

The last expression follows from Parseval's theorem. We define the **energy spectral density** to be

$$\hat{S}_x(f) \equiv |\tilde{x}(f)|^2. \quad (4.6)$$

### 4.3 Stationary and weak-sense stationary stochastic processes

A stochastic process is said to be **stationary** if its joint probability distribution does not change under time shifts. A stochastic process is said to be **weak-sense stationary (WSS)** (or **wide-sense stationary**), then its first moment is independent of time, and if its autocorrelation function depends only on  $\tau = t_1 - t_2$ . For a WSS process we can write

$$R_x(\tau) = \mathbb{E}[x(t+\tau)x^*(t)], \quad (4.7)$$

where  $t$  is arbitrary. For a WSS/stationary process, we can write the expectation of  $x$  at a given instant as the average of  $x$  over all time

$$\mathbb{E}[x] = \lim_{T \rightarrow \infty} \frac{1}{T} \int_{-\infty}^{\infty} dt x_T(t), \quad (4.8)$$

and similarly for functions of  $x(t)$ . Here  $x_T(t)$  is defined to be

$$x_T(t) \equiv w_T(t) x(t), \quad (4.9)$$

$$w_T(t) \equiv \begin{cases} 1 & |t| < T/2 \\ 0 & \text{otherwise} \end{cases}. \quad (4.10)$$

In other words, we can write

$$R_x(\tau) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_{-\infty}^{\infty} dt x_T(t+\tau) x_T^*(t). \quad (4.11)$$

We emphasize that for WSS processes we can replace ensemble averages with time averages. This is extremely useful in practice, as we can then determine the statistical properties of WSS by taking repeated time measurements of observables of the time series. While it is



common to assume a given time stream is WSS, most real world data is at best approximately WSS, typically over a short time scale.

Stationary stochastic processes have support over the entire real line, so the energy integrals defined above typically diverge. For these processes, instead one looks at the **power**

$$\begin{aligned} P_x &\equiv \lim_{T \rightarrow \infty} \frac{1}{T} \int_{-\infty}^{\infty} dt |x_T(t)|^2 \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \int_{-\infty}^{\infty} df |\tilde{x}_T(f)|^2 \\ &= \int_{-\infty}^{\infty} df S_x(f). \end{aligned} \quad (4.12)$$

On the last line we have used the **power spectral density**, which is defined to be

$$S_x(f) \equiv \lim_{T \rightarrow \infty} \frac{1}{T} |\tilde{x}_T(f)|^2. \quad (4.13)$$

If  $x(t)$  is real, then  $\tilde{x}_T(-f) = \tilde{x}_T^*(f)$ , and it is common to define the power spectral density to be

$$S_x(f) \equiv \lim_{T \rightarrow \infty} \frac{2}{T} |\tilde{x}_T(f)|^2, \quad (4.14)$$

and to write

$$P_x = \int_0^{\infty} df S_x(f). \quad (4.15)$$

Finally, we consider the expectation value of the Fourier transform of a stationary signal

$$\begin{aligned} \mathbb{E}[\tilde{x}(f') \tilde{x}^*(f)] &= \int_{-\infty}^{\infty} dt \int_{-\infty}^{\infty} dt' e^{-2\pi i(f t - f' t')} \mathbb{E}[x(t) x(t')] \\ &= \int_{-\infty}^{\infty} dt e^{-2\pi i(f - f')t} \int_{-\infty}^{\infty} d\tau e^{-2\pi i f' \tau} \mathbb{E}[x(t + \tau) x(t)] \\ &= \int_{-\infty}^{\infty} dt e^{-2\pi i(f - f')t} \tilde{R}_x(f) \\ &= \delta(f - f') S_x(f). \end{aligned} \quad (4.16)$$

On the third line we used that  $x(t)$  was stationary.

## 4.4 Wiener-Khinchin theorem

The power spectral density and the Fourier transform of the autocorrelation are equal for WSS processes. This is known as the **Wiener-Khinchin theorem**. To prove this, we set  $\tau \equiv t_1 - t_2$ . We then write

$$\tilde{R}_x(f) = \int_{-\infty}^{\infty} d\tau e^{-2\pi i f \tau} R_x(\tau)$$

$$\begin{aligned}
&= \lim_{T \rightarrow \infty} \frac{1}{T} \int_{-\infty}^{\infty} d\tau e^{-2\pi i f \tau} \int_{-\infty}^{\infty} dt x_T(t + \tau) x_T^*(t) \\
&= \int_{-\infty}^{\infty} df' \int_{-\infty}^{\infty} d\tau e^{2\pi i (f' - f) \tau} \lim_{T \rightarrow \infty} \frac{1}{T} |\tilde{x}_T(f')|^2 \\
&= S_x(f).
\end{aligned} \tag{4.17}$$

On the second line we used the formula for the autocorrelation function for WSS processes (4.11). On the third line we used the convolution theorem for Fourier transforms. On the last line we used that  $\int d\tau e^{i\tau f} = \delta(f)$ , and used the definition of  $S_x(f)$  (4.13).

## 4.5 Gaussian white noise

As a special case of a stationary stochastic process, we first consider **Gaussian white noise**. By **Gaussian**, we mean that the probability distribution for  $x(t)$  for each  $t_i$  is a Gaussian

$$x(t_i) \sim \mathcal{N}(\mu_i, \sigma_i^2). \tag{4.18}$$

By **white**, we mean that the  $x(t_i)$  are uncorrelated, and that the means  $\mu_i = 0$ . The autocorrelation function then is

$$R_x(t_i, t_j) = \sigma^2 \delta_{ij}. \tag{4.19}$$

Notice that  $\sigma^2$  does not depend on  $t_i$ . We see that white noise is stationary. The autocorrelation function can be written in terms of  $\tau \equiv t_i - t_j$  as

$$R_x(\tau) = \sigma^2 \delta(\tau). \tag{4.20}$$

By the Wiener-Kinchin theorem, we can compute the power spectral density from the Fourier transform of the autocorrelation function

$$\begin{aligned}
S_x(f) &= \int_{-\infty}^{\infty} dt e^{2\pi i f t} R_x(t) \\
&= \sigma^2.
\end{aligned} \tag{4.21}$$

We see for Gaussian white noise, the power spectral density is a constant—there is constant power across all frequencies. For real functions, the integral over frequencies goes from  $[0, \infty)$ , and we define

$$S_x(f) = 2\sigma^2. \tag{4.22}$$

## 4.6 Likelihood function for a series of measurements with colored stationary noise

Our treatment roughly follows [CA11] (see also [Fin92]). We consider a series of a continuous time stream of observations  $y(t)$ . We assume that  $y(t)$  can be related to a convolution of a

time stream drawn from Gaussian white noise

$$y(t) = \int_{-\infty}^{\infty} dt' \gamma(t-t') x(t'). \quad (4.23)$$

Here  $\gamma$  is the kernel and  $x$  is a time stream drawn from a Gaussian distribution with constant. We assume  $y, x, \gamma$  are all real functions. The Fourier transform gives us

$$\tilde{y}(f) = \tilde{\gamma}(f) \tilde{x}(f). \quad (4.24)$$

The power spectral density of  $y$  then is

$$S_y(f) = |\tilde{\gamma}(f)|^2 S_x(f) = |\tilde{\gamma}(f)|^2 \sigma. \quad (4.25)$$

The main point of adding the convolution is that we can consider processes with **colored noise**, that noise where the power spectral density can vary with frequency. We can do this by choosing some  $\sigma$ , and then choosing a  $\gamma$  such that  $|\tilde{\gamma}(f)|^2$  gives us the spectral density we desire.

We note that  $y(t)$  describes a WSS process, as

$$\begin{aligned} \mathbb{E}[y(t_1) y(t_2)] &= \int_{-\infty}^{\infty} dt'_1 \int_{-\infty}^{\infty} dt'_2 \gamma(t_1 - t'_1) \gamma(t_2 - t'_2) \mathbb{E}[x(t'_1) x(t'_2)] \\ &= \sigma^2 \int_{-\infty}^{\infty} dt' \gamma(t_1 - t') \gamma(t_2 - t') \\ &= \sigma^2 \int_{-\infty}^{\infty} dt \gamma(t_1 - t) \gamma(t_2 - t). \end{aligned} \quad (4.26)$$

On the second line we used (4.20), while on the third we set  $t = t_2 - t'$ .

We consider a discretized set of  $N$  points (evenly spaced) from  $x(t)$  and  $y(t)$ . We write the vectors  $\mathbf{x}$  and  $\mathbf{y}$  where the componets are, e.g.  $\mathbf{x}_i \equiv x(t_i)$ . We define the discretized matrix  $\Gamma_{ij} \equiv \gamma(t_i - t_j)$ . We then have

$$y_i = \Gamma_{ij} x_j \quad (4.27)$$

We define the matrix

$$\Sigma_{ij} \equiv \frac{\sigma^2}{\Delta t} \delta_{ij}, \quad (4.28)$$

where  $\Delta t \equiv T/(N-1)$ , and  $N$  is the number of discretized points. We choose this scaling so that the autocorrelation of  $x_i$  approaches the correct behavior in the continuum limit, as we show below. The probability distribution for each  $\mathbf{x}$  is

$$P_{\mathbf{x}}(\mathbf{x}) = \left( \frac{1}{\sqrt{2\pi}} \right)^N \frac{1}{\sqrt{\det \Sigma}} \exp \left[ -\frac{1}{2} \mathbf{x}_i \Sigma_{ij}^{-1} \mathbf{x}_j \right]. \quad (4.29)$$

The correlation for  $x$  is then

$$\begin{aligned}
\mathbb{E}[x_i x_j] &= \int d^N x P_{\mathbf{X}}(\mathbf{x}) x_i x_j \\
&= \int d^N x \left( \frac{1}{\sqrt{2\pi}} \right)^N \frac{1}{\sqrt{\det \Sigma}} \exp \left[ -\frac{1}{2} x_i \Sigma_{ij}^{-1} x_j \right] x_i x_j \\
&= \sigma^2 \frac{\delta_{ij}}{\Delta t}.
\end{aligned} \tag{4.30}$$

In the limit  $\Delta t \rightarrow 0$ , this approaches  $\sigma^2 \delta(t_i - t_j)$ , which is the autocorrelation function for Gaussian white noise; see (4.20).

We obtain the probability distribution for  $\mathbf{y}$  under a linear transformation of variables (note the ordering of the indices)

$$P_{\mathbf{Y}}(\mathbf{y}) = \left( \frac{1}{\sqrt{2\pi}} \right)^N \frac{1}{\sqrt{\det \Sigma \det \Gamma}} \exp \left[ -\frac{1}{2} y_i \Gamma_{mi}^{-1} \Sigma_{mn}^{-1} \Gamma_{nj}^{-1} y_j \right]. \tag{4.31}$$

This expression gives the probability the  $N$  draws. We now need to take the continuum limit. First we look at the argument of the exponential

$$\begin{aligned}
y_i \Gamma_{mi}^{-1} \Sigma_{mn}^{-1} \Gamma_{nj}^{-1} y_j &= \frac{1}{\sigma^2} x_i x_j \Delta t \\
&\rightarrow \frac{2}{S_x} \int_{t_s}^{t_f} dt |x(t)|^2 \\
&\approx \frac{2}{S_x} \int_{-\infty}^{\infty} dt |x(t)|^2 \\
&= \frac{4}{S_x} \int_0^{\infty} df |\tilde{x}(f)|^2 \\
&= 4 \int_{-\infty}^{\infty} df \frac{|\tilde{y}(f)|^2}{S_y}.
\end{aligned} \tag{4.32}$$

Here  $t_s, t_f$  are the start and end times for the series  $x(t)$ . On the first line we used  $x_i = \Gamma_{ij}^{-1} x_j$ , and that  $\Sigma_{ij}^{-1} = \sigma^{-1} \delta_{ij}$ . On the second and third lines we converted the Riemann sum to an integral (we took the continuum limit). We approximated the start/end times with  $\pm\infty$ . On the last line we use  $\tilde{x} = \tilde{y}/\tilde{\gamma}$ ,  $S_y = |\tilde{\gamma}(f)|^2 S_x$ , and that  $S_x$  is a constant, so we can pull it into the integral. Remember that we assume that  $x, y, \gamma$  are all real, so that for example  $x(-f) = x^*(-f)$ . Ignoring the constant normalization factor, we see that the probability density function (the likelihood function) for  $y(t)$  is

$$P(y(t)) \propto \exp \left[ -\frac{1}{2} (y, y) \right], \tag{4.33}$$

where we have defined the inner product

$$(a, b) \equiv 2 \int_0^{\infty} df \frac{a(f) b^*(f) + b(f) a^*(f)}{S_y(f)}. \tag{4.34}$$

Here  $S_y(f)$  is the spectral noise density for the process, and  $a, b$  can represent the Fourier transform of particular draws. We interpret (4.33) as the likelihood function (up to normalization) for colored WSS noise. We call (4.34) a **matched filter**. We define the **signal to noise ratio (SNR)** for a signal  $s$  with noise  $n$  to be

$$\rho^2 \equiv (s, s) = 4 \int_0^\infty df \frac{|s(f)|^2}{S_n(f)}. \quad (4.35)$$

## 4.7 Matched filter theorem

We next derive the optimal test statistic for extracting a signal from WSS colored noise. We consider a time series  $x(t)$  that can be written as

$$x(t) = s(t) + n(t). \quad (4.36)$$

We assume that  $n(t)$  can be written as a convolution with Gaussian white noise, as we described in Sec. (4.6). We assume we are searching for a signal  $s(t)$  that we know how to compute. Under these assumptions, we can compute the optimal test statistic to distinguish between the two following hypothesis:

Null hypothesis  $\mathcal{H}_0$ :  $x(t) = n(t)$ .

Alternative hypothesis  $\mathcal{H}_1$ :  $x(t) = s_1(t) + n(t)$ .

Here  $s_1(t)$  is a signal we are guessing is in the data. We compare the two hypothesis by computing the likelihood ratio (likelihood for short for the rest of this section)

$$\Lambda(\mathcal{H}_1|x) \equiv \frac{P(x|\mathcal{H}_1)}{P(x|\mathcal{H}_0)}. \quad (4.37)$$

We use the likelihood function (4.33). We next show that  $s_1 \propto s$  maximizes  $\Lambda$ , which is the **matched filtering theorem**.

If the null hypothesis is true, then the probability density function goes as

$$P(x|\mathcal{H}_0) \propto \exp \left[ -\frac{1}{2} (x, x) \right]. \quad (4.38)$$

If the alternative hypothesis is true, then the probability density function goes as

$$P(x|\mathcal{H}_1) \propto \exp \left[ -\frac{1}{2} (x - s_1, x - s_1) \right]. \quad (4.39)$$

We have used (4.34), with the noise power spectral density given by  $S_n(f)$ . The normalization factors cancel out in the likelihood ratio, and we are left with

$$\Lambda(\mathcal{H}_1|x) = \exp \left[ -\frac{1}{2} (x - s_1, x - s_1) + \frac{1}{2} (x, x) \right]$$

$$= \exp \left[ (x, s_1) - \frac{1}{2} (s_1, s_1) \right]. \quad (4.40)$$

The matched filtering theorem states that the likelihood ratio  $\Lambda(\mathcal{H}_1|x)$  is maximized when  $s_1 \propto s$ . To show this, we first note that likelihood ratio is maximized when the log-likelihood ratio  $L$  is maximized. The log-likelihood is

$$L(s_1) \equiv (n, s_1) + (s, s_1) - \frac{1}{2} (s_1, s_1). \quad (4.41)$$

We only consider  $s_1$  such that  $(n, s_1) = 0$  (this also holds for the “true” signal  $s$ ). Moreover, we fix  $(s_1, s_1) = c_1$ , where  $c_1$  is a constant (otherwise the likelihood could be arbitrarily big or small by rescaling the amplitude of  $s_1$ ). Maximizing the likelihood then reduces to maximizing

$$L(s_1) = (s, s_1) - \frac{1}{2} c_1. \quad (4.42)$$

By the Cauchy-Schwartz inequality, we have that

$$(s, s_1) \leq \sqrt{(s, s)} \sqrt{(s_1, s_1)}. \quad (4.43)$$

Equality only holds when  $s_1 \propto s$ . We conclude that choosing  $s_1 \propto s$  maximizes the likelihood function (up to a proportionally constant, which is fixed by the condition  $(s_1, s_1) = c_1$ ).

The task of finding a signal in colored WSS noise then reduces to finding a filtering function  $s_1$  that is orthogonal to the noise, and that maximizes the value of the matched filter  $(x, s_1)$ . In practice, we can determine the noise profile of the detector by measuring the response of the detector in the (assumed) absence of any signal.

The matched filtering theorem is powerful, but it relies on several strong assumption that are only approximately met in practice. First, it assumes that we know what we are looking for—that is, that we have a **template bank** of templates  $s_i(t)$  that we can convolve with the data. Even if we do have a template bank, it can be very computationally expensive to search for the  $s_i$  that fits the data best, especially if the parameter space for  $s_i$  is large. Efficiently evaluating the likelihood and searching through parameter space remains a topic of active research in, e.g. the gravitational wave astronomy community (for a review, see e.g. [CA11]). The matched filtering theorem also assumes the noise is stationary or WSS. Most kinds of detectors (say a phone line, or a gravitational wave detector) suffer from non-stationary noise, often called **glitches**. Provided those are well enough understood, they can be subtracted out of the signal, although in practice it can be difficult to completely remove glitches from a time stream.

# Chapter 5

## Numerical integration

We we discussed in Sec. 1, in parametric Bayesian statistics our goal is to determine the posterior probability distribution of the parameters of the model under consideration, given a set of measured data, or to determine the total evidence for the model.

A large portion of computational, parametric Bayesian statistics essentially consists of determining ways to compute high dimensional integrals. To understand why we need to compute integrals in parametric Bayesian statistics we look again at Bayes theorem

$$P(\boldsymbol{\theta}|\mathbf{x}) = \frac{\mathcal{L}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta})}{\mathcal{Z}(\mathbf{x})}, \quad (5.1)$$

where  $\pi$  is the prior and the the likelihood  $\mathcal{L}$  and evidence  $\mathcal{Z}$  are

$$\mathcal{L}(\boldsymbol{\theta}) = \prod_{i=1}^N P(\mathbf{x}_i|\boldsymbol{\theta}), \quad (5.2)$$

$$\mathcal{Z}(\boldsymbol{\theta}) = \int d^k \theta \mathcal{L}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}). \quad (5.3)$$

We have assumed the measurements of  $\mathbf{x}$  have been taken idependently on one another in the equation for the likelihood. We assume that the parameters  $\boldsymbol{\theta}$  are continuous. Clearly if we want to determine  $P(\boldsymbol{\theta}|\mathbf{x})$  directly, we need to compute the evidence  $\mathcal{Z}$ , which requires integrating over the likelihood. Beyond this though, many summary statistics of practical interest require computing an integral. For example, we may be interested in the expectation of  $\boldsymbol{\theta}$  and its covariance matrix

$$\mu_i = \mathbb{E}[\theta_i] \equiv \int d^k \theta P(\boldsymbol{\theta}|\mathbf{x}) \theta_i, \quad (5.4a)$$

$$C_{ij} \equiv \mathbb{E}[(\theta_i - \mu_i)(\theta_j - \mu_j)] \equiv \int d^k \theta P(\boldsymbol{\theta}|\mathbf{x}) (\theta_i - \mu_i)(\theta_j - \mu_j), \quad (5.4b)$$

We cover methods to compute high dimensional integrals, as in many applications  $k \gg 1$  (or at least,  $k \gtrsim 10$ ). In this regime, it is usually computationally infeasible to compute (5.3)

using traditional deterministic methods such as the trapezoid rule or Gaussian quadrature. For example if  $k = 10$ , and if we have 10 quadrature points in each parameter direction, we will need to make  $N \gtrsim 10^{10}$  evaluations for a trapezoid rule approximation of the evidence. The likelihood function is often a highly complex function with sharp peaks, so many more than 10 grid points would be needed to resolve in each direction in order to properly resolve the posterior.

As far as I am aware, the most efficient way to compute high dimensional integrals is through stochastic/Monte Carlo methods. The fact that Monte Carlo methods are the best methods to compute many high dimensional integrals is somewhat surprising, as they have very slow rates of convergence. In general, the error of Monte Carlo integrals goes as  $N^{-1/2}$ , where  $N$  is the number of points used in the approximation. In one dimension, approximations as simple as the trapezoid rule converge to the correct answer as  $1/N^2$  (e.g. [PTVF92]). This being said, the accuracy of methods such as the trapezoid rule rapidly deteriorate at higher dimension, while for Monte Carlo methods, the accuracy decreases as  $N^{-1/2}$ , regardless of the dimensionality of the problem, although the proportionality constant to this decrease strongly depends on the choice of algorithm one uses, and the dimensionality of the problem. We only consider stochastic/Monte Carlo integration methods in this chapter.

In effect, Bayesian parametric statistics reduces statistics to probability theory, and many problems in probability theory can be reduced to problems in the integration of complicated functions in high dimensional spaces. There are three main approaches to integration, **Riemann integration**, **Riemannian-Stieltjes integration**, and **Lebesgue integration**.

**Monte Carlo integration** can be thought of as providing an approximation to the Riemannian integral. We review Monte Carlo integration in Sec. 5.1. The most commonly used variant of Monte Carlo integration is **Markov chain Monte Carlo** (MCMC) integration, which can be thought of as approximating the Riemann-Stieltjes integral. We review MCMC integration in Sec. 5.2. The Monte Carlo approximation of certain kinds of Lebesgue integrals goes under the name **Nested Sampling** (NS), which we review in Sec. 5.3. There are many excellent, long discussions of all these methods on the internet and elsewhere (e.g. [BGJM11, Ski06, HFM18]), so we only outline the main ideas.

Before continuing, we mention two applications where you do not need to compute an integral (and hence do not need to use the methods discussed here). If we only need to compute the ratio of the posterior for two parameters values  $\theta_1$  and  $\theta_2$ , we only need to determine  $P(\theta_1|\mathbf{x})/P(\theta_2|\mathbf{x}) = \mathcal{L}(\theta_1)/\mathcal{L}(\theta_2)$ , which does not involve any integrals. We also do not need to compute any integrals if we only want the maximum of the posterior or the likelihood (the maximum likelihood estimator). We review some maximization methods in Chptr. 6.



## 5.1 Monte Carlo integration

Consider a function  $f(\boldsymbol{\theta})$ , and an integral over the domain  $\Omega$

$$I = \int_{\Omega} d^k \theta f(\boldsymbol{\theta}). \quad (5.5)$$

We can view (5.5) as the expectation of  $f$  over  $\Omega$ , with respect to the uniform distribution  $U(\Omega)$ . In Monte-Carlo integration, we sample points uniformly on  $\Omega$ , and then approximate  $I$  via

$$I_N = \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}_i). \quad (5.6)$$

Here  $N$  is the number of times we have sampled from  $U(\Omega)$ , and  $\mathbf{x}_i$  are the sample points. Monte Carlo integration works if we can efficiently evaluate  $f$ . From the law of large numbers,

$$\lim_{N \rightarrow \infty} I_N = I. \quad (5.7)$$

The standard error of the mean goes as  $N^{-1/2}$ , which gives us our estimate for the error of this approximation; that is we can write (e.g. [PTVF92])

$$\int_{\Omega} d^k \theta f(\boldsymbol{\theta}) \approx V(\Omega) \left( \mathbb{E}[f] \pm \sqrt{\frac{\mathbb{V}[f]}{N}} \right), \quad (5.8)$$

where  $V(\Omega)$  is the volume of  $\Omega$ . We see that the convergence of Monte Carlo integration scales as  $1/\sqrt{N}$ , regardless of the dimensionality of the integral. This is the key property of stochastic integration methods, and what makes them widespread use in computing high dimensional integrals. In one dimension, almost any other quadrature method outperforms Monte Carlo integration (for example, the error to the trapezoid rule scales as  $1/N^2$ ), but for higher dimensional integrals the convergence of most methods rapidly deteriorates.

We can think of Monte Carlo integration as an example of a stochastic approximation to the Riemann integral of  $f$ . Recall that the Riemann integral is the limit of the sum over  $f(\boldsymbol{\theta}_i)$  multiplied by the volume of a small (possibly multidimensional) rectangle centered on  $f(\boldsymbol{\theta}_i)$ , which we call  $V(\boldsymbol{\theta}_i)$ .

$$\int d^k \theta f(\boldsymbol{\theta}) = \lim_{N \rightarrow \infty} \sum_{i=1}^N V(\boldsymbol{\theta}_i) f(\boldsymbol{\theta}_i). \quad (5.9)$$

In effect, in Monte Carlo integration we approximate  $V(\boldsymbol{\theta}_i)$  with  $n/N$ , where  $n$  is the number of draws we made in that volume.

Monte Carlo integration works best if most of the integral  $I$  is not concentrated in a few small volume regions; that is if  $f(\boldsymbol{\theta})$  is not too “peaked”. Most likelihoods are strongly

peaked though—for example from the Bernstein–von Mises theorem (see Appendix A) we expect the likelihood to approximately behave as a multivariate normal function around the maximum likelihood estimator as the amount of data we collect goes to infinity. Moreover the covariance matrix elements scale as  $1/N_d$ , where  $N_d$  is the number of data points, so the distribution becomes increasingly localized near the maximum likelihood estimator, and more generally near other local maxima of the likelihood. That is, selecting  $\boldsymbol{\theta}_i$  from the uniform distribution could mean that we are mostly sampling from places where  $f(\boldsymbol{\theta}_i)$  is much smaller than near the peaks. In that case, we would be missing most what contributes to the integral (5.5), which slows down the rate of convergence (the prefactor in front of the asymptotic scaling of  $1/\sqrt{N}$ ).

This motivates the introduction of integrations methods that preferentially sample from regions near the local maxima of the integrand in (5.5). We next discuss two such adaptive methods: Markov chain Monte Carlo (MCMC) methods, which can also be thought of as an adaptive approximation to the Riemann integral, and nested sampling methods, which can be thought of as an adaptive approximation to the Lebesgue integral<sup>1</sup>

## 5.2 Markov chain Monte Carlo (MCMC)

The idea behind MCMC integration is to generate the random samples for the Monte Carlo integration of  $\Omega$  dynamically, through a Markov Chain. To do this, we rewrite the integral (5.5) as follows

$$I = \int d^k \theta p(\boldsymbol{\theta}) g(\boldsymbol{\theta}), \quad (5.10)$$

where

$$\int d^k \theta p(\boldsymbol{\theta}) = 1. \quad (5.11)$$

That is, we interpret  $p(\boldsymbol{\theta})$  as a probability distribution. We can view (5.10) as a Riemann–Stieltjes integral,

$$I = \int dF(\boldsymbol{\theta}) g(\boldsymbol{\theta}), \quad (5.12)$$

with the measure  $dF(\boldsymbol{\theta}) \equiv d^k \theta p(\boldsymbol{\theta})$ . We defined

$$F(\lambda) \equiv \int_0^{g(\boldsymbol{\theta}) < \lambda} d^k \theta p(\boldsymbol{\theta}). \quad (5.13)$$

---

<sup>1</sup>For smooth functions—which is what the posterior distribution function  $P(\boldsymbol{\theta}|\mathbf{x})$  is, provided the prior and our model  $P(\mathbf{x}|\boldsymbol{\theta})$  are smooth—there is no substantive, practical difference between the Riemann and Lebesgue integral. Nevertheless we will see that there are different strengths and weaknesses to MCMC and nested sampling, unrelated to the kinds of integrals they are approximating.

Properly speaking, MCMC is a method for drawing samples from  $p(\boldsymbol{\theta})$ , for use in calculating integrals of the form (5.10). It turns out that a histogram of our sampling of  $p(\boldsymbol{\theta})$  will begin to resemble  $p(\boldsymbol{\theta})$  as the number of draws goes to infinity. For this reason, MCMC methods are often seen as ways to determine the “shape” or functional properties of  $p(\boldsymbol{\theta})$ . Here we take the perspective of numerical integration theory, so we think of  $p(\boldsymbol{\theta}) = f(\boldsymbol{\theta})/g(\boldsymbol{\theta})$  as a weighting factor for our integration of (5.5).

First we define a few concepts about Markov Chains. We will borrow some of the concepts from the discussion on stochastic processes in Chptr. 4. A **Markov chain** is a sequence of random vectors (a stochastic process)  $\{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots\}$ , where the probability distribution of  $\boldsymbol{\theta}_{n+1}$  is dependent solely on  $\boldsymbol{\theta}_n$ . That is  $P(\boldsymbol{\theta}_{n+1} | \{\boldsymbol{\theta}_i\}) = P(\boldsymbol{\theta}_{n+1} | \boldsymbol{\theta}_n)$ . We call  $P(\boldsymbol{\theta}_{n+1} | \boldsymbol{\theta}_n)$  the **transition probability**. The probability distribution for  $\boldsymbol{\theta}_1$ ,  $P(\boldsymbol{\theta}_1)$ , is called the **initial distribution**. We will only consider Markov chains with **stationary transition probabilities**, where the transition probabilities do not depend on  $n$ .

Operationally, MCMC integration of (5.5) goes as follows.

1. We pick an initial point  $\boldsymbol{\theta}_1$ , and then generate new samples  $\boldsymbol{\theta}_n$  based on a suitably chosen transition probability.
2. For the first few iterations of the Markov Chain, the points  $\boldsymbol{\theta}_n$  will be highly correlated with our initial start point, but if one runs the Markov Chain for enough iterations, the points  $\{\boldsymbol{\theta}_n\}$  will eventually converge to the target distribution  $p(\boldsymbol{\theta})$ .
3. Integration then proceeds as in Monte Carlo integration

$$I_N = \frac{1}{N} \sum_{i=1}^N g(\boldsymbol{\theta}_i), \quad (5.14)$$

with similar convergence properties to Monte Carlo integration (the error will asymptotically go down as  $1/\sqrt{N}$ ). The hope though is that the prefactor to the leading asymptotic decay will be much smaller than it would be for regular Monte Carlo integration.

We see that we can view (5.14) as an approximation to the Riemann–Stieltjes integral (5.12).

There are whole volumes on MCMC (e.g. [BGJM11]); for a nice shorter review see [HFM18]. Here we only outline what a “suitable” Markov Chain transition probability must satisfy, the Metropolis-Hastings algorithm, and some limitations of most MCMC methods.

An MCMC chain must eventually limit to a stationary distribution that is equal to  $p(\boldsymbol{\theta})$ . A sufficient (but not necessary) conditions for a Markov chain to have a stationary distribution  $Q(\boldsymbol{\theta})$  is that the transition probabilities must satisfy the **detailed balance** condition

$$P(\boldsymbol{\theta} | \boldsymbol{\psi}) Q(\boldsymbol{\psi}) = P(\boldsymbol{\psi} | \boldsymbol{\theta}) Q(\boldsymbol{\theta}), \quad (5.15)$$

for any  $\boldsymbol{\theta}, \boldsymbol{\psi}$ . To see why detailed balance implies stationarity, we compute the probability of a transition to new step  $\boldsymbol{\theta}_n$ . The probability distribution for a new step  $\boldsymbol{\theta}$  is equal to the integral (or sum, if there a discrete number of points) over all possible earlier points  $\boldsymbol{\psi}$ . We assume those are distributed according to the probability distribution  $Q(\boldsymbol{\psi})$ . We then show that the distribution for  $\boldsymbol{\theta}$ ,  $P(\boldsymbol{\theta})$ , is equal to  $Q(\boldsymbol{\theta})$ , which implies that the chain is stationary. We have

$$\begin{aligned}
P(\boldsymbol{\theta}) &= \int d^k \boldsymbol{\psi} P(\boldsymbol{\theta}|\boldsymbol{\psi}) Q(\boldsymbol{\psi}) \\
&= \int d^k \boldsymbol{\psi} P(\boldsymbol{\psi}|\boldsymbol{\theta}) Q(\boldsymbol{\theta}) \\
&= \frac{Q(\boldsymbol{\theta})}{P(\boldsymbol{\theta})} \int d^k \boldsymbol{\psi} P(\boldsymbol{\psi}, \boldsymbol{\theta}) \\
&= Q(\boldsymbol{\theta}).
\end{aligned} \tag{5.16}$$

This proves existence of a stationary chain, but it does not prove uniqueness. Proving uniqueness of the stationary distribution is beyond the scope of these notes. Most practitioners simply ignore the question of uniqueness.

Finally, we discuss an example of a Markov Chain that satisfies the detailed balance condition for the target function  $p(\boldsymbol{\theta})$  (lower case  $p$ ; see (5.10)): the **Metropolis-Hastings algorithm**. Consider a point  $\boldsymbol{\theta}$ . We draw  $\boldsymbol{\psi}$  from the **proposal probability**  $Q(\boldsymbol{\psi}|\boldsymbol{\theta})$  (we are free to specify  $Q$ ). We then draw a random variable  $x$  from the uniform distribution  $U(0, 1)$ . We next compute the **acceptance probability**

$$r = \min \left( 1, \frac{p(\boldsymbol{\theta})}{p(\boldsymbol{\psi})} \frac{Q(\boldsymbol{\theta}|\boldsymbol{\psi})}{Q(\boldsymbol{\psi}|\boldsymbol{\theta})} \right). \tag{5.17}$$

If  $x > r$ , we jump to the point  $\boldsymbol{\psi}$ , otherwise, we stay at the point  $\boldsymbol{\theta}$ . To prove that transition probability in the Metropolis-Hastings algorithm satisfies the detailed balance condition, we rewrite the transition probability amplitude as being equal to the proposal probability times the acceptance probability

$$P(\boldsymbol{\psi}|\boldsymbol{\theta}) = r \times Q(\boldsymbol{\psi}|\boldsymbol{\theta}). \tag{5.18}$$

We then have

$$\begin{aligned}
P(\boldsymbol{\theta}|\boldsymbol{\psi}) p(\boldsymbol{\psi}) &= \min(p(\boldsymbol{\psi}) Q(\boldsymbol{\psi}|\boldsymbol{\theta}), p(\boldsymbol{\theta}) Q(\boldsymbol{\theta}|\boldsymbol{\psi})) \\
&= \min(p(\boldsymbol{\theta}) Q(\boldsymbol{\theta}|\boldsymbol{\psi}), p(\boldsymbol{\psi}) Q(\boldsymbol{\psi}|\boldsymbol{\theta})) \\
&= P(\boldsymbol{\psi}|\boldsymbol{\theta}) p(\boldsymbol{\theta}).
\end{aligned} \tag{5.19}$$

## 5.3 Nested sampling

As with MCMC integration, we consider integrals of the form

$$I = \int d^k \boldsymbol{\theta} p(\boldsymbol{\theta}) L(\boldsymbol{\theta}). \tag{5.20}$$

We can only integrate positive definite functions with the nested sampling algorithm, which is why we use the slightly different notation of  $L$  instead of  $g$  here: we restrict to functions such that  $L \geq 0$ . This notation is motivated from the following: the main application of the nested sampling integration is to compute the evidence  $\mathcal{Z}$ , which is an integral of the prior probability distribution times the likelihood

$$\mathcal{Z} = \int d^k \theta \pi(\boldsymbol{\theta}) \mathcal{L}(\boldsymbol{\theta}). \quad (5.21)$$

Before we describe the algorithm, we first need to rewrite (5.20) as an integral over the level sets of  $L(\boldsymbol{\theta})$ . To do this, we write (5.20) as Riemann-Stieltjes integral (5.12), and then integrate by parts. We define the function

$$X(\lambda) \equiv \int_{L(\boldsymbol{\theta}) > \lambda} d^k \theta p(\boldsymbol{\theta}). \quad (5.22)$$

As  $\lambda$  increases,  $X$  decreases from 1 to 0. With this, we have

$$\begin{aligned} I &= \int d^k \theta p(\boldsymbol{\theta}) L(\boldsymbol{\theta}) = - \int dX L \\ &= - X(L) L \Big|_{L=0}^{L=L_{max}} + \int_0^{L_{max}} dL X(L) \\ &= \int_0^{L_{max}} dL X(L). \end{aligned} \quad (5.23)$$

We assumed that  $L_{min} = 0$  (which holds for the likelihood function), and used that  $X = 0$  at  $L = L_{max}$ . We assume that we can invert  $X(L)$  to the function  $L(X)$ . We can then rewrite (5.23) by integrating by parts, to obtain

$$I = - \int_1^0 dX L(X) = \int_0^1 dX L(X). \quad (5.24)$$

Unlike (5.20), (5.24) is a one-dimensional integral. In (5.24) we should think of  $L$  as the parameter in  $X(L)$ , not as  $L(\boldsymbol{\theta})$ . Nested sampling provides a noisy approximation to (5.24), through a partitioning of the  $X$  interval, and hence provides a noisy approximation to the Lebesgue integral of (5.20).

To understand why (5.24) is the Lebesgue integral of (5.20), recall that the Lebesgue integral is the limit of the sum over  $g_i \equiv g(\boldsymbol{\theta}_i)$  multiplied by the Lebesgue measure of the set  $E_i$  of points  $\boldsymbol{\theta}_j$  for which  $g(\boldsymbol{\theta}_j) \approx g(\boldsymbol{\theta}_i)$

$$I_N = \sum_{i=1}^N g_i \mu(E_i). \quad (5.25)$$

We can think of  $g_i \mu(E_i)$  as the discretization of  $dX \lambda(X)$ .

The nested sampling algorithm goes as follow

1. We draw  $n$  points  $\boldsymbol{\theta}_i$  from  $p(\boldsymbol{\theta})$ , treating it as a probability distribution. Set  $X_0 = 1$ .

2. Repeat for  $N$  times, so you have the sequence  $X_1, \dots, X_n$  and  $L_{min,1}, \dots, L_{min,N}$ . For the  $j^{th}$  iteration
  - (a) Record the lowest value of  $= L_{min,j} = L(\theta_j)$ . Set  $X_j = e^{-j/n}$ , or alternatively set  $X_j = t_j X_{j-1}$ , where  $t_j$  is drawn from the beta distribution  $\text{Beta}(1, n)$ .
  - (b) Remove the value of  $\theta_j$  that minimizes  $L(\theta)$ , and then sample again from  $p(\theta)$ , until you get a point  $\theta_k$  such that  $L(\theta_k) > L_{min,j}$ .
3. The integral can then be obtained by summing the  $L$  values via, e.g. the trapezoid rule

$$I_N = \sum_{i=1}^{N-1} w_i L_{min,i}, \quad (5.26)$$

where  $w_i = (X_{i+1} - X_{i-1})/2$ .

One of the tricky things to understand about the nested sampling method is the value of the measure of the likelihood  $L_{min,i}$ ,  $X_i$ . Consider a sample from  $p(\theta)$ :  $\{\theta_j\}$ , subject to  $L(\theta_j) > L_{min,j-1}$ . We assume the values of the volumes  $X(\theta_j)$  are uniformly distributed in the interval  $[0, X_{j-1}]$ . Then  $X_j = t_j X_{j-1}$ , where  $t_j$  is the largest of  $n$  uniformly distributed numbers in the interval  $(0, 1)$ . The number  $t_j$  is called the **shrinkage factor**. Notice that we have

$$X_j = \prod_{i=1}^j t_i. \quad (5.27)$$

The cumulative probability distribution function for the maximum of  $n$  randomly distributed numbers in that interval is

$$\begin{aligned} C.D.F. (t_{max}) &= P(\max\{t_1, \dots, t_n\} < t_{max}) \\ &= (P(t < t_{max}))^n \\ &= t_{max}^n. \end{aligned} \quad (5.28)$$

The probability density function for the maximum is then the beta distribution  $\text{Beta}(1, n)$ , that is

$$P(t_{max}) = n t_{max}^{n-1}. \quad (5.29)$$

To estimate  $X_j$  then, we could take a draw from the Beta distribution,  $t$ , and multiply that by  $X_{j-1}$ . To get a (presumably) less noisy answer, we could set  $X_j$  to be its averaged expected value. We take the expectation of the log of  $X_j$ , to simplify the calculation of the expectation

$$\mathbb{E}[\log X_j] = \sum_{i=1}^j \mathbb{E}[\log t_i]. \quad (5.30)$$

As the  $t_i$  are independent, we can estimate the error of this approximation by computing the variance

$$\mathbb{V}[\log X_j] = \sum_{i=1}^j \mathbb{V}[\log t_i]. \quad (5.31)$$

The expectation value of the logarithm of  $t_{max}$  is

$$\begin{aligned} \mathbb{E}[\log t_{max}] &= \int_0^1 dt \, n t^{n-1} \log t \\ &= -\frac{1}{n}. \end{aligned} \quad (5.32)$$

The variance of the log of  $t_{max}$

$$\begin{aligned} \mathbb{V}[\log t_{max}] &= \mathbb{E}[(\log t_{max})^2] - (\mathbb{E}[\log t_{max}])^2 \\ &= \int_0^1 dt \, n t^{n-1} (\log t)^2 - \frac{1}{n^2} \\ &= \frac{1}{n^2}. \end{aligned} \quad (5.33)$$

Combining everything, we see that the shrinkage factor is approximately

$$\begin{aligned} \log X_j &\approx j \mathbb{E}[\log t_{max}] + \sqrt{j \mathbb{V}[\log t_{max}]} \\ &= -\frac{j}{n} \left( 1 \pm \frac{1}{\sqrt{j}} \right). \end{aligned} \quad (5.34)$$

This gives us

$$X_j \approx e^{-j/n}. \quad (5.35)$$

Our approximation to the integral gets better as we add more points  $n$ , and as we take more steps  $N$ . We incur the biggest relative errors in the integration for the first few small steps  $j$ , but so long as  $L$  is highly peaked, and we take very small steps, those terms contribute very little to the total integral.

# Chapter 6

## Numerical optimization

While most often we want to compute integrals of the posterior probability distribution (for example, to compute the mean or covariance matrix of the posterior), sometimes it is informative to simply compute its maximum, or even just the maximum of the likelihood. Again we consider the posterior

$$P(\boldsymbol{\theta}|\mathbf{x}) = \frac{\mathcal{L}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta})}{\mathcal{Z}(\mathbf{x})}, \quad (6.1)$$

The value of  $\boldsymbol{\theta}$  that maximizes  $\mathcal{L}(\boldsymbol{\theta})$  is the **maximum likelihood estimator (MLE)**, and the value of  $\boldsymbol{\theta}$  that maximizes  $\mathcal{L}(\boldsymbol{\theta}) \pi(\boldsymbol{\theta})$  is the **maximum a posteriori probability estimator (MAP)**. Note that the MLE and MAP do not give us any knowledge of the variance of those parameters—that requires knowledge of the full posterior probability distribution. This in turn requires integration of the likelihood.

From a numerical point of view, it is convenient to consider numerical minimizers, and to find the MLE by finding the minimum of the negative log likelihood, which we call  $\ell(\boldsymbol{\theta})$

$$\ell(\boldsymbol{\theta}) \equiv -\log \mathcal{L}(\boldsymbol{\theta}). \quad (6.2)$$

It is convenient to consider the log likelihood, as the likelihood itself can vary drastically in value between its maxima and minima, which can be hard for a computer to resolve with finite precision arithmetic. The likelihood, prior, evidence, and posterior are positive definite quantities as well, so there is no change of taking a logarithm of these quantities.

Here we review a few minimization methods. There is no best method that will work for all likelihoods, so we only review the basics of a few basic methods that underlie more complex optimization procedures. For concreteness we will focus on minimizing the negative log likelihood  $\ell(\boldsymbol{\theta})$ .



## 6.1 Convex functions

We first consider the problem of optimizing convex functions. While the posterior is almost never convex, it is still useful to review this case first as the local max/min of a strictly convex function is the global maximum/minimum (this is almost never the case for non-convex functions). Because of this, some methods (namely, Newton’s method and its extensions) used to find the minimum of functions try to convert the problem into one for finding the minimum of a convex function.

A **convex function**  $f(\boldsymbol{\theta}) : X \rightarrow \mathbb{R}$  satisfies

$$f(t\boldsymbol{\theta}_1 + (1-t)\boldsymbol{\theta}_2) \leq tf(\boldsymbol{\theta}_1) + (1-t)f(\boldsymbol{\theta}_2), \quad (6.3)$$

for  $t \in [0, 1]$  and for all  $\boldsymbol{\theta}_{1,2} \in X$  (for example,  $X = \mathbb{R}^n$ ). A **strictly convex** function satisfies (6.4) except the  $\leq$  is replaced by  $<$ , and  $t \in (0, 1)$  instead.

The local minimum of a convex function is the global minimum of the function. This is easy to show. Say  $\boldsymbol{\theta}_*$  is a local minimum, and assume that we have found  $\boldsymbol{\theta}$  such that  $f(\boldsymbol{\theta}) < f(\boldsymbol{\theta}_*)$ . Then we would have

$$\begin{aligned} f(t\boldsymbol{\theta}_* + (1-t)\boldsymbol{\theta}) &\leq tf(\boldsymbol{\theta}_*) + (1-t)f(\boldsymbol{\theta}) \\ &< tf(\boldsymbol{\theta}_*) + (1-t)f(\boldsymbol{\theta}_*) = f(\boldsymbol{\theta}_*). \end{aligned} \quad (6.4)$$

Setting  $t = 1$ , we encounter a contradiction, which concludes the argument. We see for a strictly convex function, a local minimum is a global minimum, and the global minimum is unique.

## 6.2 Gradient descent

First we consider a linear method for finding local minima—gradient descent. To understand this method, we Taylor series expand the negative log likelihood about a fiducial point  $\boldsymbol{\theta}_0$

$$\ell = \ell_0 + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{g}_0 + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{H}_0 (\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \cdots, \quad (6.5)$$

where

$$\mathbf{g}_{0,i} \equiv \nabla_i \ell(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}, \quad (6.6)$$

$$\mathbf{H}_{0,ij} \equiv \nabla_i \nabla_j \ell(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}. \quad (6.7)$$

As a local minimum  $\boldsymbol{\theta}_*$ , the gradient of the function is zero, and the Hessian is positive definite. Near a local minimum then, we expect the gradient to be pointing “away” from the local minimum. Thus if move in the opposite direction to the gradient, we move in the

direction of the local minimum. In gradient descent then, we pick a fiducial value of  $\boldsymbol{\theta}_0$ , and then iterate the following

$$\boldsymbol{\theta}_{n+1} = \boldsymbol{\theta}_n - \gamma_n \mathbf{g}_n. \quad (6.8)$$

Where  $0 < \gamma_n$  is a scalar that one can introduce to make the change between steps be less large. We stop iterating when  $|\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n| < \epsilon$ , where  $|\cdots|$  is a norm of our choosing and  $0 < \epsilon$  is a pre-set tolerance.

While gradient descent is an easy algorithm to implement, it suffers from a few problems. First, the method (if it converges) only find a local minimum, or possibly only a saddle point. Also, it can be tricky to find a good value of  $\gamma_n$ . If  $\gamma_n$  is too small, the method converges very slowly. If  $\gamma_n$  is too large, the method may never converge.

## 6.3 Newton's method

We next consider a quadratic method for finding local minima—Newton's method (this is Newton's method for optimizing a function, not for finding the root to a function). To understand this method, we Taylor series expand the negative log likelihood about a fiducial point  $\boldsymbol{\theta}_0$

$$\ell = \ell_0 + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{g}_0 + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{H}_0 (\boldsymbol{\theta} - \boldsymbol{\theta}_0) + \cdots, \quad (6.9)$$

where

$$\mathbf{g}_{0,i} \equiv \nabla_i \ell(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}, \quad (6.10)$$

$$\mathbf{H}_{0,ij} \equiv \nabla_i \nabla_j \ell(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}. \quad (6.11)$$

We choose  $\boldsymbol{\theta}$  to minimize the quadratic Taylor series expansion. *Assuming* that  $\mathbf{H}_0$  is positive definite, minimizing the quadratic Taylor series expansion is an exercise in minimizing a convex function. The minimum is then located at the zero of the gradient of the second order Taylor series. We find that

$$\mathbf{g}_0 + \mathbf{H}_0 (\boldsymbol{\theta} - \boldsymbol{\theta}_0) = 0 \implies \boldsymbol{\theta} = \boldsymbol{\theta}_0 - \mathbf{H}_0^{-1} \mathbf{g}_0. \quad (6.12)$$

This motivates **Newton's method**. Starting with a fiducial point  $\boldsymbol{\theta}_0$ , we iterate in  $\boldsymbol{\theta}_n$ , where at each iteration we set

$$\boldsymbol{\theta}_{n+1} = \boldsymbol{\theta}_n - \gamma_n \mathbf{H}_n^{-1} \mathbf{g}_n. \quad (6.13)$$

Where  $0 < \gamma_n \leq 1$  is a scalar that one can introduce to make the change between steps be less large. We stop iterating when  $|\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n| < \epsilon$ , where  $|\cdots|$  is a norm of our choosing and  $0 < \epsilon$  is a pre-set tolerance.

As with the gradient descent method, Newton's method may only find a local minimum or saddle point. There are also numerous technical problems with inversion of the Hessian. First, the Hessian matrix may be very large if there are many parameters, so it could be hard to invert (it may be ill conditioned). Additionally the Hessian could be singular, or nearly singular. We note that the Hessian may not be positive definite either at a given point (it often won't be), which in principle isn't fatal to the method, but depending on the size of the eigenvalues to the Hessian, large negative eigenvalues could dramatically change the value of  $\theta_{n+1}$  versus  $\theta_n$ .

# Appendix A

## Probability theory

### A.1 Conditional probability and Bayes theorem

The condition probability is

$$P(\boldsymbol{\theta}|\boldsymbol{\psi}) = \frac{P(\boldsymbol{\theta}, \boldsymbol{\psi})}{P(\boldsymbol{\psi})}. \quad (\text{A.1})$$

Using this, we have Bayes theorem

$$P(\boldsymbol{\theta}|\boldsymbol{\psi}) P(\boldsymbol{\psi}) = P(\boldsymbol{\psi}|\boldsymbol{\theta}) P(\boldsymbol{\theta}). \quad (\text{A.2})$$

Bayes theorem is more often written as

$$P(\boldsymbol{\theta}|\boldsymbol{x}) = \frac{P(\boldsymbol{x}|\boldsymbol{\theta}) P(\boldsymbol{\theta})}{P(\boldsymbol{x})}. \quad (\text{A.3})$$

As we discuss more in Chptr. 1, we can view  $P(\boldsymbol{\theta}|\boldsymbol{x})$  as the distribution of model parameters given a model  $P(\boldsymbol{x}|\boldsymbol{\theta})$ , and some prior knowledge of the model parameters  $P(\boldsymbol{\theta})$ .

### A.2 Change of variables

Consider a distribution  $P(\boldsymbol{\theta})$ . What is the probability distribution to  $P(\boldsymbol{\psi}(\boldsymbol{\theta}))$ , where  $\boldsymbol{\psi}$  is some function of  $\boldsymbol{\theta}$ ? The following remains unchanged under a change of variables

$$\begin{aligned} \int_V d^k \theta P(\boldsymbol{\theta}) &= \int_V d^k \psi P(\boldsymbol{\psi}) \\ &= \int_V d^k \theta |\det(J_{ij})| P(\boldsymbol{\psi}). \end{aligned} \quad (\text{A.4})$$

We viewed  $V$  as a geometric volume (that is, it is independent of the coordinate choice we use). Here

$$J_{ij} \equiv \frac{\partial \psi^i}{\partial \theta_j}, \quad (\text{A.5})$$

is the Jacobian matrix. Equating terms within the integral, we find that

$$P(\boldsymbol{\psi}) = \frac{1}{|\det(J_{ij})|} P(\boldsymbol{\theta}). \quad (\text{A.6})$$

As an example application of this formula, we consider the posterior probability distribution for  $\boldsymbol{\psi}(\boldsymbol{\theta})$ . From Bayes theorem (A.3), we have

$$\begin{aligned} P(\boldsymbol{\psi}|\mathbf{x}) &= \frac{P(\mathbf{x}|\boldsymbol{\psi}(\boldsymbol{\theta})) P(\boldsymbol{\psi})}{P(\mathbf{x})} \\ &= \frac{1}{|\det(J_{ij})|} \frac{P(\mathbf{x}|\boldsymbol{\psi}(\boldsymbol{\theta})) P(\boldsymbol{\theta})}{P(\mathbf{x})}. \end{aligned} \quad (\text{A.7})$$

That is, to find the probability distribution for some function of the distribution parameters, we only need to find the probability distribution for the prior under that change in coordinates.

## A.3 Expectation and covariance

We review a few basic definitions from probability theory, as they come up later in the notes. We define the **expectation** of a random variable  $\boldsymbol{\Theta}(\mathbf{x})$  to be

$$\mathbb{E}[\boldsymbol{\Theta}] \equiv \int d\mathbf{x} P(\mathbf{x}|\boldsymbol{\theta}) \boldsymbol{\Theta}(\mathbf{x}). \quad (\text{A.8})$$

Sometimes we denote the expectation with  $\mu_{\boldsymbol{\Theta}}$ . The **variance** is

$$\mathbb{V}_{ij}[\boldsymbol{\Theta}] \equiv \mathbb{E}[(\Theta_i - \mu_{\Theta_i})(\Theta_j - \mu_{\Theta_j})]. \quad (\text{A.9})$$

The **covariance** for two random variables  $\boldsymbol{\Theta}, \boldsymbol{\Psi}$  is

$$\mathbb{C}_{ij}[\boldsymbol{\Theta}, \boldsymbol{\Psi}] \equiv \mathbb{E}[(\Theta_i - \mu_{\Theta_i})(\Psi_j - \mu_{\Psi_j})]. \quad (\text{A.10})$$

Note that we can think of the expectation of two scalar random variables  $\Theta, \Psi$  as an inner product

$$\mathbb{E}[\Theta\Psi] = \langle \Theta, \Psi \rangle. \quad (\text{A.11})$$

It is easy to see from (A.8) that (A.11) satisfies the properties of an inner product:  $\langle \Theta, \Theta \rangle \geq 0$ ,  $\langle \Theta, \Psi \rangle = \langle \Psi, \Theta \rangle$ , and linearity.

Consider a random variable  $\Theta(\mathbf{x})$ . Say we want this variable to represent another random variable, say the parameters of the posterior  $\theta$ . We then call  $\Theta$  an **estimator** for  $\theta$ . The **bias** of  $\Theta$  then is

$$\mathbf{b}(\Theta) \equiv \mathbb{E}[\Theta] - \theta. \quad (\text{A.12})$$

If  $\mathbf{b} = 0$ , then  $\Theta$  is an **unbiased estimator** for  $\theta$ .

## A.4 Characteristic/moment generating function

To prove the central limit theorem, first we introduce the Fourier transform (or **characteristic function**) of a probability distribution. Consider a random vector  $\mathbf{X}$  with probability distribution  $f_{\mathbf{X}}(\mathbf{x})$ , that is  $\mathbf{X} \sim f_{\mathbf{X}}$ . We denote the moment generating function with  $\psi_{\mathbf{X}}$ , which is

$$\psi_{\mathbf{X}}(\mathbf{t}) \equiv \int_{-\infty}^{\infty} d\mathbf{x} e^{i\mathbf{t}^T \mathbf{x}} f_{\mathbf{X}}(\mathbf{x}) = \mathbb{E} \left[ e^{i\mathbf{t}^T \mathbf{X}} \right]. \quad (\text{A.13})$$

If we set  $\mathbf{t} = -i\tilde{\mathbf{t}}$ , then  $\psi_{\mathbf{X}}$  is called the **moment generating function**. For most probability distributions, there is no meaningful difference between using  $\mathbf{t}$  or  $-i\tilde{\mathbf{t}}$  (there could potentially only be a difference if the distribution had complex poles or branch cuts). Notice that

$$\mathbb{E}[X_{i_1} \cdots X_{i_k}] = \frac{1}{i^k} \nabla_{t_{i_1}} \cdots \nabla_{t_{i_k}} \psi_{\mathbf{X}}(\mathbf{t}). \quad (\text{A.14})$$

That is, we can obtain the moments of the probability distribution from the characteristic function (although we need to divide by  $1/i^k$ ).

Perhaps most importantly, notice that since the characteristic function for a probability distribution is the Fourier transform of the probability density, we can uniquely map a probability density to its characteristic function and back. That is, given a characteristic function, we can find the unique probability density that it corresponds to.

Consider a linear affine transformation of the random variable  $\mathbf{X}$ , which we call  $\mathbf{Y} = a\mathbf{X} + \mathbf{b}$ . We also call  $\mathbf{y} = a\mathbf{x} + \mathbf{b}$ . The probability distribution with the volume element remains unchanged  $dy f_{\mathbf{Y}}(\mathbf{x}) = dx f_{\mathbf{X}}(\mathbf{x})$ . We conclude that

$$\begin{aligned} \psi_{\mathbf{Y}}(\mathbf{t}) &= \int_{-\infty}^{\infty} d\mathbf{y} e^{i\mathbf{t}^T \mathbf{y}} f_{\mathbf{Y}}(\mathbf{y}) \\ &= e^{i\mathbf{t}^T \mathbf{b}} \int_{-\infty}^{\infty} d\mathbf{x} e^{i\mathbf{t}^T a\mathbf{x}} f_{\mathbf{X}}(\mathbf{x}) \\ &= e^{i\mathbf{t}^T \mathbf{b}} \psi_{\mathbf{X}}(a\mathbf{t}). \end{aligned} \quad (\text{A.15})$$

The characteristic function of a sum of independent variables is the product of the characteristic function for each variable. Define  $\mathbf{Y} = \sum_{i=1}^n \mathbf{X}_i$ , we then have

$$\begin{aligned}\psi_{\mathbf{Y}}(\mathbf{t}) &= \int_{-\infty}^{\infty} dx_1 \cdots \int_{-\infty}^{\infty} dx_n e^{i\mathbf{t}^T \sum_i \mathbf{x}_i} f_{\mathbf{Y}}(\mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= \prod_{i=1}^n \int_{-\infty}^{\infty} dx_i e^{i\mathbf{t}^T \mathbf{x}_i} f_{\mathbf{X}_i}(\mathbf{x}_i) \\ &= \prod_{i=1}^n \psi_{\mathbf{X}_i}(\mathbf{t}).\end{aligned}\tag{A.16}$$

The second line follows from  $f_{\mathbf{Y}}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_i f_{\mathbf{X}_i}(\mathbf{x}_i)$ , as all the variables are independent.

Consider a multivariate normal variable  $\mathbf{X} \sim \mathcal{N}_k(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . We see that

$$\begin{aligned}\psi_{\mathbf{X}}(\mathbf{t}) &= \int_{-\infty}^{\infty} dx \frac{1}{\sqrt{2\pi} |\det \boldsymbol{\Sigma}|} \exp \left[ i\mathbf{t}^T \mathbf{x} - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \\ &= \exp \left[ i\mathbf{t}^T \boldsymbol{\mu} + i^2 \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t} \right] \\ &\quad \times \int_{-\infty}^{\infty} dx \frac{1}{\sqrt{2\pi} |\det \boldsymbol{\Sigma}|} \exp \left[ -\frac{1}{2} (\mathbf{x} - i\boldsymbol{\Sigma} \mathbf{t} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - i\boldsymbol{\Sigma} \mathbf{t} - \boldsymbol{\mu}) \right] \\ &= \exp \left[ i\mathbf{t}^T \boldsymbol{\mu} + i^2 \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t} \right].\end{aligned}\tag{A.17}$$

## A.5 Central limit theorem

Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be  $n$  independent and identically distributed random vectors (of dimension  $k$  each), and let each variable have mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ . The central limit theorem states that the probability distribution of the average of these variables,

$$\bar{\mathbf{X}}_n \equiv \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i,\tag{A.18}$$

limits to a multivariate normal distribution with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}/n$  as  $n \rightarrow \infty$ . Note that we made no assumption about the probability distribution for the  $\mathbf{X}_i$ , except that the probability distribution has a finite mean and variance. We can write the central limit theorem as

$$\lim_{n \rightarrow \infty} \sqrt{n} \bar{\mathbf{X}}_n \sim \mathcal{N}_k(\boldsymbol{\mu}, \boldsymbol{\Sigma}),\tag{A.19}$$

where  $\mathcal{N}_k$  is the multivariate normal distribution. To prove this, we make use of the characteristic function for  $\mathbf{X}_n$ , which is

$$\psi_{\mathbf{X}_n}(\mathbf{t}) = \prod_{i=1}^n \psi_{\mathbf{X}_i} \left( \frac{\mathbf{t}}{n} \right)$$

$$\begin{aligned}
&= \left( 1 + i \frac{1}{n} \mathbf{t}^T \boldsymbol{\mu} + i^2 \frac{1}{2} \frac{1}{n^2} \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t} + \mathcal{O} \left( \frac{1}{n^3} \right) \right)^n \\
&= \exp \left[ i \mathbf{t}^T \boldsymbol{\mu} + i^2 \frac{1}{2} \mathbf{t}^T \tilde{\boldsymbol{\Sigma}} \mathbf{t} \right] \left( 1 + \mathcal{O} \left( \frac{1}{n^3} \right) \right),
\end{aligned} \tag{A.20}$$

where  $\tilde{\boldsymbol{\Sigma}} = \boldsymbol{\Sigma}/n$ . We used the identity

$$\lim_{n \rightarrow \infty} \left( 1 + \frac{a}{n} \right)^n = e^a. \tag{A.21}$$

to leading order, the last line of (A.20) is the characteristic function for  $\mathcal{N}_k(\boldsymbol{\mu}, \tilde{\boldsymbol{\Sigma}})$  (see (A.17)). This concludes the proof.

## A.6 Fisher information and the Bernstein–von Mises theorem

The **Fisher information** is the negative expectation of the Hessian of the log likelihood. In terms of components, we have

$$\begin{aligned}
F_{ij}(\boldsymbol{\theta}) &\equiv -\mathbb{E}_{\boldsymbol{\theta}} [\nabla_{\theta^i} \nabla_{\theta^j} \ln P(\mathbf{x}|\boldsymbol{\theta})] \\
&= -\int d\mathbf{x} P(\mathbf{x}|\boldsymbol{\theta}) \nabla_{\theta^i} \nabla_{\theta^j} \ln P(\mathbf{x}|\boldsymbol{\theta}).
\end{aligned} \tag{A.22}$$

The Fisher information can also be written as the variance of the **score function**. The score function is

$$s_i(\mathbf{x}; \boldsymbol{\theta}) \equiv \nabla_{\theta^i} \ln P(\mathbf{x}|\boldsymbol{\theta}). \tag{A.23}$$

The expectation of the score function is zero

$$\begin{aligned}
\mathbb{E}[s_i(\mathbf{x}; \boldsymbol{\theta})] &= \int d\mathbf{x} P(\mathbf{x}|\boldsymbol{\theta}) \nabla_{\theta^i} \ln P(\mathbf{x}|\boldsymbol{\theta}) \\
&= \nabla_{\theta^i} \int d\mathbf{x} P(\mathbf{x}|\boldsymbol{\theta}) = 0.
\end{aligned} \tag{A.24}$$

We then have

$$\begin{aligned}
F_{ij}(\boldsymbol{\theta}) &= -\int d\mathbf{x} P(\mathbf{x}|\boldsymbol{\theta}) \nabla_{\theta^i} \nabla_{\theta^j} \ln P(\mathbf{x}|\boldsymbol{\theta}) \\
&= \int d\mathbf{x} \left[ \frac{1}{P(\mathbf{x}|\boldsymbol{\theta})} \nabla_{\theta^i} P(\mathbf{x}|\boldsymbol{\theta}) \nabla_{\theta^j} P(\mathbf{x}|\boldsymbol{\theta}) - \nabla_{\theta^i} \nabla_{\theta^j} P(\mathbf{x}|\boldsymbol{\theta}) \right] \\
&= \int d\mathbf{x} P(\mathbf{x}|\boldsymbol{\theta}) \nabla_{\theta^i} \ln P(\mathbf{x}|\boldsymbol{\theta}) \nabla_{\theta^j} \ln P(\mathbf{x}|\boldsymbol{\theta}) \\
&= \mathbb{E}_{\boldsymbol{\theta}} [s_i s_j]
\end{aligned}$$



$$= \mathbb{V}_{\theta, ij} [\mathbf{s}(\mathbf{x}; \boldsymbol{\theta})]. \quad (\text{A.25})$$

That is, the Fisher information is the variance of the score.

Let  $\hat{\boldsymbol{\theta}}$  be the maximum likelihood estimator for  $\boldsymbol{\theta}$ . Under appropriate regularity conditions, the likelihood  $\mathcal{L}(\boldsymbol{\theta})$  tends towards a multivariate number with mean  $\hat{\boldsymbol{\theta}}$  and covariance matrix given by the inverse Fisher information divided by the number of measurements of the data  $n$ ,  $\tilde{\mathbf{F}}^{-1} = \mathbf{F}^{-1}/n$ . In equations, we have

$$\lim_{n \rightarrow \infty} P(\boldsymbol{\theta}|\mathbf{x}) = \lim_{n \rightarrow \infty} \frac{\pi(\boldsymbol{\theta}) \prod_{i=1}^n P(\mathbf{x}_i|\boldsymbol{\theta})}{\mathcal{Z}(\mathbf{x})} = \mathcal{N}(\hat{\boldsymbol{\theta}}, \tilde{\mathbf{F}}^{-1}). \quad (\text{A.26})$$

This is known as the **Bernstein–von Mises theorem** (BvM theorem for short). We provide a rough sketch of how the proof goes. For more details see [Was10]. The log likelihood is

$$\ln \mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^N \ln P(\mathbf{x}_i|\boldsymbol{\theta}) \quad (\text{A.27})$$

We Taylor series expand the derivative of the log-likelihood to linear order about a point  $\boldsymbol{\theta}_0$

$$\nabla_{\theta^i} \ln \mathcal{L}(\boldsymbol{\theta}) = \nabla_{\theta^i} \ln \mathcal{L}(\boldsymbol{\theta}) \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} + \nabla_{\theta^i} \nabla_{\theta^j} \ln \mathcal{L}(\boldsymbol{\theta}) \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} (\theta - \theta_0)^j + \dots \quad (\text{A.28})$$

Setting  $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}$ , relabeling  $\boldsymbol{\theta}_0 \rightarrow \boldsymbol{\theta}$  (and dropping the  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ , to reduce clutter), and rearranging gives us (we used that at the maximum likelihood estimator,  $\hat{\boldsymbol{\theta}}$ , the derivative of the likelihood is zero)

$$\sqrt{n}(\hat{\theta}^i - \theta^i) = - \left( \frac{1}{n} \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} \ln \mathcal{L}(\boldsymbol{\theta}) \right)_{ij}^{-1} \left( \frac{1}{\sqrt{n}} \nabla_{\theta^j} \ln \mathcal{L}(\boldsymbol{\theta}) \right). \quad (\text{A.29})$$

As  $n \rightarrow \infty$ , we see that

$$\begin{aligned} \lim_{n \rightarrow \infty} \nabla_{\theta^j} \ln \mathcal{L}(\boldsymbol{\theta}) &= \frac{1}{\sqrt{n}} \lim_{n \rightarrow \infty} \sum_{m=0}^n \nabla_{\theta^j} \ln P(\boldsymbol{\theta}|\mathbf{x}_m) \\ &\rightarrow \sim \mathcal{N}_k(\mathbf{0}, \mathbf{F}). \end{aligned} \quad (\text{A.30})$$

This follows from the central limit theorem: the mean of the score is zero, and the variance of the score is the Fisher information. By the law of large numbers we can average

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \left( -(\nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}} \ln \mathcal{L}(\boldsymbol{\theta}))_{ij} \right) &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{m=1}^n (-\nabla_{\theta^i} \nabla_{\theta^j} \ln P(\boldsymbol{\theta}|\mathbf{x}_m)) \\ &\rightarrow F_{ij}. \end{aligned} \quad (\text{A.31})$$

Thus the variance of the limit is modified to be  $\mathbf{F}^{-1} \mathbf{F} \mathbf{F}^{-1} = \mathbf{F}^{-1}$ . We can then conclude that

$$\lim_{n \rightarrow \infty} \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \sim \mathcal{N}_k(\mathbf{0}, \mathbf{F}^{-1}). \quad (\text{A.32})$$

Or in other words

$$\lim_{n \rightarrow \infty} \sqrt{n} \hat{\boldsymbol{\theta}} \sim \mathcal{N}_k \left( \hat{\boldsymbol{\theta}}, \mathbf{F}^{-1} \right). \quad (\text{A.33})$$

We have not been careful by what we mean by “ $\rightarrow$ ” and “ $\sim$ ” here—in fact there are various notions of convergence that go into the full proof (see for example [Was10]).

We can understand the BvM theorem heuristically as follows. As we collect more data, the posterior probability becomes increasingly “peaked” near the maximum likelihood estimator. We can then Taylor series about the maximum of the log-likelihood to quadratic order. Exponentiating the log-likelihood gives us a multivariate normal with the inverse Fisher information as the covariance matrix.

## A.7 Fisher information and the Cramér-Rao bound

Consider an estimator  $\boldsymbol{\Theta}$  for model parameters  $\boldsymbol{\theta}$ . Let  $\boldsymbol{\Sigma}$  is the covariance matrix for the estimator  $\boldsymbol{\Theta}$ , let  $\mathbb{E}[\boldsymbol{\Theta}] = \boldsymbol{\psi}$ , and let  $\mathbf{F}$  be the Fisher information evaluated at  $\boldsymbol{\psi}$ . The Cramér-Rao bound states that

$$\Sigma_{ij} \geq \nabla_{\theta_m} \psi_i \nabla_{\theta_n} \psi_j F_{mn}^{-1}, \quad (\text{A.34})$$

If  $\boldsymbol{\Theta}$  is an unbiased estimator ( $\boldsymbol{\psi} = \boldsymbol{\theta}$ ), then (A.34) reduces to

$$\Sigma_{ij} \geq F_{ij}^{-1}, \quad (\text{A.35})$$

If  $\boldsymbol{\Theta}$  is a biased estimator, then  $\nabla_{\theta_i} \psi_j = \delta_{ij} + \nabla_{\theta_i} b_j$ , where  $b_j$  is the bias. The Cramér-Rao bound can be used to interpret the Fisher matrix as an estimate for the lowest error one could achieve for an unbiased estimator. For biased estimators though, we see that the Fisher information does not give a lower bound on the elements of the covariance matrix, since it is possible that the bias could be negative,  $\nabla_{\theta_i} b_j < 0$ . That is, biased estimators can have *smaller* covariance matrix elements than unbiased estimators. If the error from the bias is less than the error from the covariance (for example, if one only has a few measurements of noisy data), a biased estimator can sometimes be superior to an unbiased estimator in determining  $\boldsymbol{\theta}$ .

Here we provide the outline of a proof of (A.34). First we prove a generalization of the Cauchy-Schwartz inequality. Let  $\mathbf{y}$  and  $\mathbf{z}$  be random vectors (not necessarily of the same dimensionality). Then

$$\mathbb{V}_{ij}[\mathbf{z}] \geq \mathbb{C}_{ip}[\mathbf{z}, \mathbf{y}] \mathbb{V}_{pq}[\mathbf{y}]^{-1} \mathbb{C}_{qj}[\mathbf{y}, \mathbf{z}]. \quad (\text{A.36})$$

To prove this, we define  $\mathbf{u} \equiv \mathbf{y} - \boldsymbol{\mu}_y$  and  $\mathbf{v} \equiv \mathbf{z} - \boldsymbol{\mu}_z$ , so  $\boldsymbol{\mu}_u = 0$  and  $\boldsymbol{\mu}_v = 0$ . For any matrix  $\mathbf{A}$  we have the following matrix inequality (we insert the matrix in case  $\mathbf{v}$  and  $\mathbf{u}$  have different dimensionality)

$$(\mathbf{v} + \mathbf{A}\mathbf{u})(\mathbf{v} + \mathbf{A}\mathbf{u})^T \geq 0. \quad (\text{A.37})$$

Taking the expectation of this and expanding, we have

$$\mathbb{E} [\mathbf{v}\mathbf{v}^T] + \mathbf{A}\mathbb{E} [\mathbf{u}\mathbf{v}^T] + \mathbb{E} [\mathbf{v}\mathbf{u}^T] \mathbf{A}^T + \mathbf{A}\mathbb{E} [\mathbf{u}\mathbf{u}^T] \mathbf{A}^T \geq 0. \quad (\text{A.38})$$

Set  $\mathbf{A} = -\mathbb{E} [\mathbf{u}\mathbf{v}^T] \mathbb{E} [\mathbf{u}\mathbf{u}^T]^{-1}$ . The last two terms cancel, and we are left with

$$\mathbb{E} [\mathbf{v}\mathbf{v}^T] \geq \mathbb{E} [\mathbf{u}\mathbf{v}^T] \mathbb{E} [\mathbf{u}\mathbf{u}^T]^{-1} \mathbb{E} [\mathbf{u}\mathbf{v}^T]. \quad (\text{A.39})$$

Re-introducing  $\mathbf{y}$  and  $\mathbf{z}$ , and using the definition of the covariance (A.9) and variance (A.10), we have (A.36),

$$\mathbb{V}_{ij} [\mathbf{z}] \geq \mathbb{C}_{ip} [\mathbf{y}, \mathbf{z}] \mathbb{V}_{pq} [\mathbf{y}]^{-1} \mathbb{C}_{qj} [\mathbf{y}, \mathbf{z}]. \quad (\text{A.40})$$

This completes the proof of the generalized Cauchy-Schwartz inequality.

We now prove (A.34). We use (A.36), and set

$$\mathbf{z} = \mathbf{\Theta}, \quad \mathbf{y} = \mathbf{s}, \quad (\text{A.41})$$

where  $\mathbf{\Theta}$  is an estimator for  $\boldsymbol{\theta}$ , and  $\mathbf{s}$  is the score (see (A.23)). The covariance between  $\mathbf{\Theta}$  and  $\mathbf{s}$  is

$$\begin{aligned} \mathbb{C}_{ij} [\mathbf{\Theta}, \mathbf{s}] &= \mathbb{E} [(\Theta_i - \mu_{\Theta_i}) (s_j - \mu_{s_j})] \\ &= \mathbb{E} [\Theta_i s_j] \\ &= \int dx P(\mathbf{x}; \boldsymbol{\theta}) \Theta_i \nabla_{\theta_j} \ln P(\mathbf{x}; \boldsymbol{\theta}) \\ &= \nabla_{\theta_j} \mathbb{E} [\Theta_i] = \nabla_{\theta_j} \psi_i. \end{aligned} \quad (\text{A.42})$$

We also have

$$\mathbb{V}_{ij} [\mathbf{\Theta}] = \Sigma_{ij}, \quad \mathbb{V}_{ij} [\mathbf{s}] = F_{ij}. \quad (\text{A.43})$$

We have defined  $\Sigma$  to be the covariance matrix of  $\mathbf{\Theta}$ , and used that Fisher information is the variance of the score (see (A.25)). Plugging this all into (A.36), we obtain the Cramér-Rao bound

$$\Sigma_{ij} \Big|_{\boldsymbol{\theta}=\boldsymbol{\mu}_{\boldsymbol{\Theta}}} \geq \nabla_{\theta_p} \psi_i \nabla_{\theta_1} \psi_j F_{pq}^{-1}(\boldsymbol{\theta}). \quad (\text{A.44})$$

# Appendix B

## Fourier and other transforms

### B.1 Brief review of complex analysis

For a function  $A(t)$  that is singular at infinity, the **Cauchy principal value** is defined to be

$$\text{p.v.} \int_{-\infty}^{\infty} dt A(t) = \lim_{T \rightarrow \infty} \int_{-T}^T dt A(t). \quad (\text{B.1})$$

For complex-valued functions  $A(z)$  that are singular at a point  $z_0$ , the Cauchy principal value is defined to be the limit of the deformation of the integral  $C$  by a disk of radius  $\epsilon$  centered around  $z_0$

$$\text{p.v.} \int_C dz A(z) = \lim_{\epsilon \rightarrow 0^+} \int_{C(\epsilon)} dz A(z). \quad (\text{B.2})$$

This can also be written as

$$\text{p.v.} \int_C dz A(z) = \lim_{\epsilon \rightarrow 0^+} \left( \int_{-\infty}^{z_0 - \epsilon} dz A(z) + \int_{z_0 + \epsilon}^{\infty} dz A(z) \right). \quad (\text{B.3})$$

### B.2 The Fourier transform

We briefly review Fourier transforms, along with a helpful transforms that are used in signal processing.

The one-dimensional **Fourier transform** and its inverse are

$$A(t) = \mathcal{F}^{-1} [\tilde{A}(f)](t) = \int_{-\infty}^{\infty} df e^{2\pi i f t} \tilde{A}(f), \quad (\text{B.4a})$$

$$\tilde{A}(f) = \mathcal{F}[A(t)](f) = \int_{-\infty}^{\infty} df e^{-2\pi i f t} A(t). \quad (\text{B.4b})$$

The Fourier representation of the Dirac delta function  $\delta(t)$  is

$$\tilde{\delta}(f) = \int_{-\infty}^{\infty} df e^{-2\pi i f t} \delta(t) = 1. \quad (\text{B.5})$$

The **convolution** of two functions  $A(t)$  and  $B(t)$  are

$$(A * B)(t) \equiv \int_{-\infty}^{\infty} d\tau A(\tau) B(t - \tau) = \int_{-\infty}^{\infty} d\tau A(t - \tau) B(\tau). \quad (\text{B.6})$$

The Fourier transform of the convolution is

$$\begin{aligned} \mathcal{F}[(A * B)(t)](f) &= \int_{-\infty}^{\infty} d\tau \int_{-\infty}^{\infty} df \int_{-\infty}^{\infty} df' e^{2\pi i f \tau} e^{2\pi i f'(t - \tau)} \tilde{A}(f) \tilde{B}(f') \\ &= \int_{-\infty}^{\infty} df e^{2\pi i f t} \tilde{A}(f) \tilde{B}(f). \end{aligned} \quad (\text{B.7})$$

That is, convolution in real space is multiplication in frequency space.

## B.3 The Laplace transform

The **Laplace transform** of a function  $A(t)$  is

$$A(t) = \mathcal{L}^{-1}[A(\lambda)](t) = \frac{1}{2\pi i} \lim_{T \rightarrow \infty} \int_{\gamma - iT}^{\gamma + iT} d\lambda e^{\lambda t} \tilde{A}(\lambda), \quad (\text{B.8a})$$

$$\tilde{A}(t) = \mathcal{L}[A(t)](\lambda) = \int_0^{\infty} dt e^{-\lambda t} A(t), \quad (\text{B.8b})$$

Here  $\gamma$  is a real number so that the contour path of integration is in the region of convergence of  $\tilde{A}(\lambda)$ . In effect, the inverse Laplace transform is like the inverse Fourier transform.

## B.4 The Hilbert transform

The **Hilbert transform** of a function  $A(t)$  is

$$A(t) = \mathcal{H}^{-1}[\tilde{A}(\tau)](t) = -\frac{1}{\pi} \text{p.v.} \int_{-\infty}^{\infty} d\tau \frac{\tilde{A}(\tau)}{t - \tau}, \quad (\text{B.9a})$$

$$\tilde{A}(\tau) = \mathcal{H}[A(t)](\tau) = \frac{1}{\pi} \text{p.v.} \int_{-\infty}^{\infty} dt \frac{A(t)}{\tau - t}. \quad (\text{B.9b})$$

# Appendix C

## Stationary phase approximation

### C.1 Stationary phase approximation

Here we review the **stationary phase approximation** for the Fourier transform. For more discussion see [BO99]. Consider a complex function, which we write as

$$B(t) = A(t) e^{i\phi(t)}. \quad (\text{C.1})$$

The Fourier transform is

$$\tilde{B}(f) = \int_{-\infty}^{\infty} dt A(t) e^{i\phi(t) - 2\pi i f t}. \quad (\text{C.2})$$

We imagine  $A(t)$  is a slowly varying function, while  $\phi(t)$  is rapidly varying. We then expect that the integral for  $\tilde{B}(f)$  will be dominated by the stationary points of  $\phi(t) - 2\pi f t$ , that is the points where

$$\frac{d\phi}{dt} - 2\pi f = 0. \quad (\text{C.3})$$

This can be more formally justified by the Riemann-Lebesgue lemma, which states that

$$\lim_{x \rightarrow \infty} \int_a^b dt e^{ixt} A(t) = 0, \quad (\text{C.4})$$

provided  $\int_a^b dt A(t)$  exists. We can extend  $a, b \rightarrow \pm\infty$  so long as  $A(t)$  is integrable. Going back to (C.2), we assume that  $\phi(t) - 2\pi f t$  has one stationary point for each value of  $f$ , which we call  $t_0(f)$ . That is,  $t_0(f)$  is defined to solve the stationary phase equation

$$\left. \frac{d\phi}{dt} \right|_{t=t_0} - 2\pi f = 0. \quad (\text{C.5})$$

We Taylor series expand about the stationary point to quadratic order in  $\phi$ ,

$$\phi(t) - 2\pi ft = \phi(t_0) - 2\pi ft_0 + \frac{1}{2} \frac{d^2\phi}{dt^2} \Big|_{t=t_0} (t - t_0)^2 + \mathcal{O}[(t - t_0)^3], \quad (\text{C.6})$$

insert this into (C.2), and obtain

$$\begin{aligned} \tilde{B}(f) &\approx A(t_0) e^{i\phi(t_0) - 2\pi i f t_0} \int_{-\infty}^{\infty} dt \exp \left[ i \frac{1}{2} \frac{d^2\phi}{dt^2} \Big|_{t=t_0} (t - t_0)^2 \right] \\ &= \left[ \frac{1}{2} \frac{d^2\phi}{dt^2} \Big|_{t=t_0} \right]^{-1/2} A(t_0) e^{i\phi(t_0) - 2\pi i f t_0 + i\pi/4} \int_{-\infty}^{\infty} dx e^{-x^2} \\ &= \left[ \frac{1}{2\pi} \frac{d^2\phi}{dt^2} \Big|_{t=t_0} \right]^{-1/2} A(t_0) e^{i\phi(t_0) - 2\pi i f t_0 + i\pi/4}. \end{aligned} \quad (\text{C.7})$$

Using the chain rule, we can write  $\phi(t_0)$  and  $t_0$  as integral equations in terms of the frequency. We have

$$t_0(f) = \int^f df' \frac{dt}{df}, \quad (\text{C.8a})$$

$$\phi(t_0) = 2\pi \int^f df' \frac{dt}{df} f'. \quad (\text{C.8b})$$

Defining  $\dot{f} \equiv df/dt$ , we see that we can write the phase of  $\tilde{B}(f)$  as

$$\begin{aligned} \Psi &\equiv \phi(t_0) - 2\pi f t_0 + \frac{\pi}{4} \\ &= 2\pi \int^f df' \frac{1}{\dot{f}} (f' - f) + \frac{\pi}{4}. \end{aligned} \quad (\text{C.9})$$

# Bibliography

- [BGJM11] Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. *Handbook of Markov Chain Monte Carlo*. CRC press, Boca Raton, FL, USA, 2011.
- [BO99] C.M. Bender and S.A. Orszag. *Advanced Mathematical Methods for Scientists and Engineers I: Asymptotic Methods and Perturbation Theory*. Advanced Mathematical Methods for Scientists and Engineers. Springer, 1999.
- [CA11] Jolien D. E. Creighton and Warren G. Anderson. *Gravitational-wave physics and astronomy: An introduction to theory, experiment and data analysis*. Wiley-VCH, 2011.
- [DL70] James M. Dickey and B. P. Lientz. The Weighted Likelihood Ratio, Sharp Hypotheses about Chances, the Order of a Markov Chain. *The Annals of Mathematical Statistics*, 41(1):214 – 226, 1970.
- [Fin92] Lee S. Finn. Detection, measurement and gravitational radiation. *Phys. Rev. D*, 46:5236–5249, 1992.
- [HFM18] David W. Hogg and Daniel Foreman-Mackey. Data analysis recipes: Using Markov Chain Monte Carlo. *Astrophys. J. Suppl.*, 236(1):11, 2018.
- [Mac03] David J. C. Mackay. *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, 2003.
- [PTVF92] William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery. *Numerical Recipes in C*. Cambridge University Press, Cambridge, MA, USA, second edition, 1992.
- [Ski06] John Skilling. Nested sampling for general Bayesian computation. *Bayesian Analysis*, 1(4):833 – 859, 2006.
- [Was10] Larry Wasserman. *All of Statistics : A Concise Course in Statistical Inference*. Springer, New York, 2010.