

# Suitably impressive thesis title

Lennart Golks

Department of Physics

University of Otago

*A thesis submitted for the degree of  
Doctor of Philosophy*

November 2025

## Abstract

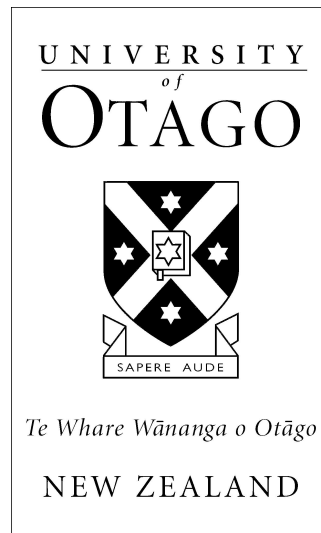
Your abstract text goes here. Check your departmental regulations, but generally this should be less than 300 words. See the beginning of [Chapter 1](#) for more.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Pellentesque sit amet nibh volutpat, scelerisque nibh a, vehicula neque. Integer placerat nulla massa, et vestibulum velit dignissim id. Ut eget nisi elementum, consectetur nibh in, condimentum velit. Quisque sodales dui ut tempus mattis. Duis malesuada arcu at ligula egestas egestas. Phasellus interdum odio at sapien fringilla scelerisque. Mauris sagittis eleifend sapien, sit amet laoreet felis mollis quis. Pellentesque dui ante, finibus eget blandit sit amet, tincidunt eu neque. Vivamus rutrum dapibus ligula, ut imperdiet lectus tincidunt ac. Pellentesque ac lorem sed diam egestas lobortis.

Suspendisse leo purus, efficitur mattis urna a, maximus molestie nisl. Aenean porta semper tortor a vestibulum. Suspendisse viverra facilisis lorem, non pretium erat lacinia a. Vestibulum tempus, quam vitae placerat porta, magna risus euismod purus, in viverra lorem dui at metus. Sed ac sollicitudin nunc. In maximus ipsum nunc, placerat maximus tortor gravida varius. Suspendisse pretium, lorem at porttitor rhoncus, nulla urna condimentum tortor, sed suscipit nisi metus ac risus.



# Suitably impressive thesis title



Lennart Golks  
Department of Physics  
University of Otago

A thesis submitted for the degree of  
*Doctor of Philosophy*  
November 2025



# Acknowledgements

## Personal

I would like to thank Alex Elliott for his wonderful help and support. None of this would be possible otherwise.

This is where you thank your advisor, colleagues, and family and friends.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Vestibulum feugiat et est at accumsan. Praesent sed elit mattis, congue mi sed, porta ipsum. In non ullamcorper lacus. Quisque volutpat tempus ligula ac ultricies. Nam sed erat feugiat, elementum dolor sed, elementum neque. Aliquam eu iaculis est, a sollicitudin augue. Cras id lorem vel purus posuere tempor. Proin tincidunt, sapien non dictum aliquam, ex odio ornare mauris, ultrices viverra nisi magna in lacus. Fusce aliquet molestie massa, ut fringilla purus rutrum consectetur. Nam non nunc tincidunt, rutrum dui sit amet, ornare nunc. Donec cursus tortor vel odio molestie dignissim. Vivamus id mi erat. Duis porttitor diam tempor rutrum porttitor. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Sed condimentum venenatis consectetur. Lorem ipsum dolor sit amet, consectetur adipiscing elit.

Aenean sit amet lectus nec tellus viverra ultrices vitae commodo nunc. Mauris at maximus arcu. Aliquam varius congue orci et ultrices. In non ipsum vel est scelerisque efficitur in at augue. Nullam rhoncus orci velit. Duis ultricies accumsan feugiat. Etiam consectetur ornare velit et eleifend.

Suspendisse sed enim lacinia, pharetra neque ac, ultricies urna. Phasellus sit amet cursus purus. Quisque non odio libero. Etiam iaculis odio a ex volutpat, eget pulvinar augue mollis. Mauris nibh lorem, mollis quis semper quis, consequat nec metus. Etiam dolor mi, cursus a ipsum aliquam, eleifend venenatis ipsum. Maecenas tempus, nibh eget scelerisque feugiat, leo nibh lobortis diam, id laoreet purus dolor eu mauris. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Nulla eget tortor eu arcu sagittis euismod fermentum id neque. In sit amet justo ligula. Donec rutrum ex a aliquet egestas.

## Institutional

If you want to separate out your thanks for funding and institutional support, I don't think there's any rule against it. Of course, you could also just remove the subsections and do one big traditional acknowledgement section.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut luctus tempor ex at pretium. Sed varius, mauris at dapibus lobortis, elit purus tempor neque, facilisis sollicitudin felis nunc a urna. Morbi mattis ante non augue blandit pulvinar. Quisque nec euismod mauris. Nulla et tellus eu nibh auctor malesuada quis imperdiet quam. Sed eget tincidunt velit. Cras molestie sem ipsum, at faucibus quam mattis vel. Quisque vel placerat orci, id tempor urna. Vivamus mollis, neque in aliquam consequat, dui sem volutpat lorem, sit amet tempor ipsum felis eget ante. Integer lacinia nulla vitae felis vulputate, at tincidunt ligula maximus. Aenean venenatis dolor ante, euismod ultrices nibh mollis ac. Ut malesuada aliquam urna, ac interdum magna malesuada posuere.

# Abstract

Your abstract text goes here. Check your departmental regulations, but generally this should be less than 300 words. See the beginning of [Chapter 1](#) for more.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Pellentesque sit amet nibh volutpat, scelerisque nibh a, vehicula neque. Integer placerat nulla massa, et vestibulum velit dignissim id. Ut eget nisi elementum, consectetur nibh in, condimentum velit. Quisque sodales dui ut tempus mattis. Duis malesuada arcu at ligula egestas egestas. Phasellus interdum odio at sapien fringilla scelerisque. Mauris sagittis eleifend sapien, sit amet laoreet felis mollis quis. Pellentesque dui ante, finibus eget blandit sit amet, tincidunt eu neque. Vivamus rutrum dapibus ligula, ut imperdiet lectus tincidunt ac. Pellentesque ac lorem sed diam egestas lobortis.

Suspendisse leo purus, efficitur mattis urna a, maximus molestie nisl. Aenean porta semper tortor a vestibulum. Suspendisse viverra facilisis lorem, non pretium erat lacinia a. Vestibulum tempus, quam vitae placerat porta, magna risus euismod purus, in viverra lorem dui at metus. Sed ac sollicitudin nunc. In maximus ipsum nunc, placerat maximus tortor gravida varius. Suspendisse pretium, lorem at porttitor rhoncus, nulla urna condimentum tortor, sed suscipit nisi metus ac risus.





# Contents

<b>List of Figures</b>	<b>ix</b>
<b>List of Abbreviations</b>	<b>xi</b>
<b>1 Background</b>	<b>1</b>
1.1 Gaussian Markov Random Field . . . . .	1
1.1.1 Markov Random Field . . . . .	2
1.1.2 Multivariate Normal Distribution . . . . .	2
1.1.3 Conditional distribution through Spatial Dependencies . . . .	3
1.2 linear-Gaussian hierarchical Bayesian model . . . . .	5
1.2.1 Monte-Carlo Method Integration . . . . .	8
1.2.2 Markov Chains . . . . .	8
1.3 Sampling from the posterior distribution . . . . .	11
1.3.1 Metropolis-Hastings - Markov chain Monte-Carlo . . . . .	11
1.3.2 Marginal and then Conditional Sampler - MTC . . . . .	12
1.3.3 Adaptive MCMC - GOMOS . . . . .	14
<b>Appendices</b>	
<b>A Posterior of Bayesian Hierarchical model</b>	<b>19</b>
<b>B Convergence of the Metropolis-Hastings</b>	<b>21</b>
<b>C Randomize then Optimize - RTO</b>	<b>23</b>
<b>D Inverting Matrices - QR factorization</b>	<b>25</b>
<b>E Taylor expansion of <math>g(\lambda)</math></b>	<b>27</b>
<b>F Radiation transfer and absorption line shape</b>	<b>29</b>
<b>G whispering gallery resonator</b>	<b>31</b>
<b>References</b>	<b>33</b>



# List of Figures

1.1	Undirected Random Field of seven nodes including maximum cliques.	4
1.2	A hierarchical Bayesian model describes how we model a physical process to observe data. . . . .	6
G.1	whispering gallery resonator . . . . .	32



## List of Abbreviations

<b>i.i.d.</b>	. . . . .	independent and identically distributed
<b>MRF</b>	. . . . .	Markov Random Field
<b>GMRF</b>	. . . . .	Gaussian Markov Random Field
<b>MTC</b>	. . . . .	Marginal Then Conditional sampler
<b>GOMOS</b>	. . . . .	Global Ozone Monitoring by Occultation of Stars
<b>MCMC</b>	. . . . .	Markov Chain Monte-Carlo
<b>MH</b>	. . . . .	Metropolis-Hastings



# 1

## Background

This chapter provides the background knowledge for the applied methods. We try to carefully guide the reader through the various topics so that we build a base to understand the used methodology. Firstly, we discuss Markov Random Fields and their key properties to describe spatially connected variables in Section 1.1.1. In Section 1.1.2 we introduce the normal distribution to describe the random variables in such a field. By specifying the local connections in Section 1.1.3 we globally observe a Gaussian Markov Random Field (GMRF). Next, in Section 1.2 we relate the variables in the GMRF to a measurement process through a Bayesian Model. Given the measurement, we try to find the most likely variables through a Markov chain, as described in Section 1.2.2. Constructing this Markov chain we present the Metropolis-Hastings algorithm in Section 1.3.1, followed by the Marginal and then conditional sampler in Section 1.3.2 and the adaptive Metropolis-Hastings algorithm in Section 1.3.3.

### 1.1 Gaussian Markov Random Field

In this section, we guide the reader through the basic principles of Gaussian Markov Random Fields (GMRF). We discuss some of the key properties of GMRFs and

how to set up such a spatial field. The discussed topics are strongly related to our specific application and for further reading, we refer to the books [1, 2].

### 1.1.1 Markov Random Field

A Markov Random Field (MRF) is based on a graphical model describing spatial dependencies of  $n$  nodes, also described as vertices. We label each node with an integer  $i$  and attach a spatial position  $s^{(i)}$  and a random variable  $x^{(i)}$ , where  $i = 1, \dots, n$ . On a set of all vertices  $\mathcal{V}$  we define an undirected graph  $\mathcal{G}(\mathcal{V}, \mathcal{E})$  with a set of edges  $\mathcal{E}$ . If an edge describes spatial dependencies of a pair of different nodes  $i$  and  $j$  then we write  $i \sim j$ , where  $(i, j) \in \mathcal{E}$ , with no preferred direction. We limit interactions to pair-wise only.

For a subset of nodes  $A = \mathcal{V}^A \subset \mathcal{V}$ , we can define a subgraph  $\mathcal{G}^A(\mathcal{V}^A, \mathcal{E}^A)$ , including edges  $\mathcal{E}^A = \{(i, j) \in \mathcal{E} \text{ and } i, j \in \mathcal{V}^A \text{ and } i \neq j\}$ . A discrete example is illustrated in Figure 1.1.

Through the Markov properties of such a field, a single node  $i$  is conditionally dependent on its neighbors  $\partial i = \{j : (i, j) \in \mathcal{E} \text{ and } j \neq i\}$  only. We denote  $-A = \mathcal{V} - A$  to all nodes not in  $A$  but part of  $\mathcal{V}$ . Similarly, we can extend this to a subset  $A \in \mathcal{V}$  where the neighborhood  $\partial A = \{j : (i, j) \in \mathcal{E} \text{ and } j \notin A\}$  is all nodes  $j \in -A$  sharing an edge to a node  $i \in \mathcal{V}^A$ . Conditioned on all its neighbours  $\mathcal{V}^A$  is independent of all other nodes  $\pi(\mathbf{x}^{(A)} | \mathbf{x}^{(-A)}) = \pi(\mathbf{x}^{(A)} | \mathbf{x}^{(\partial A)})$ . Here  $\pi(\cdot)$  denotes the probability density for its argument. Intuitively we interpret  $\mathbf{Q}$  in terms of conditional dependencies and  $\Sigma$  marginally, which reduces distribution to maximal  $n$  dimensional.

### 1.1.2 Multivariate Normal Distribution

In our case, we specifically choose a probability density that distributes each random variable  $x^{(i)}$  normally and denote this with  $x^{(i)} \sim \mathcal{N}(\mu^{(i)}, \sigma)$ . Here  $\mu_i = \mathbb{E}(x^{(i)})$  is the mean of  $x^{(i)}$  and  $\sigma^2 = \text{Var}(x^{(i)})$  is the variance of any  $x^{(i)}$ .



Constructing a Gaussian Markov Random Field (GMRF), we introduce a normal probability density function  $\pi$  for a single node  $x_i$ , with mean  $\mu$  and variance  $\sigma^2$ .

$$\pi(x^{(i)}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(x^{(i)} - \mu^{(i)})^2}{2\sigma^2} \right] \quad (1.1)$$

To describe the values of all nodes in a GMRF we introduce a multivariate Gaussian

$$\mathbf{x} = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(n)} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mu^{(1)} \\ \vdots \\ \mu^{(n)} \end{bmatrix}, \mathbf{\Sigma} \right), \quad (1.2)$$

with a covariance matrix  $\mathbf{\Sigma} = \mathbf{Q}^{-1}$ , which is the inverse of the precision-matrix. As we deal with an undirected Markov field,  $\mathbf{Q}$  is a symmetric positive-definite matrix and sparse. Hence, we can efficiently decompose  $\mathbf{Q} = \mathbf{L}\mathbf{L}^T$  into its Cholesky triangle  $\mathbf{L}$ , which is a lower triangular matrix. Within this precision matrix, we can define the interactions and correlations between the different nodes of our GMRF.

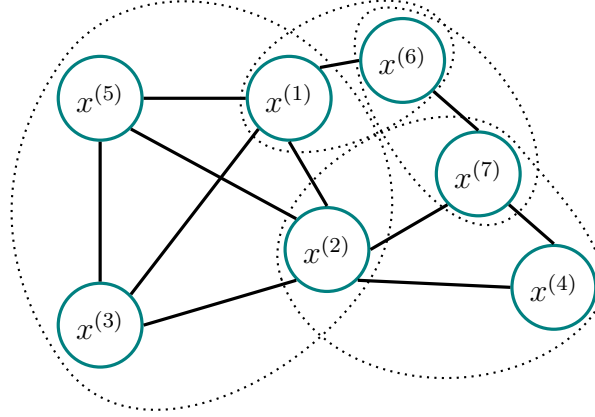
When constructing a GMRF the precision matrix  $\mathbf{Q}$  as well as its inverse  $\mathbf{\Sigma}$  does a play key role. If the GMRF is fully connected then  $\mathbf{Q}$  is completely dense, we aim to construct a GMRF so that the precision matrix is sparse. The off diagonal non-zero entries of  $\mathbf{Q}$  provide information about conditional correlations  $\text{Corr}(x^{(i)}, x^{(j)} | \mathbf{x}^{(-ij)}) = Q_{ij} / \sqrt{Q_{ii}Q_{jj}}$ . Where as the off diagonal entries of  $\mathbf{\Sigma}$  give information about marginal correlations  $\text{Corr}(x^{(i)}, x^{(j)}) = \Sigma_{ij} / \sqrt{\Sigma_{ii}\Sigma_{jj}}$ . The diagonal entries provide marginal variances  $\text{Var}(x^{(i)}) = \Sigma_{ii}$  and conditional precision  $\text{Prec}(x^{(i)}, x^{(j)} | \mathbf{x}^{(-ij)}) = Q_{ii}$ .

### 1.1.3 Conditional distribution through Spatial Dependencies

In this subsection, we present the Conditional Auto Regressive (CAR) model and the Gibbs field to define a graph  $\mathcal{G}(\mathcal{V}, \mathcal{E})$  and the corresponding precision matrix  $\mathbf{Q}$ . For curious readers we recommend the book [3].

The CAR method can include neighbouring nodes of numerous orders. We define the entries of the precision matrix as follows:

$$Q_{ij} = \begin{cases} \kappa_i & i = j \\ \kappa_i \beta_{ij} & i \neq j \end{cases}, \quad (1.3)$$



**Figure 1.1: Undirected Random Field of seven nodes including maximum cliques.** The random field displayed here can be described with a Graph  $\mathcal{G}$ . Where  $\mathcal{V}$  are all nodes seen in the figure and  $\mathcal{E}$  are the edges drawn with a line between nodes. We draw dotted lines around the set of maximum cliques, which are not a full subset of any other clique. Maximum cliques are a set of fully connected nodes with common neighbors. A neighborhood is defined through an edge between two nodes, such as  $\{x^{(1)}, x^{(6)}\} \in \mathcal{E}$ . Here  $\{\{x^{(1)}, x^{(6)}\}, \{x^{(1)}, x^{(3)}, x^{(4)}, x^{(5)}\}, \{x^{(6)}, x^{(7)}\}, \{x^{(2)}, x^{(4)}, x^{(7)}\}\}$  form the set of maximum cliques. On a set of maximum cliques, we can define an energy function, Gibbs distribution or use the [conditional autoregressive model](#) to define spatial dependencies. The precision matrix  $\mathbf{Q}$  represents those spatial dependencies weighted by the hyper parameter  $\boldsymbol{\theta}$ .

where  $\beta_{ij} = 0$ , if  $x^{(i)}$  and  $x^{(j)}$  are conditionally independent, and  $\kappa_i > 0$  to ensure positive-definiteness of  $\mathbf{Q}$  for all  $i = 1, \dots, n$ . The full conditionals are normally distributed

$$x^{(i)} | x^{(-i)} = x^{(i)} | \partial x^{(i)} \sim \mathcal{N} \left( \underbrace{\sum_{j:j \neq i} \beta_{ij} x^{(j)}}_{\text{mean}}, \underbrace{\kappa_i^{-1}}_{\text{variance}} \right) \quad (1.4)$$

and depend only on neighbouring nodes  $\partial x^{(i)}$  having an edge to  $x^{(i)}$ . We can then write the probability density over all nodes  $\mathbf{x}$  including their means  $\boldsymbol{\mu}$  as follows:

$$\pi(\mathbf{x}) = \frac{1}{(2\pi)^{n/2}} \sqrt{\det(\mathbf{Q})} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{Q} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (1.5)$$

$$= \frac{1}{(2\pi)^{n/2}} |\boldsymbol{\Sigma}|^{-1/2} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]. \quad (1.6)$$

Alternatively, we can describe a set of nodes through a Gibbs field, where we define a clique  $c \in \mathcal{C}$  as a set of fully connected nodes. In a clique each node  $j \in c$

has mutual neighbours to any other node  $i \neq j \in c$ . If for any node  $j \notin c$  and  $j \in \mathcal{V}$ , the union of  $j \cup c$  is not a clique then we call  $\tilde{c} \in \mathcal{C}$  a maximal clique. We conclude that maximal cliques are not a complete subset of another clique. The positive clique-potential  $\phi_{\tilde{c}}(\tilde{c})$  describes the interactions of all nodes within a clique. The sum over all maximum clique-potentials is often described as the Energy:

$$E = \sum_{\tilde{c} \in \mathcal{C}} \phi_{\tilde{c}}(x^{(i)} : x^{(i)} \in c). \quad (1.7)$$

We normalize the Gibbs distribution  $\pi(\mathbf{x})$  of a random field with a normalization constant  $Z$ .

$$\pi(\mathbf{x}) = \frac{1}{Z} \exp \left[ - \sum_{\tilde{c} \in \mathcal{C}} \phi_{\tilde{c}}(x^{(i)} : x^{(i)} \in \tilde{c}) \right] \quad (1.8)$$

$$Z = \sum_{i \in \mathcal{V}} \prod_{\tilde{c} \in \mathcal{C}} \exp \left[ - \phi_{\tilde{c}}(x^{(i)} : x^{(i)} \in \tilde{c}) \right] \quad (1.9)$$

The Hammersley-Clifford theorem states that if we describe any Random Field through a Gibbs distribution over maximum cliques then we deal with a MRF. First proven by Julien Besag we can show that conditional joint distribution is expressed as

$$\pi(x^{(i)} | \partial x^{(i)}) = \frac{1}{Z} \exp \left[ - \sum_{\tilde{c} \in \mathcal{C}} \phi_{\tilde{c}}(x^{(i)}, x^{(j)} : i, j \in \tilde{c} \text{ and } i \neq j) \right] \quad (1.10)$$

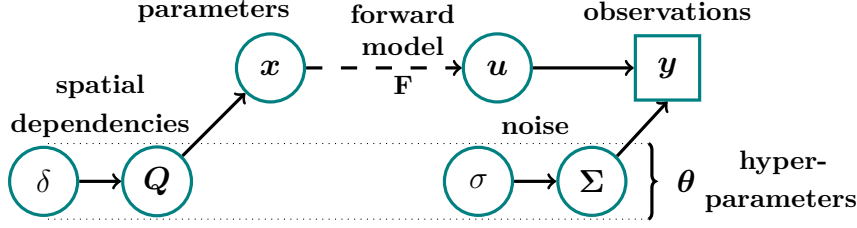
$$Z = \sum_{i \in \mathcal{V}} \prod_{\tilde{c} \in \mathcal{C}} \exp \left[ - \phi_{\tilde{c}}(x^{(i)}, x^{(j)} : i, j \in \tilde{c} \text{ and } i \neq j) \right] \quad (1.11)$$

and can be found in [4]. The potential  $\phi_{\tilde{c}}$  can describe the interactions between different nodes of a Gaussian Gibbs field if set accordingly, which leads to the precision-matrix  $\mathbf{Q}$  for a Gaussian Markov Random Field (GMRF).

## 1.2 linear-Gaussian hierarchical Bayesian model

Bayesian hierarchical models are a very helpful tool to describe a measurement process through different layers of complexity. It allows us to backtrack from our observations to the source of the measurement in a **very precise/simplified** way. Here we closely follow the terminology of [5, 6]

Usually, a measurement device delivers some data  $\mathbf{y}$  including some noise  $\sigma$  and we like to model the observation to quantify some unknown parameters  $\mathbf{x}$



**Figure 1.2: A Bayesian hierarchical model describes how we model a physical process to observe data.** Through a forward model we describe a physical process depending on some parameters  $\mathbf{x}$ . The spatial dependencies of those parameters  $\mathbf{x}$  are described through the hyper-parameters  $\boldsymbol{\theta} = (\delta, \sigma)$ . The precision matrix  $\mathbf{Q}$  is defined according to our Markov Random field and the amount of interaction between certain nodes is described by  $\delta$ . The forward map  $\mathbf{F}$  models a physical process and maps the parameters  $\mathbf{x}$  to a space of all measurable  $\mathbf{u}$ . From this space  $\mathbf{u}$ , we observe some data including some random noise  $\gamma$  with variance  $\boldsymbol{\Sigma}(\gamma)$ . A Bayesian framework allows us to find the most likely parameters and hyper-parameters  $\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}$  given our measurement .

and their unknown spatial dependencies. Through the precision matrix  $\mathbf{Q}(\delta)$ , we define neighborhoods and weight pair-wise interactions of parameters  $x^{(j)}$  and  $x^{(i)}$  according to a hyper-prior distribution  $\pi(\delta)$  in (1.12c). The inverse of the precision matrix  $\mathbf{Q}(\delta)$  is the variance of the prior distribution  $\pi(\mathbf{x} | \mathbf{Q}^{-1}(\delta))$  in (1.12b), and usually sparse. Based on this graphical model we map  $\mathbf{x}$  through a forward model  $\mathbf{F}$  **Dimensions** to space of all measurables  $\mathbf{u}$ . Next, we observe some data  $\mathbf{y} = \mathbf{F}\mathbf{x} + \gamma$  and add some unknown normally distributed random noise  $\eta$  with variance  $\boldsymbol{\Sigma}(\gamma)$ . We assume that this noise vector is independent and identically distributed (i.i.d.). The likelihood function  $\pi(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})$  in (1.12a) describes how likely/realistic the observations  $\mathbf{y}$  are. In our specific case, we choose a linear forward model and a Gaussian prior, so that our likelihood function and posterior distribution in (1.12d) are normally distributed as well. In doing so we can find an analytic expression for the uncertainty of the posterior and see that it is depending on both the prior and the likelihood distribution, for more details we refer to the Appendix A.

In Equations (1.12b) - (1.12c), we summarize a hierarchical Bayesian model in a generalized form, with the hyper-parameter  $\boldsymbol{\theta} = (\delta, \gamma)$ . Note that in this case we define two hyper-parameters, but we can increase the dimensions of  $\boldsymbol{\theta}$  to an

arbitrary number suitable to describe the underlying physical process sufficiently.

$$\mathbf{y}|\mathbf{x}, \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{F}\mathbf{x}, \boldsymbol{\Sigma}(\boldsymbol{\theta})) \quad (\text{likelihood}) \quad (1.12a)$$

$$\mathbf{x}|\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{Q}^{-1}(\boldsymbol{\theta})) \quad (\text{prior}) \quad (1.12b)$$

$$\boldsymbol{\theta} \sim \pi(\boldsymbol{\theta}) \quad (\text{hyper-prior}) \quad (1.12c)$$

The posterior has the following from

$$\mathbf{x}, \boldsymbol{\theta}|\mathbf{y} \sim \mathcal{N}(\mu_{\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}}, (\mathbf{Q} + \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F})^{-1}) \quad (\text{posterior}) \quad (1.12d)$$

with the mean  $\mu_{\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}} = \boldsymbol{\mu} + (\mathbf{Q} + \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F})^{-1} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{F}\boldsymbol{\mu})$ .

Using Bayes theorem we can find an expression for the posterior distribution  $\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$  for the parameters  $\mathbf{x}$  and the hyper-parameters  $\boldsymbol{\theta}$  given the observed data  $\mathbf{y}$ .

$$\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) = \frac{\pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}, \boldsymbol{\theta})}{\pi(\mathbf{y})} \quad (1.13)$$

If we like to compute the posterior distribution we have to compute the normalization constant  $\pi(\mathbf{y})$ , where we marginalize out  $\mathbf{x}$  and  $\boldsymbol{\theta}$  :

$$\pi(\mathbf{y}) = \int \int \pi(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} d\boldsymbol{\theta}, \quad (1.14)$$

which can yield a extensive computation. Using the posterior density we can find the expected value of any function  $h(\mathbf{x})$  though an average of the function weighted by the probability density distribution  $\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ .

$$\mathbb{E}_{\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}}[h(\mathbf{x})] = \int \int h(\mathbf{x}) \pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) d\mathbf{x} d\boldsymbol{\theta} \quad (1.15)$$

Taking a look at the posterior we can see that  $\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ .. and see that to solve this might be very tidious, where as the ratio is much easier...

$$\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) \propto \pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) \quad (1.16)$$

$$= \frac{\pi(\boldsymbol{\theta})}{\sqrt{\det(\boldsymbol{\Sigma}) \det(\mathbf{Q})^{-1}}} \exp \left[ -\frac{(\mathbf{F}\boldsymbol{\mu} - \mathbf{y})^T (\mathbf{F}\boldsymbol{\mu} - \mathbf{y})}{\boldsymbol{\Sigma}} + \frac{(\boldsymbol{\mu} - \mathbf{x})^T (\boldsymbol{\mu} - \mathbf{x})}{\mathbf{Q}^{-1}} \right] \quad (1.17)$$

If we set  $h(\mathbf{x}) = \mathbf{x}$  we get the mean of the parameter  $\mathbf{x}$  or the posterior covariance if  $h(\mathbf{x}) = (\mathbf{x} - \mathbb{E}_{\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}}[\mathbf{x}])(\mathbf{x} - \mathbb{E}_{\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}}[\mathbf{x}])^T$ . Computing that is very expensive and usually not feasible as the parameters can take too many different values, depending on the range and complexity of the problem. Instead, we can explore the parameter space for a few (e.g.  $10^4$ ) discrete values and use the Monte-Carlo method to approximate the integral.

### 1.2.1 Monte-Carlo Method Integration

The Monte-Carlo method was first developed in Los Alamos, United States of America, just after the II world war to simulate the flight path of neutrons. Later in 1949, N. Metropolis and S. Ulam formulated their ideas already suggesting to use 'Markoff' chains in Monte Carlo simulations to approximate continuous functions [7]. For more details on the fundamental theorems we use in this work, we refer to the books [8, 9]

Using the Monte-Carlo method we can approximate the expected value of any function  $h(\mathbf{x})$  by the sample mean:

$$\mathbb{E}_{\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}}[h(\mathbf{x})] = \int \int \mathbf{x} \pi(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}) d\mathbf{x} d\boldsymbol{\theta} \approx \frac{1}{N+1} \sum_{i=0}^N \mathbf{x}_i, \quad (1.18)$$

where we draw  $n$  parameter samples  $\mathbf{x}_i$  from the posterior distribution  $\pi(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})$ . We can argue with the central limit theorem that if  $N \rightarrow \infty$  the sequence of  $\{\mathbf{x}_0, \dots, \mathbf{x}_N\}$  converges the sample mean and sample variance to the true mean and true covariance of  $\pi(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})$ . In our case, we know that the posterior is normally distributed and the task now is to draw samples from this distribution to perform such an integration. The problem is that it is not feasible to compute the posterior distribution due to the size of the parameter space. Instead, we construct a Markov chain of randomly picked parameters, which converges to the posterior distribution.

### 1.2.2 Markov Chains

First formulated by Andrey Markov and published in 1906, Markov chains are a very useful statistical tool to describe real-world processes [10]. In this section, we

discuss some of the properties of Markov chains to make sure that we sample from a unique equilibrium distribution, which is in our case the posterior distribution  $\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ . We like to show aperiodicity, irreducibility, and convergence towards a stationary distribution by proofing the detailed balance condition. If all of these properties hold we call a chain ergodic with a unique equilibrium distribution, e. g.  $\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ . For more detailed information we refer to [3, 11].

In a finite parameter space  $\Omega(\mathcal{X}, \theta)$  we draw a sequence of random variables  $\mathcal{M} = \{\{\mathbf{x}_0, \boldsymbol{\theta}_0\}, \{\mathbf{x}_1, \boldsymbol{\theta}_1\}, \dots, \{\mathbf{x}_N, \boldsymbol{\theta}_N\}\}$  from distributions  $\pi^{(0)}, \dots, \pi^{(N)}$ . The Markov condition requires that each transition probability  $\Pr(\cdot)$  of a new state  $\{\mathbf{x}_{N+1}, \boldsymbol{\theta}_{N+1}\}$  only depends on the current state  $\{\mathbf{x}_N, \boldsymbol{\theta}_N\}$ :

$$\Pr(\{\mathbf{x}_0, \boldsymbol{\theta}_0\}, \{\mathbf{x}_1, \boldsymbol{\theta}_1\}, \dots, \{\mathbf{x}_N, \boldsymbol{\theta}_N\}) = \Pr(\{\mathbf{x}_{N+1}, \boldsymbol{\theta}_{N+1}\} | \{\mathbf{x}_N, \boldsymbol{\theta}_N\}). \quad (1.19)$$

For reasons of simplicity in terminology we denote the state  $\{\mathbf{x}_N, \boldsymbol{\theta}_N\}$  as vector  $(\mathbf{x}_N^T, \boldsymbol{\theta}_N^T)^T = \mathbf{i}$ . The probability for the Markov chain to be in that state  $\mathbf{i}$  after  $N$  steps is denoted as  $\pi_i^{(N)} = \Pr((\mathbf{x}_N^T, \boldsymbol{\theta}_N^T)^T = \mathbf{i})$ . We arrange the probability distribution of the  $N$ th step as a row vector representing the probabilities to be in different states  $\mathbf{i}, \mathbf{j} \in \Omega(\mathcal{X}, \theta)$

$$\pi^{(N)} = [\dots \pi_i^{(N)} \pi_j^{(N)} \dots] \quad (1.20)$$

A Markov chain  $\mathcal{M}$  is homogeneous if the transition probability only depends on the value of the current and proposed state and not on  $n$ , the position in the chain.

$$\Pr(\{\mathbf{x}_{N+M+1}, \boldsymbol{\theta}_{N+M+1}\} | \{\mathbf{x}_{N+M}, \boldsymbol{\theta}_{N+M}\}) = \Pr(\{\mathbf{x}_{N+1}, \boldsymbol{\theta}_{N+1}\} | \{\mathbf{x}_N, \boldsymbol{\theta}_N\}) \quad \forall M \in \mathbb{Z} \quad (1.21)$$

Then we denote the probability to transition from state  $\{\mathbf{x}_N, \boldsymbol{\theta}_N\} = \mathbf{i}$  to  $\{\mathbf{x}_{N+1}, \boldsymbol{\theta}_{N+1}\} = \mathbf{j}$  as an element of the transition matrix  $\mathbf{P}$ :

$$P_{ij} = \Pr((\mathbf{x}_{N+1}^T, \boldsymbol{\theta}_{N+1}^T)^T = \mathbf{j} | (\mathbf{x}_N^T, \boldsymbol{\theta}_N^T)^T = \mathbf{i}) \quad (1.22)$$

This matrix represents the transitions for one step from As we draw more and more states, some of those states are more likely to occur than others. The Chain will reach an equilibrium distribution  $\pi^{(N)} \rightarrow \pi$  as  $N \rightarrow \infty$ . We call  $\pi$  the stationary

distribution for the transition matrix  $\mathbf{P}$ . Here,  $\pi$  is the left eigenvector of  $\mathbf{P}$  with the eigenvalue one such that

$$\pi = \pi \mathbf{P}. \quad (1.23)$$

Once we draw values of that chain from the stationary distribution, the distribution does not change.

A Markov chain is irreducible if all states in  $\Omega$  intercommunicate. That means for any two states  $(\mathbf{x}_N^T, \boldsymbol{\theta}_N^T)^T = \mathbf{i}$  and  $(\mathbf{x}_{N+1}^T, \boldsymbol{\theta}_{N+1}^T)^T = \mathbf{j}$  we can find a path with non-zero probability linking  $\mathbf{i} \rightarrow \mathbf{j}$  and a path linking  $\mathbf{i} \leftarrow \mathbf{j}$ . Then the state space  $\Omega$  is irreducible under  $\mathbf{P}$ . Within an irreducible chain, all states have the same period.

A period is the number of steps for a chain to revisit a set of states starting from the same set of states. An aperiodic chain has period one. If we are allowed to stay in a state  $i$ , so that the self-transition  $P_{i,i} > 0$ , we can break any periodic pattern and the chain is aperiodic.

If a Markov chain is aperiodic and irreducible then this chain is ergodic. For an ergodic Markov chain on a finite state space  $\Omega$ , there exists a stationary distribution  $\pi$ . This stationary distribution is unique if it satisfies the detailed balance condition.

$$\pi(\{\mathbf{x}_{N+1}, \boldsymbol{\theta}_{N+1}\} = \mathbf{j})P_{j,i} = \pi(\{\mathbf{x}_N, \boldsymbol{\theta}_N\} = \mathbf{i})P_{i,j} \quad (1.24)$$

$$\pi_j P_{j,i} = \pi_i P_{i,j} \quad (1.25)$$

In a large class of Markov chains the detailed balance condition is very helpful to find the stationary distribution.

As a consequence of the ergodic theorem, we observe an unique equilibrium distribution  $\pi^{(N)} \rightarrow \pi$  as  $N \rightarrow \infty$  independent of the initial distribution  $\pi^{(0)}$ . Now, we sample from the posterior distribution and can calculate the sample mean, as seen in Equation (1.18).

Next, we like to generate such a chain and define the transition probabilities using the Metropolis-Hastings algorithm.



## 1.3 Sampling from the posterior distribution

In this Section we present the major sampling algorithms which allow us to generate a Markov chain of posterior samples and to characterize the posterior distribution. We introduce the Metropolis-Hastings, Gibbs, MTC sampler, T-walk, adaptive MCMC.

### 1.3.1 Metropolis-Hastings - Markov chain Monte-Carlo

Here we introduce the Metropolis-Hastings algorithm, which is a Markov-chain Monte Carlo (MCMC) algorithm. First published in 1970 and based on the work of Metropolis et. al in 1953 this algorithm provides a framework to find samples from the posterior efficiently [12, 13]. Generating a Markov-Chain  $\{\mathbf{x}_0, \boldsymbol{\theta}_0\}, \dots, \{\mathbf{x}_j, \boldsymbol{\theta}_j\}, \dots, \{\mathbf{x}_N, \boldsymbol{\theta}_N\}$  we accept and reject proposed samples to accurately calculate the sample mean.

**Algorithm 1:** Metropolis-Hastings step to generate a new candidate in a Markov chain

**Let**  $\{\mathbf{x}_j, \boldsymbol{\theta}_j\}$  be the current state , then we **generate**  $\{\mathbf{x}_{j+1}, \boldsymbol{\theta}_{j+1}\}$  as follows:  
**1 Draw:** new state  $\{\mathbf{x}', \boldsymbol{\theta}'\} \sim g(\mathbf{x}', \boldsymbol{\theta}' | \mathbf{x}_j, \boldsymbol{\theta}_j)$   
**2 Acceptance probability:**  

$$\alpha(j+1|j) \equiv \min \left\{ 1, \frac{\pi(\mathbf{x}', \boldsymbol{\theta}' | \mathbf{y}) g(\mathbf{x}_j, \boldsymbol{\theta}_j | \mathbf{x}', \boldsymbol{\theta}')}{\pi(\mathbf{x}_{j+1}, \boldsymbol{\theta}_{j+1} | \mathbf{y}) g(\mathbf{x}', \boldsymbol{\theta}' | \mathbf{x}_j, \boldsymbol{\theta}_j)} \right\}$$
  
**3 Draw:**  $u \sim \mathcal{U}(0, 1)$   
**4 if**  $u \leq \alpha(j+1|j)$  **then**  
**5 |   Accept:**  $\{\mathbf{x}_{j+1}, \boldsymbol{\theta}_{j+1}\} = \{\mathbf{x}', \boldsymbol{\theta}'\};$   
**6 |   else Reject:**  $\{\mathbf{x}_{j+1}, \boldsymbol{\theta}_{j+1}\} = \{\mathbf{x}_j, \boldsymbol{\theta}_j\};$   
**7 end**

We generate a new state  $\{\mathbf{x}_{j+1}, \boldsymbol{\theta}_{j+1}\}$  in our Markov-Chain by proposing a state  $\{\mathbf{x}', \boldsymbol{\theta}'\}$  as described in Algorithm 1 we introduce the acceptance probability  $\alpha(j+1|j)$  and the probability to generate a new state  $g(\mathbf{x}', \boldsymbol{\theta}' | \mathbf{x}_j, \boldsymbol{\theta}_j)$  given a current state  $\{\mathbf{x}_j, \boldsymbol{\theta}_j\}$ . The transition probability  $P_{j,j+1}$  becomes:

$$\Pr(\{\mathbf{x}_{j+1}, \boldsymbol{\theta}_{j+1}\} | \{\mathbf{x}_j, \boldsymbol{\theta}_j\}) = g(\mathbf{x}_{j+1}, \boldsymbol{\theta}_{j+1} | \mathbf{x}_j, \boldsymbol{\theta}_j) \alpha(\mathbf{x}_{j+1}, \boldsymbol{\theta}_{j+1} | \mathbf{x}_j, \boldsymbol{\theta}_j) \quad (1.26)$$

Further, we define the acceptance probability  $\alpha(j+1|j)$  as:

$$\alpha(j+1|j) \equiv \min \left\{ 1, \frac{\pi(\mathbf{x}', \boldsymbol{\theta}' | \mathbf{y}) g(\mathbf{x}_j, \boldsymbol{\theta}_j | \mathbf{x}', \boldsymbol{\theta}')}{\pi(\mathbf{x}_{j+1}, \boldsymbol{\theta}_{j+1} | \mathbf{y}) g(\mathbf{x}', \boldsymbol{\theta}' | \mathbf{x}_j, \boldsymbol{\theta}_j)} \right\} \quad (1.27)$$

Here we observe that the acceptance probability is constructed in such a way that we do not need to compute the full posterior distribution. It is sufficient to sample from a not normalized posterior distribution  $\tilde{\pi}(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$  if we assume that the normalization constant  $\pi(\mathbf{y}) > 0$ .

$$\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) \propto \pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})\pi(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta}) = \tilde{\pi}(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) \quad (1.28)$$

Next, we draw a random number from a uniform distribution between 0 and 1. If this random number is smaller or equal then the acceptance probability of the proposed candidate, we accept this proposed state and set  $\{\mathbf{x}_{j+1}, \boldsymbol{\theta}_{j+1}\} = \{\mathbf{x}', \boldsymbol{\theta}'\}$ . If the uniformly drawn number is larger then we reject the proposal to stay in the current state. We can repeat this step for large enough  $n$  so that we generate an ergodic Markov chain. Sampling the hyper-parameters  $\boldsymbol{\theta}$  and the parameters  $\mathbf{x}$  is computationally very time consuming, to speed up the process we introduce the Marginal and then Conditional sample in the next Section.

Hence we like to sample from a very specific posterior distribution we show that the Markov chain, generated by the Metropolis-Hastings algorithm, fulfills the conditions for ergodicity in the Appendix B .

### 1.3.2 Marginal and then Conditional Sampler - MTC

Dealing with a hierarchical Bayesian model we have to sample  $\mathbf{x}$  and  $\boldsymbol{\theta}$  from the posterior distribution  $\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$  to find the mode of this distribution. Instead of computing full posterior samples  $\{\mathbf{x}_i, \boldsymbol{\theta}_i\}$ , we can speed up the sampling process by integration out  $\mathbf{x}$  and independently sample hyper-parameters from the marginal posterior distribution  $\boldsymbol{\theta}_i \sim \pi(\boldsymbol{\theta}|\mathbf{y})$  directly. Then we sample  $\mathbf{x}_i \sim \pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_i)$  to characterize the posterior distribution  $\{\mathbf{x}_i, \boldsymbol{\theta}_i\} \sim \pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ .

### Sampling from the marginal posterior $\pi(\boldsymbol{\theta}|\mathbf{y})$ of the hyper-parameters $\boldsymbol{\theta}$

For the linear Bayesian hierarchical model we can eliminate  $\mathbf{x}$ , so that the marginal posterior distribution is given by: [fix sqrt line](#)

$$\pi(\boldsymbol{\theta}|\mathbf{y}) = \int \pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) d\mathbf{x} \quad (1.29)$$

$$\propto \sqrt{\frac{\det(\boldsymbol{\Sigma}^{-1}) \det(\mathbf{Q})}{\det(\mathbf{Q} + \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F})}} \times \exp \left[ -\frac{1}{2}(\mathbf{y} - \mathbf{F}\boldsymbol{\mu})^T \mathbf{Q}_{\boldsymbol{\theta}|\mathbf{y}} (\mathbf{y} - \mathbf{F}\boldsymbol{\mu}) \right] \pi(\boldsymbol{\theta}), \quad (1.30)$$

with the precision matrix

$$\mathbf{Q}_{\boldsymbol{\theta}|\mathbf{y}} = \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} + \mathbf{Q})^{-1} \mathbf{F}^T \boldsymbol{\Sigma}^{-1}. \quad (1.31)$$

Note that this distribution is not Gaussian, [5].

Next, we generate a Markov chain  $\{\boldsymbol{\theta}_0, \dots, \boldsymbol{\theta}_j, \dots, \boldsymbol{\theta}_{K'}\}$  utilizing an MCMC algorithm on the conditional posterior distribution  $\pi(\boldsymbol{\theta}|\mathbf{y})$ . At state  $\boldsymbol{\theta}_j = \boldsymbol{\theta}$  we propose a new hyper-parameter sample  $\boldsymbol{\theta}'$  from the proposal distribution  $g(\boldsymbol{\theta}'|\boldsymbol{\theta})$  and accept the new state with the probability according to the Metropolis-Hastings ratio:

$$1 \wedge \frac{\pi(\boldsymbol{\theta}'|\mathbf{y})g(\boldsymbol{\theta}|\boldsymbol{\theta}')}{\pi(\boldsymbol{\theta}|\mathbf{y})g(\boldsymbol{\theta}'|\boldsymbol{\theta})} \quad (1.32)$$

Note that we have to compute the ratio  $\pi(\boldsymbol{\theta}'|\mathbf{y})/\pi(\boldsymbol{\theta}|\mathbf{y})$ . The MTC sampler is especially powerful in case of cheap evaluation of the determinants of the precision matrices, see Equation 1.30.

Once the Markov Chain  $\{\boldsymbol{\theta}_0, \dots, \boldsymbol{\theta}_j, \dots, \boldsymbol{\theta}_{K'}\}$  is long enough we calculate the integrated auto-correlation time  $\tau_{\text{int}}$  of that chain. According to this measure and an appropriate burn in period we can refine the previous Markov chain to  $\{\boldsymbol{\theta}_0, \dots, \boldsymbol{\theta}_j, \dots, \boldsymbol{\theta}_K\}$ , so that we have  $K < K'$  independent samples of the marginal posterior  $\pi(\boldsymbol{\theta}|\mathbf{y})$ .

### Sampling from the full conditional $\pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}^{(i)})$

We can draw a sample  $\boldsymbol{\theta}_i$  from the marginal by randomly choosing a state of just generated Markov chain  $\{\boldsymbol{\theta}_0, \dots, \boldsymbol{\theta}_j, \dots, \boldsymbol{\theta}_K\} \sim \pi(\boldsymbol{\theta}|\mathbf{y})$ . Conditioned on  $\pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}_i)$  we use the Randomize Then Optimize (RTO) method by Bardsley et. al. to

draw a parameter sample  $\mathbf{x}_i|\mathbf{y}, \boldsymbol{\theta}_i$  [14, 15]. Here, we like to point out that the method has been introduced under various names and refer to Oliver et. al. and Oriuex et. al. for further reading [16, 17].

The conditional posterior is defined through the linear-Gaussian hierarchical Bayesian model as:

$$\mathbf{x}|\mathbf{y}, \boldsymbol{\theta} \sim \mathcal{N}\left(\boldsymbol{\mu} + (\mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} + \mathbf{Q})^{-1} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{F} \boldsymbol{\mu}), (\mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} + \mathbf{Q})^{-1}\right), \quad (1.33)$$

for more details we refer to [2, 5, 6]. As the full conditional distribution for  $\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}$  is a normal distribution we can rewrite to:

$$\pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}) \propto \pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}|\boldsymbol{\theta}) \quad (1.34)$$

$$= \exp\|\hat{\mathbf{F}}\mathbf{x} - \hat{\mathbf{y}}\|^2, \quad (1.35)$$

where

$$\hat{\mathbf{F}} = \begin{bmatrix} \boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\theta}) \mathbf{F} \\ \mathbf{Q}^{1/2}(\boldsymbol{\theta}) \end{bmatrix}, \quad \hat{\mathbf{y}} = \begin{bmatrix} \boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\theta}) \mathbf{y} \\ \mathbf{Q}^{1/2}(\boldsymbol{\theta}) \boldsymbol{\mu} \end{bmatrix}. \quad (1.36)$$

One sample from the posterior can be computed by minimizing the following equation with respect to  $\hat{\mathbf{x}}$  :

$$\mathbf{x}_i = \arg \min_{\hat{\mathbf{x}}} \|\hat{\mathbf{F}}\hat{\mathbf{x}} - (\hat{\mathbf{y}} + \boldsymbol{\eta})\|^2, \quad \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (1.37)$$

where we add a randomized perturbation  $\boldsymbol{\eta}$ . Next, we substitute  $-\hat{\mathbf{F}}^T \boldsymbol{\eta} = \mathbf{v}_1 + \mathbf{v}_2$  we can rewrite the argument of Eq. 1.36 to

$$(\mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} + \mathbf{Q}) \mathbf{x}_i = \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} + \mathbf{Q} \boldsymbol{\mu} + \mathbf{v}_1 + \mathbf{v}_2, \quad (1.38)$$

where  $\mathbf{v}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F})$  and  $\mathbf{v}_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$  are independent random variables. Finally, we can draw an independent sample from the posterior  $(\mathbf{x}_i, \boldsymbol{\theta}_i) \sim \pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ .

### 1.3.3 Adaptive MCMC - GOMOS

**Algorithm 2:** Marginal and then Conditional (MTC) Sampler - Linear Gaussian Model

```

1 foreach  $\boldsymbol{\theta} = \boldsymbol{\theta}_j \in \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{K'}\}$  do
2   Propose new state:  $\boldsymbol{\theta}' \sim g(\boldsymbol{\theta}'|\boldsymbol{\theta})$ 
3   Acceptance probability:  $\alpha(j+1|j) \equiv \min \left\{ 1, \frac{\pi(\boldsymbol{\theta}'|\mathbf{y})g(\boldsymbol{\theta}|\boldsymbol{\theta}')}{\pi(\boldsymbol{\theta}|\mathbf{y})g(\boldsymbol{\theta}'|\boldsymbol{\theta})} \right\}$ 
4   Draw:  $u \sim \mathcal{U}(0, 1)$ 
5   if  $u \leq \alpha(j+1|j)$ , then
6     Accept:  $\boldsymbol{\theta}^{(j+1)} = \boldsymbol{\theta}'$ ;
7     else Reject:  $\boldsymbol{\theta}_{j+1} = \boldsymbol{\theta}_j$ ;
8   end
9 end
10 Refine:  $\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{K'}\}$  to  $\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K\}$  according to integrated
    auto-correlation time  $\tau_{\text{int}}$  for large enough  $K'$ , where  $K < K'$ 
11 Draw:  $\boldsymbol{\theta}_i \in \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K\} \sim \pi(\boldsymbol{\theta}|\mathbf{y})$ 
12 Draw:  $\mathbf{x}_i|\mathbf{y}, \boldsymbol{\theta}_i$  by solving  $\mathbf{x}_i = \arg \min_{\hat{\mathbf{x}}} \|\hat{\mathbf{F}}\hat{\mathbf{x}} - (\hat{\mathbf{y}} + \boldsymbol{\eta})\|^2$  with  $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 

```



# Appendices







## Posterior of Bayesian Hierarchical model

Here we show how to obtain the posterior covariance and mean of our hierarchical Bayesian model in 1.12b - 1.12d. We do not consider the hyper-parameters and start with the joint probability distribution of  $(\mathbf{x}^T, \mathbf{y}^T)^T$ , where  $\mathbf{x} \in \mathcal{X}$  and  $\mathbf{y} \in \mathcal{Y}$  do not intersect. For more details we refer to Chapter 2 in [18] and to the book of Rue and Held [1].

The exponent of the normal Gaussian can be rewritten into:

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{Q}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{x}^T \mathbf{Q} \boldsymbol{\mu} + \text{const.} \quad (\text{A.1})$$

We like to bring the joint distribution into a similar form so that we can compare the linear and second order terms and find the precision matrix and mean of the joint distribution.

In general the joint distribution to find the expression for the posterior distribution

We can express this posterior through the likelihood and prior probability by Bayesian theorem, with a constant and positive normalization constant:

$$\pi(\mathbf{x}|\mathbf{y}) \propto \pi(\mathbf{y}|\mathbf{x})\pi(\mathbf{x}) \quad (\text{A.2})$$

Taking the logarithmic function of this formulation we can find an expression for

the the posterior covariance, with the  $\text{Var}(\mathbf{x}) = \mathbf{Q}_x^{-1}$  and  $\text{Var}(\mathbf{y}) = \mathbf{Q}_y^{-1}$ .

$$\ln \pi(\mathbf{x}|\mathbf{y}) \propto \ln \pi(\mathbf{y}|\mathbf{x}) + \ln \pi(\mathbf{x}) \quad (\text{A.3})$$

$$= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{Q}_x (\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2}(\mathbf{y} - \mathbf{A}\mathbf{x})^T \mathbf{Q}_y (\mathbf{y} - \mathbf{A}\mathbf{x}) \quad (\text{A.4})$$

$$= -\frac{1}{2} \left[ \mathbf{x}^T [\mathbf{Q}_x + \mathbf{A}^T \mathbf{Q}_y \mathbf{A}] \mathbf{x} + \mathbf{x}^T [-\mathbf{A}^T \mathbf{Q}_y] \mathbf{y} \right. \quad (\text{A.5})$$

$$\left. + \mathbf{y}^T [-\mathbf{Q}_y \mathbf{A}] \mathbf{x} + \mathbf{y}^T [\mathbf{Q}_y] \mathbf{y} - 2\mathbf{x}^T \mathbf{Q}_x \boldsymbol{\mu} \right] + \text{const.} \quad (\text{A.6})$$

Hence we deal with a Gaussian distribution, we consider second order terms only and rearrange to the precision matrix.

$$-\frac{1}{2} \begin{bmatrix} \mathbf{x}^T [\mathbf{Q}_x + \mathbf{F}^T \mathbf{Q}_y \mathbf{F}] + \mathbf{y}^T [-\mathbf{Q}_y \mathbf{F}] & \mathbf{y}^T [\mathbf{Q}_y] + \mathbf{x}^T [-\mathbf{F}^T \mathbf{Q}_y] \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \quad (\text{A.7})$$

$$= \begin{bmatrix} \mathbf{x}^T & \mathbf{y}^T \end{bmatrix} \underbrace{\begin{bmatrix} \mathbf{Q}_x + \mathbf{F}^T \mathbf{Q}_y \mathbf{F} & -\mathbf{F}^T \mathbf{Q}_y \\ -\mathbf{Q}_y \mathbf{F} & \mathbf{Q}_y \end{bmatrix}}_{\text{precision matrix}} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \quad (\text{A.8})$$

We denote the precision matrix of the joint field as:

$$\mathbf{Q}_{xy} = \begin{bmatrix} \mathbf{Q}_{aa} & \mathbf{Q}_{ab} \\ \mathbf{Q}_{ba} & \mathbf{Q}_{bb} \end{bmatrix} = \begin{bmatrix} \mathbf{Q}_x + \mathbf{F}^T \mathbf{Q}_y \mathbf{F} & -\mathbf{F}^T \mathbf{Q}_y \\ -\mathbf{Q}_y \mathbf{F} & \mathbf{Q}_y \end{bmatrix} \quad (\text{A.9})$$

The mean is defined through the linear term.

$$\frac{-2\mathbf{x}^T \mathbf{Q}_x \boldsymbol{\mu}}{-2} = \begin{bmatrix} \mathbf{x}^T & 0 \end{bmatrix} \begin{bmatrix} \mathbf{Q}_x \boldsymbol{\mu} \\ 0 \end{bmatrix} \quad (\text{A.10})$$

Comparing to the linear term of Equation A.1 we can formulate an expression for the joint mean:

$$\Rightarrow \boldsymbol{\mu}_{xy} = \mathbf{Q}_{xy}^{-1} \begin{bmatrix} \mathbf{Q}_x \boldsymbol{\mu} \\ 0 \end{bmatrix} \quad (\text{A.11})$$

The mean of the conditional distribution  $\mathbf{x}|\mathbf{y}$  is given by:

$$\boldsymbol{\mu}_{x|\mathbf{y}} = \boldsymbol{\mu}_x + \mathbf{Q}_{ba}^{-1} \mathbf{Q}_{ab} (\mathbf{x} - \boldsymbol{\mu}_y) \quad (\text{A.12})$$

$$\boldsymbol{\mu}_{x|\mathbf{y}} = \boldsymbol{\mu} + (\mathbf{Q}_x + \mathbf{F}^T \mathbf{Q}_y \mathbf{F})^{-1} \mathbf{F}^T \mathbf{Q}_y (\mathbf{x} - \mathbf{F}\boldsymbol{\mu}), \quad (\text{A.13})$$

and the covariance of  $\mathbf{x}|\mathbf{y}$  is given by:

$$\mathbf{Q}_{x|\mathbf{y}} = \mathbf{Q}_{aa} = \mathbf{Q}_x + \mathbf{F}^T \mathbf{Q}_y \mathbf{F}, \quad (\text{A.14})$$

as illustrated through Theorem 2.5 in [1].

# B

## Convergence of the Metropolis-Hastings

If we show that the detailed balance condition holds and that the state space is irreducible and aperiodic under the transition matrix  $\mathbf{P}$ , we generate a Markov chain with a unique stationary distribution proportional to  $\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ . Since the posterior is strictly positive  $\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) \geq 0$  on the finite state space  $\Omega(\mathcal{X}, \theta)$  the generated chain is irreducible. Further, it is possible to reject any proposed state and stay in the current state, which leads to aperiodicity. The detailed balance holds for the case that  $\mathbf{j} = \mathbf{i}$ , but if  $\mathbf{j} \neq \mathbf{i}$  it is not trivial. In case we accept  $\{\mathbf{x}, \boldsymbol{\theta}\}^{(n+1)} = \mathbf{j}$  as the new state we have  $\pi(\mathbf{j}|\mathbf{y})g(\mathbf{i}|\mathbf{j}) > \pi(\mathbf{i}|\mathbf{y})g(\mathbf{j}|\mathbf{i})$ . This gives us  $\alpha(\mathbf{j}|\mathbf{i}) = 1$  and  $\alpha(\mathbf{i}|\mathbf{j}) = \frac{\pi_{\mathbf{i}}g(\mathbf{j}|\mathbf{i})}{\pi_{\mathbf{j}}g(\mathbf{i}|\mathbf{j})}$  and satisfies the detailed balance:

$$\cancel{\pi_{\mathbf{j}}} \frac{\pi_{\mathbf{i}}}{\cancel{\pi_{\mathbf{j}}}} g(\mathbf{j}|\mathbf{i}) = \pi_{\mathbf{i}} g(\mathbf{j}|\mathbf{i}) \quad .$$

If  $\pi(\mathbf{j}|\mathbf{y})g(\mathbf{i}|\mathbf{j}) < \pi(\mathbf{i}|\mathbf{y})g(\mathbf{j}|\mathbf{i})$  then  $\alpha(\mathbf{i}|\mathbf{j}) = 1$  and  $\alpha(\mathbf{j}|\mathbf{i}) = \frac{\pi_{\mathbf{j}}g(\mathbf{i}|\mathbf{j})}{\pi_{\mathbf{i}}g(\mathbf{j}|\mathbf{i})}$ , this satisfies the detailed balance as well.

In conclusion the Metropolis-Hastings algorithm samples from a unique distribution proportional to the posterior distribution.





## Randomize then Optimize - RTO

$$\pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}) \propto \pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})\pi(\mathbf{x}|\boldsymbol{\theta}) \quad (\text{C.1})$$

$$\propto \exp \left[ (\mathbf{F}\mathbf{x} - \mathbf{y})^T \boldsymbol{\Sigma}^{-1} (\mathbf{F}\mathbf{x} - \mathbf{y}) + (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{Q} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (\text{C.2})$$

$$= \exp \|\hat{\mathbf{F}}\mathbf{x} - \hat{\mathbf{y}}\|^2 \quad (\text{C.3})$$

where

$$\hat{\mathbf{F}} = \begin{bmatrix} \boldsymbol{\Sigma}^{-1/2} \mathbf{F} \\ \mathbf{Q}^{1/2} \end{bmatrix}, \quad \hat{\mathbf{y}} = \begin{bmatrix} \boldsymbol{\Sigma}^{-1/2} \mathbf{y} \\ \mathbf{Q}^{1/2} \boldsymbol{\mu} \end{bmatrix} \quad (\text{C.4})$$

One sample from the posterior can be computed by minimizing the following with respect to  $\mathbf{x}$

$$\mathbf{x} = \arg \min_{\mathbf{x}} \|\hat{\mathbf{F}}\mathbf{x} - (\hat{\mathbf{y}} + \boldsymbol{\eta})\|^2, \quad \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (\text{C.5})$$

We can solve this and rewrite to

$$\frac{\partial}{\partial \mathbf{x}} \left[ (\hat{\mathbf{F}}\mathbf{x} - (\hat{\mathbf{y}} + \boldsymbol{\eta}))^T (\hat{\mathbf{F}}\mathbf{x} - (\hat{\mathbf{y}} + \boldsymbol{\eta})) \right] = 0 \quad (\text{C.6})$$

$$\Leftrightarrow \mathbf{x}^T \hat{\mathbf{F}}^T \hat{\mathbf{F}} + \hat{\mathbf{F}}^T \hat{\mathbf{F}} \mathbf{x} - \hat{\mathbf{F}}^T (\hat{\mathbf{y}} + \boldsymbol{\eta}) - (\hat{\mathbf{y}} + \boldsymbol{\eta})^T \hat{\mathbf{F}} \mathbf{x} = 0 \quad (\text{C.7})$$

We can argue through the symmetry of the inner product that and the symmetry of the precision matrix

$$\hat{\mathbf{F}}^T \hat{\mathbf{F}} \mathbf{x} = \hat{\mathbf{F}}^T (\hat{\mathbf{y}} - \boldsymbol{\eta}) \quad (\text{C.8})$$

$$\Leftrightarrow (\mathbf{F}^T \mathbf{Q}_y \mathbf{F} + \mathbf{Q}) \mathbf{x} = \mathbf{F}^T \mathbf{Q}_y \mathbf{y} + \mathbf{Q} \boldsymbol{\mu} - \hat{\mathbf{F}}^T \boldsymbol{\eta} \quad (\text{C.9})$$

If we substitute  $-\hat{\mathbf{F}}^T \boldsymbol{\eta} = \mathbf{v}_1 + \mathbf{v}_2$  we end up with

$$(\mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} + \mathbf{Q}) \mathbf{x} = \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} + \mathbf{Q} \boldsymbol{\mu} + \mathbf{v}_1 + \mathbf{v}_2 \quad (\text{C.10})$$

where  $\mathbf{v}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F})$  and  $\mathbf{v}_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$  are independent random variables.  
 mayeb introduce...  $x^2$  time nomral variubale

D

Inverting Matrices - QR factorization





# E

## Taylor expansion of $g(\lambda)$

We Taylor expand the function  $g(\lambda)$  around  $\lambda = \lambda' - \Delta\lambda$

$$g(\lambda) = \ln \det \underbrace{(\mathbf{F}^T \mathbf{F} + \lambda \mathbf{L})}_{\mathbf{B}} \quad (\text{E.1})$$

$$g(\lambda') - g(\lambda) = \ln \det(\mathbf{F}^T \mathbf{F} + \lambda' \mathbf{L}) - \ln \det(\mathbf{F}^T \mathbf{F} + \lambda \mathbf{L}) \quad (\text{E.2})$$

$$= \ln \det \left[ \frac{(\mathbf{F}^T \mathbf{F} + (\lambda + \Delta\lambda) \mathbf{L})}{(\mathbf{F}^T \mathbf{F} + \lambda \mathbf{L})} \right] \quad (\text{E.3})$$

$$= \ln \det \left[ 1 + \frac{\Delta\lambda \mathbf{L}}{\mathbf{B}} \right] \quad (\text{E.4})$$

$$= \sum_{r=1}^{\infty} \frac{(-1)^{r+1}}{r!} \text{tr}((\mathbf{B}^{-1} \mathbf{L})^r) (\Delta\lambda)^r \quad (\text{E.5})$$

, where we use the identity from [19] at page 29. So the derivatives of  $g(\lambda)$  are:

$$g^{(r)}(\lambda) = (-1)^{r+1} \text{tr}((\mathbf{B}^{-1} \mathbf{L})^r) \quad (\text{E.6})$$

$$\approx (-1)^{r+1} \sum_{k=1}^p \mathbf{z}_k^T (\mathbf{B}^{-1} \mathbf{L})^r \mathbf{z}_k \quad (\text{E.7})$$

Here we use a Monte Carlo estimate and draw  $p$  vectors  $\mathbf{z}_k \in \mathbb{R}^n$ , where each vector element  $z_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}(\{-1, 1\})$  and  $i = 1, \dots, n$ .



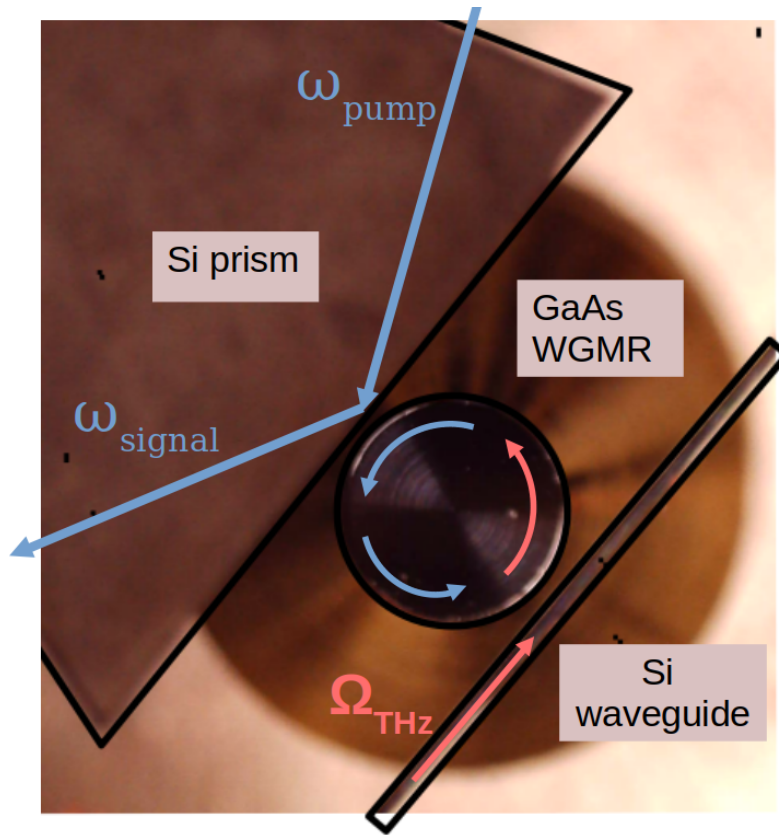
F

Radiation transfer and absorption line  
shape



G

whispering gallery resonator



**Figure G.1:** whispering gallery resonator

# References

- [1] Havard Rue and Leonhard Held. *Gaussian Markov random fields: theory and applications*. Chapman and Hall/CRC, 2005.
- [2] Dave Higdon. “A primer on space-time modeling from a Bayesian perspective”. In: *Monographs on Statistics and Applied Probability* 107 (2006), p. 217.
- [3] Pierre Brémaud. *Markov chains: Gibbs fields, Monte Carlo simulation, and queues*. Vol. 31. Springer Science & Business Media, 2013.
- [4] Julian Besag. “Spatial interaction and the statistical analysis of lattice systems”. In: *Journal of the Royal Statistical Society: Series B (Methodological)* 36.2 (1974), pp. 192–225.
- [5] Colin Fox and Richard A Norton. “Fast sampling in a linear-Gaussian inverse problem”. In: *SIAM/ASA Journal on Uncertainty Quantification* 4.1 (2016), pp. 1191–1218.
- [6] Daniel Simpson, Finn Lindgren, and Håvard Rue. “Think continuous: Markovian Gaussian models in spatial statistics”. In: *Spatial Statistics* 1 (2012), pp. 16–29.
- [7] Nicholas Metropolis and Stanislaw Ulam. “The monte carlo method”. In: *Journal of the American statistical association* 44.247 (1949), pp. 335–341.
- [8] John Michael Hammersley and David Christopher Handscomb. “General principles of the Monte Carlo method”. In: *Monte Carlo Methods*. Springer, 1964, pp. 50–75.
- [9] Paula A Whitlock and MH Kalos. *Monte Carlo Methods*. Wiley, 1986.
- [10] AA Markov. “Extension of the law of large numbers to quantities, depending on each other (1906). Reprint.” In: *Journal Électronique d’Histoire des Probabilités et de la Statistique [electronic only]* 2.1b (2006), Article–10.
- [11] Colin Fox, Geoff K Nicholls, and Sze M Tan. “An Introduction to Inverse Problems”. In: *Course notes for ELEC 404* (2010).
- [12] W Keith Hastings. “Monte Carlo sampling methods using Markov chains and their applications”. In: (1970).
- [13] Nicholas Metropolis et al. “Equation of state calculations by fast computing machines”. In: *The journal of chemical physics* 21.6 (1953), pp. 1087–1092.
- [14] Johnathan M Bardsley. “MCMC-based image reconstruction with uncertainty quantification”. In: *SIAM Journal on Scientific Computing* 34.3 (2012), A1316–A1332.
- [15] Johnathan M Bardsley et al. “Randomize-then-optimize for sampling and uncertainty quantification in electrical impedance tomography”. In: *SIAM/ASA Journal on Uncertainty Quantification* 3.1 (2015), pp. 1136–1158.

- [16] D S. Oliver, Nanqun He, and A C. Reynolds. *Conditioning permeability fields to pressure data*. 1996, cp–101.
- [17] François Orieux, Olivier Féron, and J-F Giovannelli. “Sampling high-dimensional Gaussian distributions for general linear inverse problems”. In: *IEEE Signal Processing Letters* 19.5 (2012), pp. 251–254.
- [18] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*. Vol. 4. 4. Springer, 2006.
- [19] Israel Gohberg, Seymour Goldberg, and Nahum Krupnik. *Traces and determinants of linear operators*. Vol. 116. Birkhäuser, 2012.