

# Contents

|   |            |
|---|------------|
| <b>List of Figures</b>  | <b>iii</b> |
| <b>1 Introduction</b>   | <b>3</b>   |
| 1.1 What is going on?, 3 facts, What is new in this thesis? . . . . . | 3          |
| 1.2 Thesis Outline . . . . .  | 3          |
| <b>2 Theoretical and Technical Background</b>                         | <b>5</b>   |
| 2.1 Forward Model . . . . .   | 5          |
| 2.2 Affine Map . . . . .  | 7          |
| 2.3 Bayesian Inference . . . . .                                      | 8          |
| 2.3.1 Marginal and then Conditional . . . . .                         | 10         |
| 2.4 Regularisation . . . . .  | 11         |
| 2.5 Sampling Methods . . . . .  | 13         |
| 2.5.1 Metropolis . . . . .  | 14         |
| 2.5.2 Gibbs Sampling . . . . .  | 15         |
| 2.5.3 t-walk . . . . .  | 15         |
| 2.5.4 Draw a sample from a multivariate normal distribution . . .     | 16         |
| 2.6 Numerical Approxiamtion Methods - Tensor Train . . . . .          | 16         |
| 2.6.1 Marginal Functions . . . . .                                    | 18         |
| <b>Appendices</b>   |            |
| <b>References</b>   | <b>23</b>  |



# List of Figures

|     |                                    |    |
|-----|------------------------------------|----|
| 2.1 | Schematics of Affine Map . . . . . | 8  |
| 2.2 | Bayesian Inference DAG . . . . .   | 9  |
| 2.3 | text . . . . .                     | 17 |
| 2.4 | nice matrices picture . . . . .    | 17 |



columnwidth 421.10046pt



# 1

## Introduction

### **1.1 What is going on?, 3 facts, What is new in this thesis?**

- hierachical Bayesian model, sampling to TT approx
- RTE as an example
- nonLinear to Linear Affine funciton (affine RTO)

### **1.2 Thesis Outline**

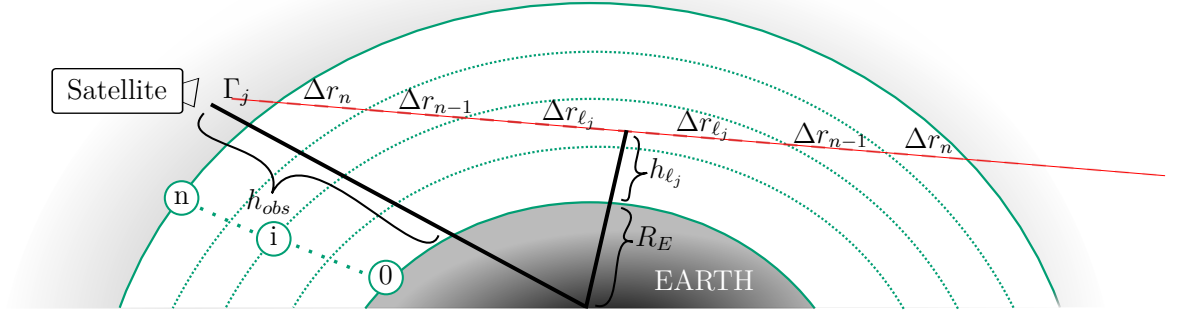




# 2

## Theoretical and Technical Background

### 2.1 Forward Model



In this section we describe the forward model which we use to simulate data and base the Bayesian inference on.

As shown in Figure ??, one measurement of a stationary satellite can be describes as the path integral along the line of sight  $\Gamma_j$  for  $j = 1, 2, \dots, m$ . For each measurement we can define a tangent height  $h_{\ell_j}$  as the shortest distance along the line of sight to the earth.

The  $j^{\text{th}}$  measurement, taken on line of sight  $\Gamma_j$  is modelled by the the radiative

transfer equation (RTE) [1]

$$y_j = \int_{\Gamma_j} B(\nu, T) k(\nu, T) \frac{\mathbf{p}(T)}{k_B \mathbf{T}(r)} \mathbf{x}(r) \tau(r) dr + \eta_j \quad (2.1)$$

$$\tau(r) = \exp \left\{ - \int_{r_{\text{obs}}}^r k(\nu, T) \frac{\mathbf{p}(T)}{k_B \mathbf{T}(r')} \mathbf{x}(r') dr' \right\} \quad (2.2)$$

where the path from the satellite along the line-of-sight of the  $j^{\text{th}}$  pointing direction is  $\Gamma_j$  and the ozone concentration  $\mathbf{x}(r)$  at distance  $r$  from the radiometer. The factor  $\tau(r) \leq 1$  accounts for re-absorption of the radiation along the line-of-sight, which makes the RTE non linear. The noise  $\eta_j$  is added to each path integral, where the noise vector  $\boldsymbol{\eta} \sim \mathcal{N}(0, \gamma^{-1} \mathbf{I})$  is normally distributed around zero with the noise precision  $\gamma$ . The absorption constant  $k(\nu, T)$  for a single gas molecule at a specific wavenumber  $\nu$  is given by the HITRAN database [2] and acts as a source function when multiplied with the black body radiation  $B(\nu, T)$ , given by Planck's law. Within the stratosphere the number density  $p(T)/(k_B T(r))$  of molecules is dependent on the pressure  $p(T)$ , the temperature  $T(r)$ , and the Boltzmann constant  $k_B$ . For fundamentals on the Radiative transfer equation we recommend 79BOOKRadiativeProcess.

We parametrize the ozone profile as a function of height, discretized into the  $n$  values in each of  $n$  layers of the discretized stratosphere where the  $i^{\text{th}}$  layer is defined by two spheres of radii  $h_{i-1} < h_i$ ,  $i = 1, \dots, n$ , with  $h_0$  and  $h_n$ . In between the heights  $h_{i-1}$  and  $h_i$ , each of the ozone concentration  $x_i$ , the pressure  $p_i$ , the temperature  $T_i$ , and thermal radiation is assumed to be constant. Above  $h_n$  and below  $h_0$ , the ozone concentration is set to zero, so no signal can be obtained. Then depending on the parameter of interest, which is either the ozone volume mixing ratio  $\mathbf{x} = \{x_1, x_2, \dots, x_n\} \in \mathbb{R}^n$  or the fraction of pressure and temperature  $\mathbf{p}/\mathbf{T} = \{p_1/T_1, p_2/T_2, \dots, p_n/T_n\} \in \mathbb{R}^n$ , we can rewrite the integral in Eq. (2.2) as e.g.  $\mathbf{A}_j(\mathbf{x}, \mathbf{p}, \mathbf{T}) \mathbf{x}$ , where the absorption  $\tau(r)$  induces non-linearity. Here, the row vector  $\mathbf{A}_j(\mathbf{x}, \mathbf{p}, \mathbf{T}) \in \mathbb{R}^n$  defines a Kernel for each measurement so that the data vector

$$\mathbf{y} = \mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T}) \mathbf{x} + \boldsymbol{\eta} = \mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T}) \frac{\mathbf{p}}{\mathbf{T}} + \boldsymbol{\eta}. \quad (2.3)$$

can be written as a matrix vector multiplication, where the matrix  $\mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T}) \in \mathbb{R}^{m \times n}$  and the noise vector  $\boldsymbol{\eta} \in \mathbb{R}^m$ .

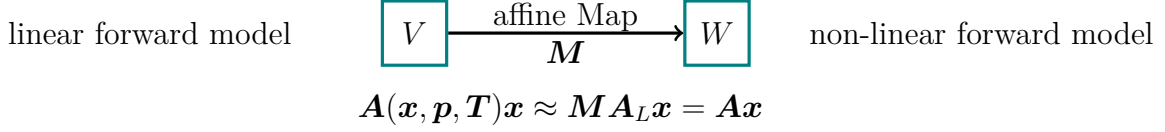
Since the absorption  $\tau(r)$  reduces measurements by of order 1%, or less, making the inverse problem only weakly non-linear. We use that to approximate the non-linear forward model  $\mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T})$  with a map  $\mathbf{M}$  so that  $\mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T}) \approx \mathbf{M}\mathbf{A}_L$ . Where each row  $\mathbf{A}_{L,j}$  of matrix as  $\mathbf{A}_L \in \mathbb{R}^{m \times n}$  is defined by the linear forward model, where absorption is neglected, e.g.  $\tau = 1$ . Then  $\mathbf{A}_{L,j}$  is either defined by  $B(\nu, T)S(\nu, T)\frac{p(T)}{k_B T(r)}dr$  or  $B(\nu, T)S(\nu, T)\frac{x}{k_B}dr$ , as in Eq.. (2.2), depending on the parameter of interest. This poses a linear inverse problem with the forward map defined by the matrix  $\mathbf{A} = \mathbf{M}\mathbf{A}_L$ , where  $\mathbf{M}$  is, more specifically, an affine map.

## 2.2 Affine Map

To approximate the non-linear forward model we use an affine map  $M : \mathbf{A}_L \mathbf{x} \rightarrow \mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T})\mathbf{x}$ , which maps the linear forward model  $\mathbf{A}_L \mathbf{x}$  onto the non-linear forward model  $\mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T})\mathbf{x}$ .

An affine map is any linear map in between two vector spaces is or affine spaces, where in affine space does not need to have a zero origin. 2.3.1. PROPOSITION AND DEFINITION On Berge book[]. In other words an affine map does not need to preserve the origin, or is a linear map on vector spaces including translation, or in the words of my supervisor, C. F., an affine map is a Taylor series of first order. For more information on affine spaces and maps we refer to [two books]

We generate two affine subspaces spaces  $V = \{\mathbf{A}(\mathbf{x}^{(1)}, \mathbf{p}, \mathbf{T}), \dots, \mathbf{A}(\mathbf{x}^{(m)}, \mathbf{p}, \mathbf{T})\}$  and  $W = \{\mathbf{A}\mathbf{x}^{(1)}, \dots, \mathbf{A}\mathbf{x}^{(m)}\}$  over the same field, with fixed  $\mathbf{p}, \mathbf{T}$ . The parameter  $\mathbf{x}$  is distributed as the so-called posterior distribution  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\} \sim \pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})$ , with hyper-parameters  $\boldsymbol{\theta}$ , according to a Bayesian hierarchical model.



**Figure 2.1:** Schematics of Affine Map, which approximates the linear forward model to the non-linear forward model.

## 2.3 Bayesian Inference

In this section we give a short introduction to Bayesian inference for a general parameter  $\mathbf{x}$  given some data  $\mathbf{y}$ , later in section ?? we set up a more sophisticated Bayesian framework applied to the forward model in section ??.

We can visualise the correlation structure of a measurement process through a hierarchially ordered directed acyclic graph (DAG), see Figure 2.2. As an observatory process naturally includes some random noise we include that in our DAG and classify the noise as a hyper-parameter in  $\boldsymbol{\theta}$ . Other hyper-parameters influence the parameters  $\mathbf{x}$  deterministically, which are then mapped through the forward model onto the space of all measurables  $\mathbf{u}$ , from which we observe some data  $\mathbf{y}$  including noise as previously mentioned. Drawing a DAG can help us to dependences within the measurement and modelling process. Given some data we infer the distribution of the underlying parameters and hyper-parameters by following the arrows in Figure ?? backwards and set up a Bayesian hierarchically ordered model.

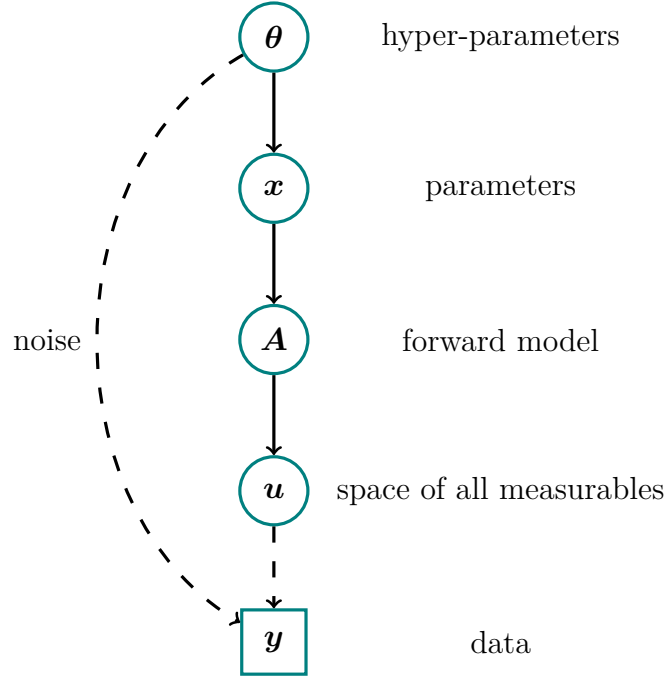
Within a linear Bayesian hierarchical model we need to define a likelihood function as well as distribution over the unknown parameters  $\mathbf{x}$  and hyper-parameters  $\boldsymbol{\theta}$ .

$$\mathbf{y}|\mathbf{x}, \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{A}\mathbf{x}, \boldsymbol{\Sigma}(\boldsymbol{\theta})) \quad (2.4a)$$

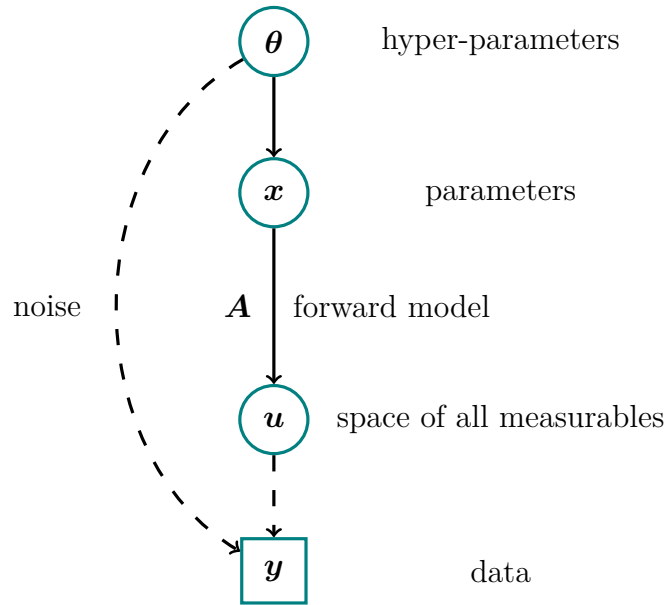
$$\mathbf{x}|\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{Q}^{-1}(\boldsymbol{\theta})) \quad (2.4b)$$

$$\boldsymbol{\theta} \sim \pi(\boldsymbol{\theta}), \quad (2.4c)$$

with the noise covariance matrix  $\boldsymbol{\Sigma}(\boldsymbol{\theta})$ , so that  $\boldsymbol{\eta} \sim \mathcal{N}(0, \boldsymbol{\Sigma}(\boldsymbol{\theta}))$  as in Eq. ??, the prior precision matrix  $\mathbf{Q}(\boldsymbol{\theta})$ , prior mean  $\boldsymbol{\mu}$  and some prior distribution over the hyper-parameters  $\pi(\boldsymbol{\theta})$ . Through sensibly choosing the prior distributions



**Figure 2.2:** The directed acyclic graph (DAG) for a typical linear inverse problem visualises forward dependencies as solid line arrows for deterministic dependencies and dotted arrows for statistical dependencies. Naturally the data  $y$  has some noise described through included in some hyper-parameters  $\theta$ . The parameters  $x$  have some dependency of those hyper-parameters  $\theta$ . The parameter  $x$  is mapped onto the space of all measurables  $u$  through the linear forward model  $A$ , so that  $Ax$  is a linear operation. From the space of all measurables we can observe some data  $y$ , statistically, where as previously mentioned some random noise is added. We set up a more sophisticated Bayesian model in chapter ?? explicitly including all hyper-parameters and parameters of interest according to the forward model in section ??.



$\pi(\mathbf{x}|\mathbf{y})$  as well as the hyper-parameters  $\boldsymbol{\theta}$  and their prior distribution  $\pi(\boldsymbol{\theta})$ , we can incorporate functional dependencies as well physical properties of the parameters. The likelihood function  $\pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$  is a measure of how the parameters and hyper-parameters fit to the data according to our forward model, including information of the measurement process.

With a normally distributed prior and likelihood function this becomes a linear-Gaussian Bayesian hierarchical model. For more detailed Bayesian analysis we recommend [].

The posterior distribution, the function of interest,

$$\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) = \frac{\pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})\pi(\mathbf{x}, \boldsymbol{\theta})}{\pi(\mathbf{y})}, \quad (2.5)$$

is given according to Bayes' theorem [], with the prior distribution  $\pi(\mathbf{x}, \boldsymbol{\theta})$  and the normalising constant  $\pi(\mathbf{y})$ . If the normalising constant is finite and non-zero we can approximate the posterior distribution

$$\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) \propto \pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})\pi(\mathbf{x}, \boldsymbol{\theta}). \quad (2.6)$$

Then the expectation of any a function  $h(\mathbf{x}, \boldsymbol{\theta})$  can be described as

$$\mathbb{E}_{\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}}[h(\mathbf{x})] = \int \int h(\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) d\mathbf{x} d\boldsymbol{\theta}, \quad (2.7)$$

which is usually a high dimensional integral and computationally not feasible to solve.

One way to work around the high dimensionality is to parameterise  $\mathbf{x}$  using hyper-parameters  $\boldsymbol{\theta}$  so that  $\mathbf{x}(\boldsymbol{\theta})$ . Another way is to separate the posterior distribution over latent field  $\mathbf{x}$  and the hyper-parameters  $\boldsymbol{\theta}$ . This is particular beneficial, when  $\mathbf{x}$  is high dimensional, e.g.  $\mathbf{x} \in \mathbb{R}^n$  with  $n = 45$  and can not be parametrised, and  $\boldsymbol{\theta}$  is low dimensional, e.g. two dimensional.

### 2.3.1 Marginal and then Conditional

The marginal and then conditional (MTC) method factorises the full posterior distribution

$$\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) = \pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})\pi(\boldsymbol{\theta}|\mathbf{y}) \quad (2.8)$$

into the marginal posterior distribution  $\pi(\boldsymbol{\theta}|\mathbf{y})$  and conditional posterior distribution  $\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})$ .

For the in Eq. ?? specified linear-Gaussian Bayesian hierarchical model the marginal posterior distribution is given as

$$\pi(\boldsymbol{\theta}|\mathbf{y}) = \int \pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) d\mathbf{x} \quad (2.9)$$

$$\propto \sqrt{\frac{\det(\boldsymbol{\Sigma}^{-1}) \det(\mathbf{Q})}{\det(\mathbf{Q} + \mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A})}} \times \exp \left[ -\frac{1}{2}(\mathbf{y} - \mathbf{A}\boldsymbol{\mu})^T \mathbf{Q}_{\boldsymbol{\theta}|\mathbf{y}}(\mathbf{y} - \mathbf{A}\boldsymbol{\mu}) \right] \pi(\boldsymbol{\theta}), \quad (2.10)$$

with

$$\mathbf{Q}_{\boldsymbol{\theta}|\mathbf{y}} = \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{A}(\mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A} + \mathbf{Q})^{-1} \mathbf{A}^T \boldsymbol{\Sigma}^{-1}. \quad (2.11)$$

See lemma [].

Then conditioned on the hyper-parameters  $\boldsymbol{\theta}$  we can draw samples of the conditional posterior distribution

$$\mathbf{x}|\mathbf{y}, \boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu} + (\mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A} + \mathbf{Q})^{-1} \mathbf{A}^T \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \mathbf{A}\boldsymbol{\mu}), (\mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A} + \mathbf{Q})^{-1}), \quad (2.12)$$

see section ?? or calculate weighted expectations of a function  $h(\mathbf{x})$

$$\mathbb{E}_{\mathbf{x}|\mathbf{y}}[h(\mathbf{x})] = \int \mathbb{E}_{\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}}[h(\mathbf{x})] \pi(\boldsymbol{\theta}|\mathbf{y}) d\boldsymbol{\theta}, \quad (2.13)$$

with weights given by  $\pi(\boldsymbol{\theta}|\mathbf{y})$ . [] Note that the noise covariance  $\boldsymbol{\Sigma}^{-1} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\theta})$  and the prior precision  $\mathbf{Q} = \mathbf{Q}(\boldsymbol{\theta})$  are depending on hyper-parameters  $\boldsymbol{\theta}$ .

In this thesis we will use sampling and deterministic methods to characterise the posterior distribution over the hyper-parameters and present the basics of those in the following sections.

## 2.4 Regularisation

Data

linear problem

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\eta} \quad (2.14)$$

$$\|\mathbf{y} - \mathbf{Ax}\|^2 \quad (2.15)$$

data misfit norm

solution semi norm

$$\lambda \|\mathbf{xLx}\|^2 \quad (2.16)$$

$$\lambda \geq 0$$

$$\arg \min \|\mathbf{y} - \mathbf{Ax}\|^2 + \lambda \|\mathbf{Tx}\|^2 \quad (2.17)$$

Derivative of solution With  $\mathbf{T}^T \mathbf{T} = \mathbf{L}$

$$\arg \min \|\mathbf{y} - \mathbf{Ax}\|^2 + \lambda \mathbf{x}^T \mathbf{Lx} \quad (2.18)$$

minizer function

L graph lapacian for example see later what we choose

Typically, L is the identity matrix or a banded matrix approximation to the (n p)th derivative.

to the 2nd derivative

The primary difficulty with linear ill-posed problems is that the inverse image is undetermined due to small (or zero) singular values of A. Actually the situation is a little worse in practice because A depends on our model of the measurement process and that is typically not precisely known, leading to a slight imprecision in the singular values. Usually that is not significant for the large singular values, but may lead to ambiguity in the small singular values so that we do not know if they are small or zero. As an introduction to regularization, which is one method for surmounting the problems associated with small singular vectors, we consider a framework for describing the quality of a reconstruction in an inverse problem.

3.1 The data misfit and the solution semi-norms In the last chapter, we considered



the linear problem  $\mathbf{d} = \mathbf{A}\mathbf{f}$  and focused on the structure of the operator  $\mathbf{A}$ . As far as the data are concerned, a reconstructed image is good provided that it gives rise to ‘mock data’  $\mathbf{A}$  which are close to the observed data. Thus, one of the quantities for measuring the quality of  $\mathbf{f}$  is the data misfit function which is usually the square of the residual norm

Since the data do not give us any information about some aspects of  $\mathbf{f}$ , it is necessary to include additional information which allows us to select from among several feasible reconstructions. Analytical solutions are available if we choose sufficiently simple criteria. One way of doing this is to introduce a second function representing our aversion to a particular reconstruction. For example, we may decide that the solution having minimum norm should be chosen from among the feasible set. This can be done by choosing

## 2.5 Sampling Methods

In this section we present the sampling based methods used in this thesis to generate an ergodic Markov-Chain  $(\mathbf{x}, \boldsymbol{\theta})^{(0)}, \dots, (\mathbf{x}, \boldsymbol{\theta})^{(k)}, \dots, (\mathbf{x}, \boldsymbol{\theta})^{(N)} \sim \pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ , where the samples  $(\mathbf{x}, \boldsymbol{\theta})^{(k)}$  are distributed as  $\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$ . Ergodicity is implied if an aperiodic and irreducible chain  $\pi$  proves to be reversible, then the chain converges and has a unique equilibrium distribution. In other words if from a state in the chain we can reach every other state in the sampling space and the previous state, and we do not get stuck in periodic loop, then it converges. Instead of proving ergodicity we can in practice look e.g. at the output samples and their trace  $\pi(\mathbf{x}^{(k)}, \boldsymbol{\theta}^{(k)}|\mathbf{y})$  to see convergence.

Ergodicity is important so that we can from the sample based estimate

Then for large enough  $N$  the samples based estimate of Eq. ?? and of any function  $h(\mathbf{x}, \boldsymbol{\theta})$  is

$$\mathbb{E}_{\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}}[h(\mathbf{x}, \boldsymbol{\theta})] \approx \frac{1}{N} \sum_{k=1}^N h(\mathbf{x}^{(k)}, \boldsymbol{\theta}^{(k)}). \quad (2.19)$$

In practise of this thesis we use Markov-chain Monte-Carlo (MCMC) methods on target distributions such as  $\pi(\boldsymbol{\theta}|\mathbf{y})$  and so we will illustrate the sampling procedures for the following three subsections on that distribution.

### 2.5.1 Metropolis

The Metropolis algorithm is special case of the Metropolis-Hastings algorithm, with a symmetric proposal distribution  $q(i|j) = q(j|i)$  [].

The Metropolis-Hastings algorithm starts with a initial sample  $\boldsymbol{\theta}^{(t)}$  at  $t = 0$ . We propose a new sample  $\boldsymbol{\theta} \sim q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t-1)})$  according to the proposal distribution  $q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t-1)})$  given the previous state. Then accept and set  $\boldsymbol{\theta}^{(t)} = \boldsymbol{\theta}$  with

$$\alpha(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t-1)}) = \min \left\{ 1, \frac{\pi(\boldsymbol{\theta}|\mathbf{y})q(\boldsymbol{\theta}^{(t-1)}|\boldsymbol{\theta})}{\pi(\boldsymbol{\theta}^{(t-1)}|\mathbf{y})q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t-1)})} \right\} \quad (2.20)$$

or reject and keep  $\boldsymbol{\theta}^{(t)} = \boldsymbol{\theta}^{(t-1)}$ , which we do by comparing  $\alpha$  to a uniform random number  $u \sim \mathcal{U}(0, 1)$ . Note that symmetrical proposal distribution  $q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t-1)}) = q(\boldsymbol{\theta}^{(t-1)}|\boldsymbol{\theta})$  cancel, then we call that a Metropolis-algorithm []. To assure convergence independent of the intial sample we discard samples after the so-called burn-in period and effectivally generate a Markov-Chain of length  $N - N_{\text{burn-in}}$ .

#### Algorithm 1: Metropolis

- 1: Initialize  $\boldsymbol{\theta}^{(0)}$
- 2: **for**  $k = 1, \dots, N$  **do**
- 3:   Propose  $\boldsymbol{\theta} \sim q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t-1)}) = q(\boldsymbol{\theta}^{(t-1)}|\boldsymbol{\theta})$
- 4:   Compute
 
$$\alpha(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t-1)}) = \min \left\{ 1, \frac{\pi(\boldsymbol{\theta}|\mathbf{y})q(\boldsymbol{\theta}^{(t-1)}|\boldsymbol{\theta})}{\pi(\boldsymbol{\theta}^{(t-1)}|\mathbf{y})q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t-1)})} \right\}$$
- 5:   Draw  $u \sim \mathcal{U}(0, 1)$
- 6:   **if**  $\alpha \geq u$  **then**
- 7:     Accept and set  $\boldsymbol{\theta}^{(t)} = \boldsymbol{\theta}$
- 8:   **else**
- 9:     Reject and keep  $\boldsymbol{\theta}^{(t)} = \boldsymbol{\theta}^{(t-1)}$
- 10:   **end if**
- 11: **end for**
- 12: Output:  $\boldsymbol{\theta}^{(0)}, \dots, \boldsymbol{\theta}^{(k)}, \dots, \boldsymbol{\theta}^{(N)} \sim \pi(\boldsymbol{\theta}|\mathbf{y})$

Since the proposal distribution is symmetric, Metropolis chains are reversible and as you can see in  $\alpha(\boldsymbol{\theta}|\boldsymbol{\theta}^{(t-1)})$  converge towards  $\pi(\boldsymbol{\theta}|\mathbf{y})$ . A more mathematical prove of ergodicity for the general Metropolis-Hastings algorithm can be found in [1].

### 2.5.2 Gibbs Sampling

Gibbs sampling is a special case of the metropolis hastings algorithm with the acceptance probability of 1. In Gibbs sampling according to [1] we move one dimension by fixing all other dimensions.

Assume  $\boldsymbol{\theta} \in \mathbb{R}^d$  is d-dimensional and  $\{\boldsymbol{\theta}_{<j}, \boldsymbol{\theta}_{>j}\} = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{j-1}, \boldsymbol{\theta}_{j+1}, \dots, \boldsymbol{\theta}_d\}$  denotes all dimensions except  $\boldsymbol{\theta}_j$ . Then to draw a new sample  $\boldsymbol{\theta}^{(t)}$ , we iterate through each dimension to draw a sample  $\boldsymbol{\theta}_j^{(t)} \sim \pi(\boldsymbol{\theta}_j|\boldsymbol{\theta}_{<j}, \boldsymbol{\theta}_{>j}, \mathbf{y})$  from the conditional distribution.

#### Algorithm 2: Gibbs

```

1: Initialize  $\boldsymbol{\theta}^{(0)} = \{\boldsymbol{\theta}_1^{(0)}, \dots, \boldsymbol{\theta}_j^{(0)}, \dots, \boldsymbol{\theta}_d^{(0)}\}$ .
2: for  $k = 1, \dots, N$  do
3:   for  $j = 1, \dots, d$  do
4:     Draw  $\boldsymbol{\theta}_j^{(t)} \sim \pi(\boldsymbol{\theta}_j|\boldsymbol{\theta}_{<j}, \boldsymbol{\theta}_{>j}, \mathbf{y})$ 
5:   end for
6: end for
7: Output:  $\boldsymbol{\theta}^{(0)}, \dots, \boldsymbol{\theta}^{(k)}, \dots, \boldsymbol{\theta}^{(N)} \sim \pi(\boldsymbol{\theta}|\mathbf{y})$ 
```

If the target distribution is well behaved  $\pi(\boldsymbol{\theta}_j|\boldsymbol{\theta}_{<j}, \boldsymbol{\theta}_{>j}, \mathbf{y}) > 0$  in the sampling space the chain produced by the Gibbs sampling algorithm is aperiodic and irreducible. We can also show that the criterion, the detailed balance condition, for reversibility is met [1]. Then a Gibbs sampler produces an ergodic Markov chain.

### 2.5.3 t-walk

We use the t-walk sampler developed by Christens and Fox as a black box sampler [1]. Within the t-walk there are three different moves in the sampling space available and by construction of the t-walk is ergodic, see section [1].

### 2.5.4 Draw a sample from a multivariate normal distribution

for Linear Gaussian Bayesian hierarchical model We need to draw a sample from the multivariate normal distribution. We can do this using the radomize then optimise (RTO) method [1], with a perturb exponential.

In this thesis we use the RTO method to draw a sample from the full conditional distribution  $\pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta})$ , see section ??, which is a normal distribution we can rewrite to:

$$\pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}) \propto \pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})\pi(\mathbf{x}|\boldsymbol{\theta}) \quad (2.21)$$

$$= \exp\|\hat{\mathbf{A}}\mathbf{x} - \hat{\mathbf{y}}\|^2, \quad (2.22)$$

where

$$\hat{\mathbf{A}} = \begin{bmatrix} \boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\theta})\mathbf{A} \\ \mathbf{Q}^{1/2}(\boldsymbol{\theta}) \end{bmatrix}, \quad \hat{\mathbf{y}} = \begin{bmatrix} \boldsymbol{\Sigma}^{-1/2}(\boldsymbol{\theta})\mathbf{y} \\ \mathbf{Q}^{1/2}(\boldsymbol{\theta})\boldsymbol{\mu} \end{bmatrix}. \quad (2.23)$$

Then one sample can be computed by minimising the following equation with respect to  $\hat{\mathbf{x}}$  :

$$\mathbf{x}_i = \arg \min_{\hat{\mathbf{x}}} \|\hat{\mathbf{A}}\hat{\mathbf{x}} - (\hat{\mathbf{y}} + \boldsymbol{\eta})\|^2, \quad \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (2.24)$$

where we add a randomised perturbation  $\boldsymbol{\eta}$ . Next, we substitute  $-\hat{\mathbf{A}}^T\boldsymbol{\eta} = \mathbf{v}_1 + \mathbf{v}_2$  we can rewrite the argument of Eq. 2.23 to

$$(\mathbf{A}^T\boldsymbol{\Sigma}^{-1}(\boldsymbol{\theta})\mathbf{A} + \mathbf{Q}(\boldsymbol{\theta}))\mathbf{x}_i = \mathbf{A}^T\boldsymbol{\Sigma}^{-1}(\boldsymbol{\theta})\mathbf{y} + \mathbf{Q}(\boldsymbol{\theta})\boldsymbol{\mu} + \mathbf{v}_1 + \mathbf{v}_2, \quad (2.25)$$

where  $\mathbf{v}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^T\boldsymbol{\Sigma}^{-1}(\boldsymbol{\theta})\mathbf{A})$  and  $\mathbf{v}_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}(\boldsymbol{\theta}))$  are independent random variables [1].

## 2.6 Numerical Approxiamtion Methods - Tensor Train

Using the tensor train format to approximate a  $d$ -dimensional function  $\pi(x)$  enables us to compute marginal posterior probability distribution cheaply. As the name

suggest the tensor train format is a train of tensors, more specifically two and three dimensional tensors which we call cores  $\pi_k \in \mathbb{R}^{r_{k-1} \times n \times r_k}$  and which are connected through a ranks  $r_k$  and  $r_{k-1}$  for the  $k$ th dimension and defined by the number of gridpoints  $n$ , as in Figure ?? displayed. For the first and last dimensional core the outer ranks are  $r_0 = r_d = 1$ .

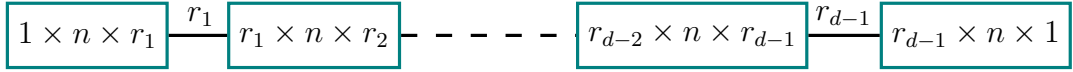


Figure 2.3: text

Figure 2.4: nice matrices picture

The approximated marginal target function

$$f_{X_k}(x_k) = \frac{1}{z} \left| \left( \int_{\mathbb{R}} \pi_1(x_1) \lambda_1(x_1) dx_1 \right) \cdots \left( \int_{\mathbb{R}} \pi_{k-1}(x_{k-1}) \lambda_{k-1}(x_{k-1}) dx_{k-1} \right) \right. \\ \left. \pi_k(x_k) \lambda_k(x_k) \left( \int_{\mathbb{R}} \pi_{k+1}(x_{k+1}) \lambda_{k+1}(x_{k+1}) dx_{k+1} \right) \cdots \left( \int_{\mathbb{R}} \pi_d(x_d) \lambda_d(x_d) dx_d \right) \right|, \quad (2.26)$$

is given by integration over each core, where  $k$ th is in  $\mathbb{R}^{r_{k-1} \times r_k}$  and  $z$  is some normalising constant. Here we introduce some Lebesgue measurable weight function  $\lambda(x) = \prod_{i=1}^d \lambda_i(x_i)$ . Why? [].

From here the notation and procedure is taken mostly from []. For numerical stability we can approximate the square root of

$$\sqrt{\pi(x)} \approx \tilde{g}(x) = \mathbf{G}_1(x_1), \dots, \mathbf{G}_k(x_k), \dots, \mathbf{G}_d(x_d) \quad (2.27)$$

where the TT-core

$$G_k^{(\alpha_{k-1}, \alpha_k)}(x_k) = \sum_{i=1}^{n_k} \phi_k^{(i)}(x_k) \mathbf{A}_k[\alpha_{k-1}, i, \alpha_k], \quad k = 1, \dots, d, \quad (2.28)$$

with the associated  $k$ th coefficient tensor  $\mathbf{A}_k \in \mathbb{R}^{r_{k-1} \times n_k \times r_k}$  and the  $k$ -th basis functions  $\phi_k^{(i)}(x_k)$ .

We assume the function

$$\pi(x) \approx \gamma' + g^2(x), \quad (2.29)$$

where  $g(x)$  is defined through the tensor train decomposition plus an error  $\gamma'$  according to the l2 norm. Then the normalised target function is

$$f_X(x) = \frac{1}{z} \pi(x) \lambda(x) = \frac{1}{z} (\gamma' \lambda(x) + g^2(x) \lambda(x)) \quad (2.30)$$

with a normalisation constant  $z$ . Consequently the approximated marginal functions can be expressed as

$$\begin{aligned} f_{X_k}(x_k) = \frac{1}{z} & \left( \gamma' \prod_{i=1}^{k-1} \lambda_i(\mathcal{X}_i) \prod_{i=k+1}^d \lambda_i(\mathcal{X}_i) \right. \\ & + \left( \int_{\mathbb{R}} \mathbf{G}_1^2(x_1) \lambda_1(x_1) dx_1 \right) \cdots \left( \int_{\mathbb{R}} \mathbf{G}_{k-1}^2(x_{k-1}) \lambda_{k-1}(x_{k-1}) dx_{k-1} \right) \\ & \mathbf{G}_k^2(x_k) \lambda_k(x_k) \\ & \left. \left( \int_{\mathbb{R}} \mathbf{G}_{k+1}^2(x_{k+1}) \lambda_{k+1}(x_{k+1}) dx_{k+1} \right) \cdots \left( \int_{\mathbb{R}} \mathbf{G}_d^2(x_d) \lambda_d(x_d) dx_d \right) \right), \end{aligned} \quad (2.31)$$

where  $\lambda_k(\mathcal{X}_k) = \int_{\mathcal{X}_k} \lambda_k(x_k) dx_k$ .

To effeciently calculate these marginals on can us a procedure which is called left and right orthogonalization of cores [] [32]. To do so we define the mass matrix  $\mathbf{M}_k \in \mathbb{R}^{n_k \times n_k}$  by

$$\mathbf{M}_k[i, j] = \int_{X_k} \phi_k^{(i)}(x_k) \phi_k^{(j)}(x_k) \lambda(x_k) dx_k, \quad i = 1, \dots, n_k, \quad j = 1, \dots, n_k, \quad (2.32)$$

where  $\{\phi_k^{(i)}(x_k)\}_{i=1}^{n_k}$  is the set of basis functions for the  $k$ -th coordinate.

### 2.6.1 Marginal Functions

We calculate the marginal functions through procedures, which we call backward marginalisation [] and forward marginalisation. We gain the coefficient matrices  $\mathbf{B}_k$  through backward marginalisation and the coefficient matrices  $\mathbf{B}_{pre,n}$  through forward marginalisation, which enables us to calculate marginal fuunction similar to [].

The proposition 1 to caculte  $\mathbf{B}_k$  is taken from [].

**Proposition 1** (Backward Marginalisation): Starting with the last coordinate  $k = d$ , we set  $\mathbf{B}_d = \mathbf{A}_d$ . The following procedure can be used to obtain the coefficient tensor  $\mathbf{B}_{k-1} \in \mathbb{R}^{r_{k-2} \times n_{k-1} \times r_{k-1}}$ , which we need for defining the marginal function  $f_{X_k}(x_k)$ :

1. Use the Cholesky decomposition of the mass matrix,  $\mathbf{L}_k \mathbf{L}_k^\top = \mathbf{M}_k \in \mathbb{R}^{n_k \times n_k}$ , to construct a tensor  $\mathbf{C}_k \in \mathbb{R}^{r_{k-1} \times n_k \times r_k}$ :

$$\mathbf{C}_k[\alpha_{k-1}, \tau, l_k] = \sum_{i=1}^{n_k} \mathbf{B}_k[\alpha_{k-1}, i, l_k] \mathbf{L}_k[i, \tau]. \quad (2.33)$$

2. Unfold  $\mathbf{C}_k$  along the first coordinate and compute the thin QR decomposition, so that  $\mathbf{C}_k^{(R)} \in \mathbb{R}^{r_{k-1} \times (n_k r_k)}$ :

$$\mathbf{Q}_k \mathbf{R}_k = (\mathbf{C}_k^{(R)})^\top. \quad (2.34)$$

3. Compute the new coefficient tensor:

$$\mathbf{B}_{k-1}[\alpha_{k-2}, i, l_{k-1}] = \sum_{\alpha_{k-1}=1}^{r_{k-1}} \mathbf{A}_{k-1}[\alpha_{k-2}, i, \alpha_{k-1}] \mathbf{R}_k[l_{k-1}, \alpha_{k-1}]. \quad (2.35)$$

Then we need to do this the other way as well.

In addiotn we also have to do forward marginalistion starig with the first dimen-

sion

**Proposition 2** (Forward Marginalistaion): Starting with the first coordinate  $k = 1$ , we set  $\mathbf{B}_{pre,1} = \mathbf{A}_1$ . The following procedure can be used to obtain the coefficient tensor  $\mathbf{B}_{pre,k+1} \in \mathbb{R}^{r_k \times n_{k+1} \times r_{k+1}}$  for defining the marginal function  $f_{X_k}(x_k)$ :

1. Use the Cholesky decomposition of the mass matrix,  $\mathbf{L}_k \mathbf{L}_k^\top = \mathbf{M}_k \in \mathbb{R}^{n_k \times n_k}$ , to construct a tensor  $\mathbf{C}_k \in \mathbb{R}^{r_{k-1} \times n_k \times r_k}$ :

$$\mathbf{C}_{pre,k}[\alpha_{k-1}, \tau, l_k] = \sum_{i=1}^{n_k} \mathbf{L}_k[i, \tau] \mathbf{B}_{pre,k}[\alpha_{k-1}, i, l_k]. \quad (2.36)$$

2. Unfold  $\mathbf{C}_{pre,k}$  along the first coordinate and compute the thin QR decomposition, so that  $\mathbf{C}_{pre,k}^{(R)} \in \mathbb{R}^{(r_{k-1} n_k) \times r_k}$ :

$$\mathbf{Q}_{pre,k} \mathbf{R}_{pre,k} = (\mathbf{C}_{pre,k}^{(R)}). \quad (2.37)$$

3. Compute the new coefficient tensor  $\mathbf{B}_{pre,k+1} \in \mathbb{R}^{r_{k-1} \times n_k \times r_k}$ :

$$\mathbf{B}_{pre,k+1}[l_{k+1}, i, \alpha_{k+1}] = \sum_{\alpha_k=1}^{r_k} \mathbf{R}_{pre,k}[l_{k+1}, \alpha_k] \mathbf{A}_{k+1}[\alpha_k, i, \alpha_{k+1}]. \quad (2.38)$$

The marginal PDF of  $X_k$  can be expressed as

$$f_{X_k}(x_k) = \frac{1}{z} \left( \gamma' \prod_{i=1}^{k-1} \lambda_i(X_i) \prod_{i=k+1}^d \lambda_i(X_i) + \sum_{l_{k-1}=1}^{r_{k-1}} \sum_{l_k=1}^{r_k} \left( \sum_{i=1}^{n_k} \phi_k^{(i)}(x_k) \mathbf{D}_k[l_{k-1}, i, l_k] \right)^2 \right) \lambda_k(x_k), \quad (2.39)$$

where  $\mathbf{D}_k \in \mathbb{R}^{r_{k-1} \times n_k \times r_k}$  and  $\mathbf{R}_{pre,k-1} \in \mathbb{R}^{r_{k-1} \times r_{k-1}}$  and  $\mathbf{B}_k \in \mathbb{R}^{r_{k-1} \times n_k \times r_k}$

$$\mathbf{D}_k[l_{k-1}, i, l_k] = \sum_{\alpha_{k-1}=1}^{r_{k-1}} \mathbf{R}_{pre,k-1}[l_{k-1}, \alpha_{k-1}] \mathbf{B}_k[\alpha_{k-1}, i, l_k]. \quad (2.40)$$

Special Cases The marginal PDF of  $X_1$  can be expressed as

$$f_{X_1}(x_1) = \frac{1}{z} \left( \gamma' \prod_{i=2}^d \lambda_i(\mathcal{X}_i) + \sum_{l_1=1}^{r_1} \left( \sum_{i=1}^{n_1} \phi_1^{(i)}(x_1) \mathbf{D}_1[i, l_1] \right)^2 \right) \lambda_1(x_1), \quad (2.41)$$

where  $\mathbf{D}_1[i, l_1] = \mathbf{B}_1[\alpha_0, i, l_1]$  and  $\alpha_0 = 1$ .

The marginal PDF of  $X_n$  can be expressed as

$$f_{X_n}(x_n) = \frac{1}{z} \left( \gamma' \prod_{i=1}^{d-1} \lambda_i(\mathcal{X}_i) + \sum_{l_{n-1}=1}^{r_{n-1}} \left( \sum_{i=1}^{n_1} \phi_1^{(i)}(x_1) \mathbf{D}_n[l_{n-1}, i] \right)^2 \right) \lambda_n(x_n), \quad (2.42)$$

where  $\mathbf{D}_n[l_{n-1}, i] = \mathbf{B}_{pre,n}[l_{n-1}, i, \alpha_n]$  and  $\alpha_n = 1$ .



# Appendices



## References

- [1] *Handbook for the Montreal protocol on substances that deplete the ozone layer*. Nairobi: The Secretariat of The Vienna Convention for the Protection of the Ozone Layer and The Montreal Protocol on Substances that Deplete the Ozone Layer, United Nations Environment Programme, 2006.
- [2] Iouli E Gordon et al. “The HITRAN2020 molecular spectroscopic database”. In: *Journal of Quantitative Spectroscopy and Radiative Transfer* 277 (2022), p. 107949.