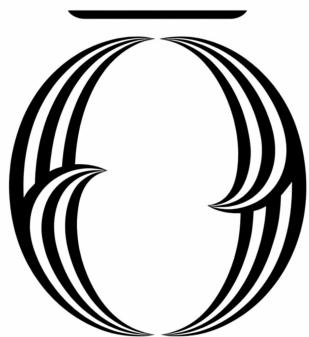


Suitably impressive thesis title



University  
of Otago

ŌTĀKOU WHAKAIHU WAKA

---

NEW ZEALAND

Lennart Golks  
Department of Physics

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

October 2025



# Acknowledgements

## **Personal**

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

## **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.



## Abstract

Your abstract text goes here. Check your departmental regulations, but generally this should be less than 300 words. See the beginning of Chapter ?? for more.



# Contents

<b>List of Figures</b>	<b>ix</b>
<b>List of Abbreviations</b>	<b>xi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Research Gap and Contribution . . . . .	2
1.3 Thesis Structure . . . . .	3
<b>2 Theoretical and Technical Background</b>	<b>5</b>
2.1 Hierarchical Bayesian Inference . . . . .	5
2.2 Sampling Methods . . . . .	9
2.2.1 Metropolis within Gibbs . . . . .	10
2.2.2 T-walkSampler as a Black Box . . . . .	10
2.2.3 Randomise then Optimise . . . . .	11
2.3 Numerical Approximation – Tensor-Train (TT) . . . . .	12
2.3.1 Marginal Functions . . . . .	14
2.3.2 Sampling from a TT Approximation . . . . .	16
2.3.3 On the Error of a TT Approximation . . . . .	18
2.4 Affine Map . . . . .	19
2.5 Regularisation . . . . .	20
<b>3 Forward Model</b>	<b>23</b>
3.1 Singular Value Decomposition of the Forward Model . . . . .	25
<b>4 Results and Conclusions</b>	<b>31</b>
4.1 Simulate Data based on a Ground Truth . . . . .	32
4.2 Hierarchical Bayesian Framework for Ozone . . . . .	35
4.2.1 Prior Modelling . . . . .	36
4.2.2 Posterior Distribution – linear Model . . . . .	38
4.3 Approximate non-linear Forward Model with an Affine Map . . . . .	46
4.4 Regularisation Solution vs. Bayesian Approach – approximated Model . .	49
4.4.1 Posterior Distribution for Ozone . . . . .	49
4.4.2 Solution by Regularisation . . . . .	53

4.5	Hierarchical Bayesian Framework for Ozone, Pressure and Temperature . . . . .	55
4.5.1	Prior Modelling . . . . .	57
4.5.2	Posterior Distribution . . . . .	61
<b>5</b>	<b>Summary and Outlook</b>	<b>77</b>
5.1	Regularisation Solution vs. Hierarchical Bayesian Approach . . . . .	77
5.2	Sampling Methods vs. TT Approximation . . . . .	77
5.3	Atmospheric Physics . . . . .	79
5.3.1	Measurement Device . . . . .	79
5.3.2	Model . . . . .	79
<b>References</b>		<b>81</b>
<b>Appendices</b>		
<b>A</b>	<b>Theoretical and technical background</b>	<b>87</b>
A.1	Correlation Structure . . . . .	87
A.2	On the Monte-Carlo Error and Integrated Autocorrelation time . . . . .	88
A.3	Measure theory . . . . .	89
A.3.1	Probability Measure . . . . .	90
A.3.2	$\sigma$ -Algebra . . . . .	90
A.4	Python Code . . . . .	92
<b>B</b>	<b>Additional Figures</b>	<b>95</b>
B.1	Ozone . . . . .	95
B.1.1	Ozone Prior . . . . .	95
B.1.2	Integrated Autocorrelation plots . . . . .	95
B.2	Pressure and Temperature . . . . .	100
B.2.1	Priors . . . . .	100
B.2.2	T-walk Trace . . . . .	100
B.2.3	Integrated Autocorrelation Plots . . . . .	100

# List of Figures

2.1	Hierarchical Bayesian Inference . . . . .	6
2.2	Visualisation of a tensor train . . . . .	13
2.3	Schematics of the affine map . . . . .	20
3.1	Schematic of measurement and analysis geometry. . . . .	23
3.2	Tangent heights for different sequences of measurements. . . . .	26
3.3	Singular values of linear forward model matrix for different sequences of measurements. . . . .	27
3.4	First 10 right singular vectors of forward model. . . . .	28
3.5	Right singular vectors 11 to 19 of forward model. . . . .	29
3.6	Last 5 right singular vectors of forward model. . . . .	30
4.1	Logarithmic plot of data points at different tangent height. . . . .	34
4.2	Directed acyclic graph for ozone retrieval and MTC scheme. . . . .	35
4.3	Plot of the functions $f(\lambda)$ and $g(\lambda)$ for marginal posterior. . . . .	39
4.4	IATC of $\lambda$ samples from $\pi(\gamma, \lambda   \mathbf{y})$ , for linear model. . . . .	41
4.5	Scatter plot of samples from marginal posterior, including weighting from TT approximation; trace plot of the marginal posterior samples. . . . .	42
4.6	Ozone samples of the full posterior. . . . .	44
4.7	Assessment of Monte-Carlo error. . . . .	45
4.8	Strategy to find affine map. . . . .	46
4.9	Assessment of affine map. . . . .	48
4.10	Marginal posterior histograms and TT approximation as well as hyper-prior distribution. . . . .	50
4.11	Ozone posterior mean and variance and the regularised solution compared to the ground truth. . . . .	51
4.12	Singular values of the posterior covariance matrix . . . . .	52
4.13	Plot of the L-curve to find the regularised solution. . . . .	54
4.14	Directed acyclic graph of Bayesian model for ozone, pressure and temperature. .	55
4.15	Prior Samples of $\mathbf{T}$ according to the respective hyper-prior distribution. .	58
4.16	Prior Samples of $\mathbf{p}$ according to the respective hyper-prior distribution. .	59
4.17	Prior Samples of $\mathbf{p}/\mathbf{T}$ according to the respective hyper-prior distribution. .	60
4.18	Correlation plot of samples from TT-approximation . . . . .	63

4.19 Optimal rank and number of grid points for TT approximation . . . . .	66
4.20 Histograms and TT approximation of posterior distribution as well as hyper-prior distribution . . . . .	67
4.21 Histograms and TT approximation of posterior distribution as well as hyper-prior distribution . . . . .	68
4.22 Histograms and TT approximation of posterior distribution as well as hyper-prior distribution . . . . .	69
4.23 Histograms and TT approximation of posterior distribution as well as hyper-prior distribution . . . . .	70
4.24 Histograms and TT approximation of posterior distribution as well as hyper-prior distribution . . . . .	71
4.25 Histograms and TT approximation of posterior distribution as well as hyper-prior distribution . . . . .	72
4.26 Temperature posterior samples . . . . .	73
4.27 Pressure posterior samples . . . . .	75
4.28 Pressure posterior samples . . . . .	76
 A.1 Correlation structure in between parameters and hyper-parameters . . . . .	88
 B.1 Samples from ozone prior distribution . . . . .	96
B.2 IACT of $\gamma$ samples from $\pi(\gamma, \lambda   \mathbf{y})$ , for linear model . . . . .	97
B.3 IACT and autocorrelation function for samples $\lambda \sim \pi(\cdot   \gamma, \mathbf{y})$ . . . . .	98
B.4 IACT and autocorrelation function for samples $\gamma \sim \pi(\cdot   \lambda, \mathbf{y})$ . . . . .	99
B.5 Prior distributions $\pi(\mathbf{h}_T)$ . . . . .	100
B.6 Prior samples of $1/T$ . . . . .	101
B.7 T-walk trace . . . . .	101
B.8 IACT and autocorrelation function for $h_{T,1}$ samples . . . . .	102
B.9 IACT and autocorrelation function for $h_2$ samples . . . . .	103
B.10 IACT and autocorrelation function for $h_3$ samples . . . . .	104
B.11 IACT and autocorrelation function for $h_4$ samples . . . . .	105
B.12 IACT and autocorrelation function for $h_5$ samples . . . . .	106
B.13 IACT and autocorrelation function for $h_6$ samples . . . . .	107
B.14 IACT and autocorrelation function for $a_0$ samples . . . . .	108
B.15 IACT and autocorrelation function for $a_1$ samples . . . . .	109
B.16 IACT and autocorrelation function for $a_2$ samples . . . . .	110
B.17 IACT and autocorrelation function for $a_3$ samples . . . . .	111
B.18 IACT and autocorrelation function for $a_4$ samples . . . . .	112
B.19 IACT and autocorrelation function for $a_5$ samples . . . . .	113
B.20 IACT and autocorrelation function for $a_6$ samples . . . . .	114
B.21 IACT and autocorrelation function for $T_0$ samples . . . . .	115
B.22 IACT and autocorrelation function for $b$ samples . . . . .	116
B.23 IACT and autocorrelation function for $p_0$ samples . . . . .	117

## List of Abbreviations

<b>DAG</b>	Directed Acyclic Graph
<b>HITRAN</b>	High Resolution Transmission
<b>IATC</b>	Integrated Autocorrelation Time
<b>IRT</b>	Inverse Rosenblatt Transform
<b>L</b>	Linear
<b>MCMC</b>	Markov Chain Monte-Carlo
<b>MH</b>	Metropolis-Hastings
<b>MLS</b>	Microwave Limb Sounder
<b>MTC</b>	Marginal and Then Conditional
<b>MVN</b>	Multivariate Normal
<b>MWG</b>	Metropolis Wihtin Gibbs
<b>NASA</b>	National Aeronautics and Space Administration
<b>NL</b>	Non-Linear
<b>RMS</b>	Root Mean Square
<b>RTE</b>	Radiative Transfer Equation
<b>RTO</b>	Randomise Then Optimise
<b>SIRT</b>	Squared Inverse Rosenblatt Transform
<b>STD</b>	Standard Deviation
<b>SVD</b>	Singular Value Decomposition
<b>TT</b>	Tensor-Train
<b>VMR</b>	Volume Mixing Ratio



# 1

## Introduction

Here, we briefly describe the currently used standard to retrieve atmospheric trace gas concentrations and what motivates us to employ a hierarchical Bayesian framework on an atmospheric limb sounder measuring ozone, where we contribute to existing methods and how we improve those. Lastly, we provide the reader with the thesis structure.

### 1.1 Motivation

Since the only currently operating ozone limb sounder, the Microwave Limb Sounder (MLS) on NASA’s Aura satellite, is gradually drifting from its orbit and scheduled to be phased out by 2026 [1], a group led by Harald Schwefel has proposed an alternative approach to fill this observational gap using a much smaller platform, such as a disk-shaped resonator mounted on a 6U CubeSat [2]. The proposed system targets a narrow frequency band and converts the thermal radiation emitted by ozone from the terahertz region to the optical domain [3, 4].

This conversion offers a cost-effective and energy-efficient solution, as it circumvents the need for large, energy-hungry cooling devices that are traditionally required to capture terahertz signals. Instead, signal acquisition in the optical domain can be implemented by using compact, cheap, and low-power photonic technologies.

Currently, the inverse problem to retrieve any trace gas from limb-sounding data is approached by the atmospheric physics community using optimisation and regularisation techniques developed in the 1970s [5, 6]. These approaches are based on a “best fit to data but not the best fit to parameters”[7]. Instead, we employ a hierarchically structured Bayesian framework to infer ozone concentrations, where we find a distribution of parameters given some data. This probabilistic approach provides estimates of parameters and their true uncertainties.

## 1.2 Research Gap and Contribution

As already mentioned, currently the MLS retrieval algorithm [8] is based on the “optimal estimation” method from [5]. This approach iteratively minimises a squared residual norm by fitting parameters to a set of data and penalises against a chosen regularisation. This does not provide comprehensive information about the parameters, the underlying correlation structures, and can lead to unphysical results, e.g. negative ozone concentration values [9]. The errors provided are based on a local derivative of the forward map around one optimal solution, which is inherently highly sensitive to its location. Additionally, these retrievals are conditioned on external estimates of other parameters, such as temperature or pressure [8]. This results in biased solutions, where the bias is then removed based on empirical decisions [10, 11]. Current machine learning efforts condition on one single noise value in the developed model, which is trained for about one month, and do not include noise as a retrieval parameter. Additionally they compare to a “ground truth” provided by the previously described optimal estimation approach [12, 13].

We address these limitations by including measurement noise as a hyper-parameter as well as the smoothness of the ozone profile explicitly in the modelling and inversion process. Then we are able to provide a range of ozone profiles, which are all fitting to the data. Naturally, noise (hyper-parameter) is a random process and according to that noise we deal with distributions over hyper-parameters and parameters (e.g. ozone concentrations) and can provide errors according to those distributions, instead of one “optimal” solution. This approach is called *hierarchical* Bayesian modelling. Livesey et al. [8] report “unexpected spectrally correlated noise” on the MLS aura, so here is another real reason why one should model and estimate noise.

To solve this inverse problem we employ a linear-Gaussian hierarchical Bayesian framework, we apply the marginal and then conditional (MTC) method [14]. Within second we can evaluate distributions over both hyper-parameters (e.g. noise and ozone smoothness) and parameters. This is a fairly new method within the Bayesian community, and we are the first, to our knowledge, to apply it to a forward model based on the radiative transfer equation (RTE). Then, instead of sampling from those posterior distributions, we are the first to utilise a tensor-train (TT) to approximate the posterior distribution, which enables us to provide estimates and uncertainties via quadrature or the inverse Rosenblatt transform (IRT).

Since the RTE is weakly non-linear, we approximate the RTE with an affine map, which seems to be another novelty in the field of atmospheric remote sensing. Additionally, we provide a new approach to tackle this inverse problem by jointly inferring pressure, temperature and ozone profiles given one set of measurements.

### 1.3 Thesis Structure

In Ch. 2, we give a brief overview of the methods used and provide references for more details. Then, in Ch. 3, we provide the forward model based on a simplified RTE, and discuss how to measure most effectively. Using our findings, we simulate some noisy data for an idealised limb sounder within a simplified atmosphere based on the RTE. Then, in Ch. 4, we setup our linear hierarchical Bayesian model and discuss some prior modelling. Given the simulated data, we provide posterior distributions of our Bayesian framework based on the linearised RTE to then approximate the non-linear forward model with an affine map. We compare a regularisation solution with the posterior distributions of the approximated linear Bayesian model against a ground truth ozone profile, where we also provide posterior distributions over hyper-parameters. Additionally, we condition on an ozone profile and noise sample to give joint pressure and temperature posterior profiles. Furthermore, we assess and discuss some errors of the approximation used to provide arguments for choices made regarding those approximations. Lastly, we discuss our results and provide an outlook, see Ch. 5.



# 2

## Theoretical and Technical Background

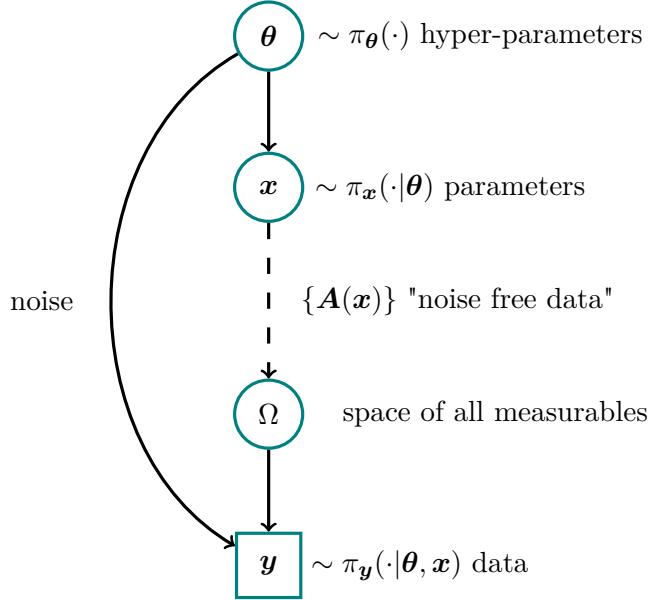
In this chapter, we provide introductions and brief derivations of the methods used in this thesis, as well as references for more details. We keep it as general as possible, as the expressions specifically tailored towards the forward map will be presented in the results Chapter 4. We begin by introducing a general hierarchical Bayesian approach to an inverse problem. Next, we provide the basics of Metropolis-Hastings sampling, more specifically, the essentials of Markov-Chain Monte Carlo (MCMC) methods. Further, we explain how we approximate functions using a Tensor-Train (TT) approach, which enables us to calculate marginals from the posterior distribution cheaply. Then, we elaborate on the Wasserstein distance for assessing upper error bounds. Lastly, we provide some background information on affine maps and the Tikhonov regularisation method.

### 2.1 Hierarchical Bayesian Inference

Assume we observe some data

$$\mathbf{y} = \mathbf{A}(\mathbf{x}) + \boldsymbol{\eta}, \quad (2.1)$$

based on a forward model  $\mathbf{A}(\mathbf{x})$ , which may be non-linear, a unknown parameter  $\mathbf{x}$  and some additive random noise  $\boldsymbol{\eta}$ . Naturally, due to the noise, which we classify through a hyper-parameter, we deal with a random process, and we wish to include that in our hierarchically ordered modelling. Then define the likelihood function  $\pi(\mathbf{y}|\boldsymbol{\theta}, \mathbf{x})$  according to the nature of the noise as well as all relevant information about the measurement process, captured by the model  $\mathbf{A}(\mathbf{x})$ . We read  $\pi(\mathbf{y}|\boldsymbol{\theta}, \mathbf{x})$  as the distribution over  $\mathbf{y}$  conditioned on  $\mathbf{x}$  and  $\boldsymbol{\theta}$ , and  $\boldsymbol{\eta} \sim \pi_{\boldsymbol{\eta}}(\cdot|\boldsymbol{\theta})$  as  $\boldsymbol{\eta}$  is distributed as  $\pi_{\boldsymbol{\eta}}(\cdot|\boldsymbol{\theta})$ . Here  $\boldsymbol{\theta}$  may account for multiple hyper-parameters, e.g. describing the distribution  $\pi_{\boldsymbol{\eta}}(\cdot|\boldsymbol{\theta})$  over the



**Figure 2.1:** The directed acyclic graph (DAG) for an inverse problem visualises statistical dependencies as solid line arrows and deterministic dependencies as dotted arrows. The hyper-parameters  $\theta$  are distributed as the hyper-prior distribution  $\pi(\theta)$ . The prior distribution  $\pi_x(\cdot|\theta)$  for the parameter  $x$  and the noise are statistically dependent on some of those hyper-parameters. Then a parameter  $x \sim \pi_x(\cdot|\theta)$  is mapped onto the space of all measurables  $\Omega = A(x)$  deterministically through the forward model. From the space of all measurable noise free data we observe a data set  $y = A(x) + \eta$  with some random noise  $\eta \sim \pi_\eta(\cdot|\theta)$ , which determines the likelihood function  $\pi(y|\theta, x)$ .

noise vector  $\eta$ , and describing physical properties or functional dependencies of  $x$ , e.g. smoothness, through the prior distribution  $\pi(x|\theta)$ . Consequently we define a hyper-prior distribution  $\pi(\theta)$ , where  $\pi(x, \theta) = \pi(x|\theta)\pi(\theta)$ . Choosing these prior distributions is ultimately a modeller's choice and crucial, as it shall not affect the posterior distribution

$$\pi(x, \theta|y) = \frac{\pi(y|x, \theta)\pi(x, \theta)}{\pi(y)} \propto \pi(y|x, \theta)\pi(x, \theta), \quad (2.2)$$

which according to Bayes' theorem, gives us a distribution of  $x$  and  $\theta$  given (conditioned on) the data. Note that here we include the hyper-parameters within the posterior distribution, which is the key idea of hierarchical Bayesian modelling, as we not only aim to quantify the posterior distribution over the parameters  $x$ , but also the posterior distribution over the hyper-parameter  $\theta$ . We can visualise this hierarchically-ordered correlation structure between parameters as well as how distributions progress through a measurement process, using a directed acyclic graph (DAG), see Figure 2.1.

The expectation of any function  $h(x_\theta)$ , where  $x$  may depend on  $\theta$ , is described as

$$E_{x, \theta|y}[h(x_\theta)] = \underbrace{\int \int h(x_\theta) \pi(x, \theta|y) dx d\theta}_{\mu_{\text{int}}}. \quad (2.3)$$

If that is a high-dimensional integral and computationally not feasible to solve, we approximate

$$\mathbb{E}_{\mathbf{x}, \boldsymbol{\theta} | \mathbf{y}}[h(\mathbf{x}_{\boldsymbol{\theta}})] \approx \underbrace{\frac{1}{N} \sum_{k=1}^N h(\mathbf{x}_{\boldsymbol{\theta}}^{(k)})}_{\boldsymbol{\mu}_{\text{samp}}}, \quad (2.4)$$

with an unbiased sample-based Monte Carlo estimate [15] for large enough  $N$  (law of large numbers [16, Chapter 17]). Here, the samples  $\{\mathbf{x}^{(k)}, \boldsymbol{\theta}^{(k)}\} \sim \pi_{\mathbf{x}, \boldsymbol{\theta}}(\cdot | \mathbf{y})$ , for  $k = 1, \dots, N$ , form a sample set  $\mathcal{M} = \{(\mathbf{x}, \boldsymbol{\theta})^{(1)}, \dots, (\mathbf{x}, \boldsymbol{\theta})^{(N)}\}$ .

Generating a representative sample set quickly from the posterior distribution often presents a significant challenge. This is mainly due to the strong correlations that usually exist between the parameters and hyper-parameters, as discussed by Rue and Held in [17] and illustrated in Appendix A.1. If  $\mathbf{x}$  cannot be parametrised directly in terms of the hyper-parameters  $\boldsymbol{\theta}$ , so that  $\mathbf{x}(\boldsymbol{\theta})$  is function of  $\boldsymbol{\theta}$ , it is beneficial to factorise the posterior distribution as

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

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

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

over the hyper-parameters  $\boldsymbol{\theta}$  (see [14, Lemma 2]). In [18], they classify inverse problems into problems with known or unknown conditional posterior distributions, and conclude that if  $\pi(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y}) = \pi(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x} | \boldsymbol{\theta}) / \pi(\mathbf{y} | \boldsymbol{\theta})$  has a known form, the normalising constant of  $\pi(\boldsymbol{\theta} | \mathbf{y})$  is available

$\int \pi(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} = \pi(\mathbf{y} | \boldsymbol{\theta}) \propto \pi(\boldsymbol{\theta} | \mathbf{y}) / \pi(\boldsymbol{\theta})$  and one can almost surely determine the  $\boldsymbol{\theta}$ -dependence of the marginal posterior  $\pi(\boldsymbol{\theta} | \mathbf{y})$ .

This approach, known as the marginal and then conditional (MTC) method, is particularly advantageous when  $\mathbf{x} \in \mathbb{R}^n$  is high-dimensional, while  $\boldsymbol{\theta} \in \mathbb{R}^{n_{\boldsymbol{\theta}}}$  is low-dimensional, so that  $n_{\boldsymbol{\theta}} \ll n$  and evaluation of  $\pi(\boldsymbol{\theta} | \mathbf{y})$  is cheap. Applying the law of total expectation [19], Eq. (2.3) becomes

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

where, in the case of a linear-Gaussian hierarchical Bayesian model, both the marginal distribution and the inner expectation  $\mathbb{E}_{\mathbf{x} | \boldsymbol{\theta}, \mathbf{y}}[h(\mathbf{x}_{\boldsymbol{\theta}})]$  are well defined (see next subsection). Furthermore, the central limit theorem states that the sample mean  $\boldsymbol{\mu}_{\text{samp}}^{(i)}$ , of independent sample sets  $\mathcal{M}_i$  for  $i = 1, \dots, n$  of any distribution, converges to be normally distributed, so that

$$\sqrt{n}(\boldsymbol{\mu}_{\text{samp}}^{(i)} - \boldsymbol{\mu}_{\text{int}}) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma^2)[20], \quad (2.8)$$

and if  $\sigma^2 < \infty$  the Monte-Carlo error  $\boldsymbol{\mu}_{\text{samp}}^{(i)} - \boldsymbol{\mu}_{\text{int}}$  is bounded.

### Integrated Autocorrelation time

To assess the error  $\sigma^2$  of chain  $\mathcal{M}_i$ , we ignore systematic error due to initialisation bias (burn-in period), but we have to take into account that samples produced by any system or algorithm are correlated. To derive the integrated autocorrelation time (IATC), we follow the lecture notes [21]. In general, the error of a Monte-Carlo-based estimate from a sample set  $\mathcal{M}_i = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}, \dots, \mathbf{x}^{(s)}, \dots, \mathbf{x}^{(N)}\} \sim \pi(\mathbf{x}|\mathbf{y})$  is:

$$(\sigma^{(i)})^2 = \text{var}(\boldsymbol{\mu}_{\text{samp}}^{(i)}) = \text{var}(\mathbb{E}_{\mathbf{x}|\mathbf{y}}[h(\mathbf{x})]) = \left( \frac{1}{N} \sum_{k=1}^N h(\mathbf{x}^{(k)}) - \boldsymbol{\mu}^{(i)} \right)^2. \quad (2.9)$$

Expanding this summation, we see that

$$\left( \frac{1}{N} \sum_{k=1}^N h(\mathbf{x}^{(k)}) - \boldsymbol{\mu}^{(i)} \right)^2 = \frac{\text{var}(\boldsymbol{\mu}_{\text{samp}}^{(i)})}{N^2} \sum_{k,s=1}^N \rho(k-s), \quad (2.10)$$

with the normalised auto correlation coefficient  $\rho(k-s) = \Gamma(k-s)/\Gamma(0)$  at lag  $k-s$ , where the auto correlation coefficient  $\Gamma(k-s) = (h(\mathbf{x}^{(k)}) - \boldsymbol{\mu}^{(i)})(h(\mathbf{x}^{(s)}) - \boldsymbol{\mu}^{(i)})$  and  $\Gamma(0) = \text{var}(\boldsymbol{\mu}_{\text{samp}}^{(i)})$  for  $k=s$ . Typically  $\Gamma(t)$  decays exponentially so that, for  $N \gg \tau$ ,  $\Gamma(t) \xrightarrow{t \rightarrow \infty} \exp\{-|t|/\tau\}$  and we can approximate

$$(\sigma^{(i)})^2 \approx \frac{\text{var}(h(\mathbf{x}))}{N} \underbrace{\sum_{t=-\infty}^{\infty} \rho(t)}_{:=2\tau_{\text{int}}} = \text{var}(h(\mathbf{x})) \frac{2\tau_{\text{int}}}{N}, \quad (2.11)$$

where we define the IATC as in [21, pp. 103-105]. See Appendix A.2 and [22] for a more detailed derivation. The IACT provides a good estimate of how many steps the sampling algorithm needs to take to produce one independent sample, accordingly we define the effective sample size as  $\frac{2\tau_{\text{int}}}{N}$ . Since this is an estimate [23] provides a way to not only calculate the IATC but also to quantify the errors of the estimate of the IATC.

### Linear-Gaussian hierarchical Bayesian model

In case of normally distributed noise  $\boldsymbol{\eta} \sim \mathcal{N}(0, \boldsymbol{\Sigma}(\boldsymbol{\theta}))$ , with zero mean and covariance  $\boldsymbol{\Sigma}(\boldsymbol{\theta})$ , and a linear model  $\mathbf{A}$ , the data is given as

$$\mathbf{y} = \mathbf{Ax} + \boldsymbol{\eta}, \quad (2.12)$$

and we can derive the marginal and conditional posterior distribution explicitly. We define our hierarchical Bayesian model as

$$\mathbf{y}|\mathbf{x}, \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{Ax}, \boldsymbol{\Sigma}(\boldsymbol{\theta})) \quad (2.13a)$$

$$\mathbf{x}|\boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{Q}(\boldsymbol{\theta})^{-1}) \quad (2.13b)$$

$$\boldsymbol{\theta} \sim \pi(\boldsymbol{\theta}), \quad (2.13c)$$

with a Gaussian likelihood function  $\pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$ , the prior mean  $\boldsymbol{\mu}$ , prior precision  $\mathbf{Q}(\boldsymbol{\theta})$  and a hyper-prior distribution  $\pi(\boldsymbol{\theta})$ . To derive the marginal posterior and the conditional posterior distribution, we consider the joint multivariate Gaussian distribution

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N} \left[ \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} \end{pmatrix}, \begin{pmatrix} \mathbf{Q}(\boldsymbol{\theta}) + \mathbf{A}^T \Sigma(\boldsymbol{\theta})^{-1} \mathbf{A} & -\mathbf{A}^T \Sigma(\boldsymbol{\theta})^{-1} \\ \Sigma(\boldsymbol{\theta})^{-1} \mathbf{A} & \Sigma(\boldsymbol{\theta})^{-1} \end{pmatrix}^{-1} \right], \quad (2.14)$$

where we provide the joint precision matrix as in [24], see also [14, 17]. Immediately, we formulate the conditional posterior as

$$\mathbf{x}|\mathbf{y}, \boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\mu} + (\mathbf{Q}(\boldsymbol{\theta}) + \mathbf{A}^T \Sigma(\boldsymbol{\theta})^{-1} \mathbf{A})^{-1}(\mathbf{y} - \mathbf{A}\boldsymbol{\mu}), (\mathbf{Q}(\boldsymbol{\theta}) + \mathbf{A}^T \Sigma(\boldsymbol{\theta})^{-1} \mathbf{A})^{-1}). \quad (2.15)$$

Then the marginal posterior distribution over the hyper-parameters can be derived as in Eq. 2.6, where, as noted in [14], the parameter  $\mathbf{x}$  cancels and we arrive at

$$\begin{aligned} \pi(\boldsymbol{\theta}|\mathbf{y}) \propto & \sqrt{\frac{\det(\Sigma(\boldsymbol{\theta})^{-1}) \det(\mathbf{Q}(\boldsymbol{\theta}))}{\det(\mathbf{Q}(\boldsymbol{\theta}) + \mathbf{A}^T \Sigma(\boldsymbol{\theta})^{-1} \mathbf{A})}} \exp \left\{ -\frac{1}{2}(\mathbf{y} - \mathbf{A}\boldsymbol{\mu})^T \right. \\ & \left. [\Sigma(\boldsymbol{\theta})^{-1} - \Sigma(\boldsymbol{\theta})^{-1} \mathbf{A}(\mathbf{Q}(\boldsymbol{\theta}) + \mathbf{A}^T \Sigma(\boldsymbol{\theta})^{-1} \mathbf{A})^{-1} \mathbf{A}^T \Sigma(\boldsymbol{\theta})^{-1}] (\mathbf{y} - \mathbf{A}\boldsymbol{\mu}) \right\} \pi(\boldsymbol{\theta}). \end{aligned} \quad (2.16)$$

Having the marginal posterior distribution available breaks up the correlation structure between  $\mathbf{x}$  and  $\boldsymbol{\theta}$  (see Appendix A.1), and makes the MTC approach very efficient [14]. This scheme evaluates the marginal posterior values first and then conditions on hyper-parameters to draw posterior samples  $\mathbf{x} \sim \pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta})$  or to evaluate expectation and variance of  $\pi(\mathbf{x}|\mathbf{y})$  by integration over the marginal posterior.

## 2.2 Sampling Methods

In this section we present the underlying methodology of the sampling methods used in this thesis and show how these methods draw samples

$\mathcal{M} = \{(\mathbf{x}, \boldsymbol{\theta})^{(1)}, \dots, (\mathbf{x}, \boldsymbol{\theta})^{(k)}, \dots, (\mathbf{x}, \boldsymbol{\theta})^{(N)}\} \sim \pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$  from the desired target distribution, so that we can calculate sample-based estimates as in Eq. 2.4. Here,  $\mathcal{M}$  denotes a Markov chain, where each new sample  $(\mathbf{x}, \boldsymbol{\theta})^{(k)}$  is only affected by the previous one,  $(\mathbf{x}, \boldsymbol{\theta})^{(k-1)}$ . Markov chain Monte Carlo (MCMC) methods generate such a chain  $\mathcal{M}$  using random (Monte Carlo) proposals  $(\mathbf{x}, \boldsymbol{\theta})^{(k)} \sim q(\cdot | (\mathbf{x}, \boldsymbol{\theta})^{(k-1)})$  according to a proposal distribution conditioned on the previous sample (Markov), where ergodicity of the chain  $\mathcal{M}$  is a sufficient criterion for using sample-based estimates [7, 15].

The ergodicity theorem in [7] states that, if a Markov chain  $\mathcal{M}$  is aperiodic, irreducible, and reversible, then it converges to a unique stationary equilibrium distribution. In other words, if the chain can reach any state from any other state (irreducibility), is not stuck

in periodic cycles (aperiodicity), and is reversible (detailed balance condition [7]). Then the chain converges to the desired target distribution with  $\mathcal{M} \sim \pi(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})$ . In practice, one can inspect the trace  $\pi(\mathbf{x}^{(k)}, \boldsymbol{\theta}^{(k)} | \mathbf{y})$  for  $k = 1, \dots, N$  and visually assess convergence and mixing properties of the chain to evaluate ergodicity. The sampling methods used in this thesis possess proven ergodic properties, and we therefore refer the reader to the corresponding literature for further details.

### 2.2.1 Metropolis within Gibbs

As in Eq. 2.5, when using the MTC method we sample from  $\pi(\boldsymbol{\theta} | \mathbf{y})$  first and then determine the full conditional  $\pi(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})$  as in Eq. 2.7. To sample from  $\pi(\boldsymbol{\theta} | \mathbf{y})$ , we use a Metropolis-within-Gibbs (MWG) sampler as described in [14]. We apply the MWG sample for the two-dimensional case only, with  $\boldsymbol{\theta} = (\theta_1, \theta_2)$ , where we perform a Metropolis step in the  $\theta_1$  direction and a Gibbs step in the  $\theta_2$  direction. Ergodicity for this approach is proven in [25].

The Metropolis-within-Gibbs algorithm begins with an initial guess  $\boldsymbol{\theta}^{(t)}$  at  $t = 0$ . We then propose a new sample  $\theta_1 \sim q(\theta_1 | \theta_1^{(t-1)})$ , conditioned on the previous state, using a symmetric proposal distribution  $q(\theta_1 | \theta_1^{(t-1)}) = q(\theta_1^{(t-1)} | \theta_1)$ , which is a special case of the Metropolis-Hastings algorithm [25]. We accept and set  $\theta_1^{(t)} = \theta_1$  with the acceptance probability

$$\alpha(\theta_1 | \theta_1^{(t-1)}) = \min \left\{ 1, \frac{\pi(\theta_1 | \theta_2^{(t-1)}, \mathbf{y}) \underline{q(\theta_1^{(t-1)} | \theta_1)}}{\pi(\theta_1^{(t-1)} | \theta_2^{(t-1)}, \mathbf{y}) \underline{q(\theta_1 | \theta_1^{(t-1)})}} \right\} \quad (2.17)$$

or reject and keep  $\theta_1^{(t)} = \theta_1^{(t-1)}$ , which we do by comparing  $\alpha$  to a uniform random number  $u \sim \mathcal{U}(0, 1)$ .

Next, we perform a Gibbs step in the  $\theta_2$  direction, where Gibbs sampling is again a special case of the Metropolis-Hastings algorithm with acceptance probability equal to one, and draw the next sample  $\theta_2^{(t)} \sim \pi(\cdot | \theta_1^{(t)}, \mathbf{y})$ , conditioned on the current value  $\theta_1^{(t)}$ .

We repeat this procedure  $N'$  times and ensure convergence independently of the initial sample (irreducibility) by discarding the initial  $N_{\text{burn-in}}$  samples after a so-called burn-in period, resulting in a Markov chain of length  $N = N' - N_{\text{burn-in}}$ .

### 2.2.2 T-walkSampler as a Black Box

If the parameters  $\mathbf{x}$  are functionally dependent on the hyper-parameters  $\boldsymbol{\theta}$ , i.e.,  $\mathbf{x} = \mathbf{x}(\boldsymbol{\theta})$ , we can sample directly from the marginal posterior  $\pi(\boldsymbol{\theta} | \mathbf{y})$  using the **t-walk** sampler as by Christen and Fox [26]. The **t-walk** is employed as a black-box sampler, requiring the specification of the number of samples, burn-in period, support region, and the target distribution. Convergence to the target distribution is guaranteed by construction of the algorithm.

**Algorithm 1:** Metropolis within Gibbs

```

1: Initialise and suppose two dimensional vector  $\boldsymbol{\theta}^{(0)} = (\theta_1^{(0)}, \theta_2^{(0)})$ 
2: for  $k = 1, \dots, N'$  do
3:   Propose  $\theta_1 \sim q(\cdot | \theta_1^{(t-1)}) = q(\theta_1^{(t-1)} | \cdot)$ 
4:   Compute

$$\alpha(\theta_1 | \theta_1^{(t-1)}) = \min \left\{ 1, \frac{\pi(\theta_1 | \theta_2^{(t-1)}, \mathbf{y}) q(\theta_1^{(t-1)} | \theta_1)}{\pi(\theta_1^{(t-1)} | \theta_2^{(t-1)}, \mathbf{y}) q(\theta_1 | \theta_1^{(t-1)})} \right\}$$

5:   Draw  $u \sim \mathcal{U}(0, 1)$ 
6:   if  $\alpha \geq u$  then
7:     Accept and set  $\theta_1^{(t)} = \theta_1$ 
8:   else
9:     Reject and keep  $\theta_1^{(t)} = \theta_1^{(t-1)}$ 
10:  end if
11:  Draw  $\theta_2^{(t)} \sim \pi(\cdot | \theta_1^{(t)}, \mathbf{y})$ 
12: end for
13: Output:  $\boldsymbol{\theta}^{(0)}, \dots, \boldsymbol{\theta}^{(k)}, \dots, \boldsymbol{\theta}^{(N')} \sim \pi(\boldsymbol{\theta} | \mathbf{y})$ 

```

**2.2.3 Randomise then Optimise**

If we can not evaluate the mean and variance of the full posterior through integration, e.g. if it is computationally not feasible to solve the integral in Eq. 2.7 due to a large number of hyper-parameters, we need an alternative way to characterise  $\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})$ . For the linear-Gaussian Bayesian model we can draw samples from the normal distribution  $\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})$  using the randomise then optimise (RTO) [27] scheme after sampling from the marginal posterior  $\pi(\boldsymbol{\theta}|\mathbf{y})$ .

The full conditional distribution can be rewritten as

$$\pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}) \propto \pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}|\boldsymbol{\theta}) \quad (2.18)$$

$$= \exp \left( - \left\| \hat{\mathbf{A}}\mathbf{x} - \hat{\mathbf{y}} \right\|_{L^2}^2 \right), \quad (2.19)$$

where we define

$$\hat{\mathbf{A}} := \begin{bmatrix} \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1/2} \mathbf{A}_{\boldsymbol{\theta}} \\ \mathbf{Q}_{\boldsymbol{\theta}}^{1/2} \end{bmatrix}, \quad \hat{\mathbf{y}} := \begin{bmatrix} \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1/2} \mathbf{y} \\ \mathbf{Q}_{\boldsymbol{\theta}}^{1/2} \boldsymbol{\mu} \end{bmatrix} \quad [28]. \quad (2.20)$$

Here we write  $\mathbf{A}(\boldsymbol{\theta}) = \mathbf{A}_{\boldsymbol{\theta}}$ ,  $\mathbf{Q}(\boldsymbol{\theta}) = \mathbf{Q}_{\boldsymbol{\theta}}$  and  $\boldsymbol{\Sigma}^{-1}(\boldsymbol{\theta}) = \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1}$ , which are all dependent on the hyper-parameters  $\boldsymbol{\theta}$ . A sample  $\mathbf{x}^{(k)} \sim \pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})$  from the conditional posterior is obtained by minimising the following equation with respect to  $\hat{\mathbf{x}}$ :

$$\mathbf{x}^{(k)} = \arg \min_{\hat{\mathbf{x}}} \| \hat{\mathbf{A}}\hat{\mathbf{x}} - (\hat{\mathbf{y}} + \mathbf{b}) \|_{L^2}^2, \quad \mathbf{b} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (2.21)$$

where we add a randomised perturbation  $\mathbf{b}$ . Similar to Section 2.5, this expression can be rewritten as

$$\left( \mathbf{A}_{\boldsymbol{\theta}}^T \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \mathbf{A}_{\boldsymbol{\theta}} + \mathbf{Q}_{\boldsymbol{\theta}} \right) \mathbf{x}^{(k)} = \mathbf{A}_{\boldsymbol{\theta}}^T \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \mathbf{y} + \mathbf{Q}_{\boldsymbol{\theta}} \boldsymbol{\mu} + \mathbf{v}_1 + \mathbf{v}_2, \quad (2.22)$$

where the term  $-\hat{\mathbf{A}}^T \mathbf{b}$  is decomposed as  $\mathbf{v}_1 + \mathbf{v}_2$ , with  $\mathbf{v}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{A}_{\boldsymbol{\theta}}^T \boldsymbol{\Sigma}_{\boldsymbol{\theta}}^{-1} \mathbf{A}_{\boldsymbol{\theta}})$  and  $\mathbf{v}_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_{\boldsymbol{\theta}})$ , representing independent Gaussian random variables [14, 27].

If the Markov chain over the marginal posterior  $\pi(\boldsymbol{\theta}|\mathbf{y})$  is ergodic, and the conditional samples  $\mathbf{x}^{(k)} \sim \pi(\mathbf{x}|\boldsymbol{\theta}^{(k)}, \mathbf{y})$  are drawn independently, then the resulting joint chain  $\{(\mathbf{x}, \boldsymbol{\theta})^{(1)}, \dots, (\mathbf{x}, \boldsymbol{\theta})^{(N)}\} \sim \pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y}) = \pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})\pi(\boldsymbol{\theta}|\mathbf{y})$  is also ergodic [29].

## 2.3 Numerical Approximation – Tensor-Train (TT)

Instead of sampling from a target distribution  $\pi(\mathbf{x})$  we can approximate that distribution on a d-dimensional grid with far fewer function evaluation compared to sampling methods using a tensor-train (TT) approximation  $\tilde{\pi}(\mathbf{x}) \approx \pi(\mathbf{x})$ , with  $\mathbf{x} \in \mathbb{R}^d$ . Here, we explain how to calculate marginal distribution from an approximated probability density in the TT format and generate samples via the inverse Rosenblatt transform (IRT), following the notation of [30],

As in [30], we can define the parameter space as the product space  $\mathcal{X} = \mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_d$  with  $x_k \in \mathcal{X}_k \subseteq \mathbb{R}$  and  $\mathbf{x} = (x_1, \dots, x_k, \dots, x_d)$ . Then marginal density function for the  $k$ -th component is then given by

$$f_{X_k}(x_k) = \frac{1}{z} \int_{\mathcal{X}_1} \dots \int_{\mathcal{X}_d} \lambda(\mathbf{x}) \pi(\mathbf{x}) \, dx_1 \dots dx_{k-1} \, dx_{k+1} \dots dx_d, \quad (2.23)$$

where we integrate over all dimensions except the  $k$ -th, and  $z$  is a normalisation constant. Here, we introduce a weight function  $\lambda(x)$ , which can be useful for quadrature rules [31], to which [30] refer to as a "product-form Lebesgue-measurable weighting function" and define it as

$$\lambda(\mathcal{X}) = \prod_{i=1}^d \lambda_i(\mathcal{X}_i), \quad \text{where } \lambda_i(\mathcal{X}_i) = \int_{\mathcal{X}_i} \lambda_i(x_i) \, dx_i. \quad (2.24)$$

In the TT format, the integral in Eq. 2.23 for the marginal probability can be computed at a low computational cost as  $\pi(\mathbf{x})$  is approximated by

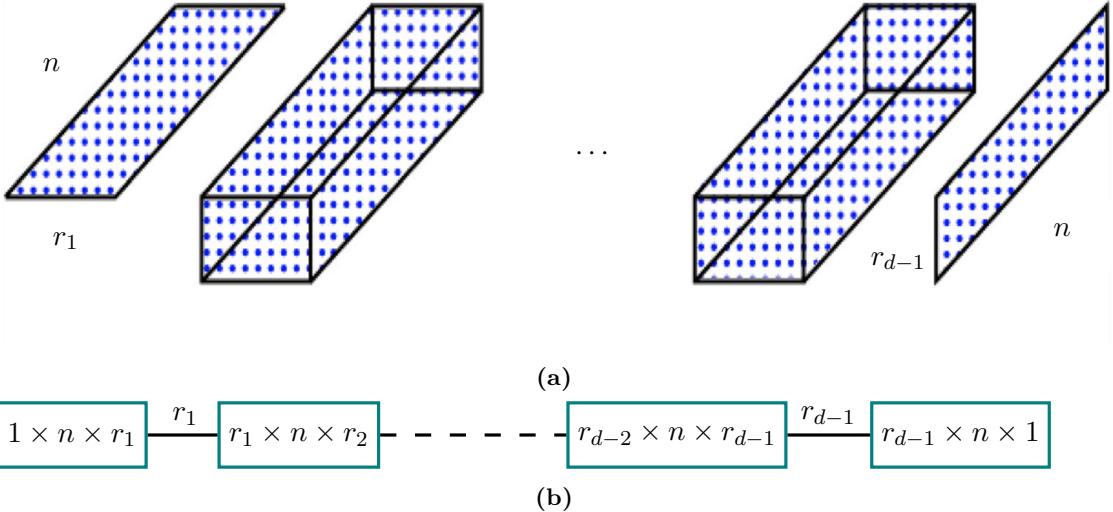
$$\tilde{\pi}(\mathbf{x}) = \tilde{\pi}_1(x_1) \tilde{\pi}_2(x_2) \dots \tilde{\pi}_d(x_d) \in \mathbb{R},$$

which is a sequence of matrix multiplications, with  $\tilde{\pi}_k(x_k) \in \mathbb{R}^{r_{k-1} \times r_k}$  for a fixed point  $\mathbf{x} = (x_1, \dots, x_d)$  on a predefined  $d$ -dimensional discrete univariate grid over the parameter space  $\mathcal{X}$ . We call  $\tilde{\pi}_k \in \mathbb{R}^{r_{k-1} \times n \times r_k}$  a TT-core with ranks  $r_{k-1} = r_k = r$ , where the outer ranks are  $r_0 = r_d = 1$ , representing each dimension on  $n$  grid points and connecting to neighbouring

dimensions through its ranks. This enables us to approximate  $\pi(\mathcal{X}) \approx \tilde{\pi}_1 \tilde{\pi}_2 \cdots \tilde{\pi}_d$  over the parameter space  $\mathcal{X}$  using  $2nr + (d - 2)nr^2$  evaluation points, as illustrated in Figure 2.2, instead of  $n^d$  function evaluation. Consequently, the marginal target distribution

$$\begin{aligned} f_{X_k}(x_k) &\approx \frac{1}{z} \left| \left( \int_{\mathcal{X}_1} \lambda_1(x_1) \tilde{\pi}_1(x_1) dx_1 \right) \cdots \left( \int_{\mathcal{X}_{k-1}} \lambda_{k-1}(x_{k-1}) \tilde{\pi}_{k-1}(x_{k-1}) dx_{k-1} \right) \right. \\ &\quad \left. \lambda_k(x_k) \tilde{\pi}_k(x_k) \right. \\ &\quad \left( \int_{\mathcal{X}_{k+1}} \lambda_{k+1}(x_{k+1}) \tilde{\pi}_{k+1}(x_{k+1}) dx_{k+1} \right) \cdots \left( \int_{\mathcal{X}_d} \lambda_d(x_d) \tilde{\pi}_d(x_d) dx_d \right) \right| \end{aligned} \quad (2.25)$$

is computed by integrating over all TT cores except  $\pi_k$ , as in [32], including a normalisation constant  $z$  [30].



**Figure 2.2:** Here, we visualise the tensor train cores as two- and three-dimensional matrices. Each core has a length  $n$ , corresponding to the number of grid points in one dimension, and the cores are connected through ranks  $r_k$ . More specifically, a core  $\tilde{\pi}_k$  has dimensions  $r_{k-1} \times n \times r_k$ , with outer ranks  $r_0 = r_d = 1$ . Using the TT-format enables us to represent a  $d$ -dimensional grid with only  $2nr + (d - 2)nr^2$  evaluation points instead of  $n^d$  grid points. Figure (a) is adapted from [33].

In practice, tensor train approximations may suffer from numerical instability, in particular because it is not advantageous to approximate the target function  $\pi(\mathbf{x})$  in e.g. the logarithmic space. Hence, Cui et al. [30] approximate the square root of the probability density

$$\sqrt{\pi(\mathbf{x})} \approx \tilde{g}(\mathbf{x}) = \mathbf{G}_1(x_1), \dots, \mathbf{G}_k(x_k), \dots, \mathbf{G}_d(x_d) \quad [30, \text{Eq. 18}], \quad (2.26)$$

which ensures positivity. Here, each TT-core is given by

$$G_k^{(\alpha_{k-1}, \alpha_k)}(x_k) = \sum_{i=1}^{n_k} \phi_k^{(i)}(x_k) \mathbf{A}_k[\alpha_{k-1}, i, \alpha_k], \quad k = 1, \dots, d, \quad [30, \text{Eq. 21}], \quad (2.27)$$

where  $\mathbf{A}_k \in \mathbb{R}^{r_{k-1} \times n_k \times r_k}$  is the  $k$ -th coefficient tensor and  $\{\phi_k^{(i)}(x_k)\}_{i=1}^{n_k}$  are the basis functions corresponding to the  $k$ -th coordinate. The approximated unnormalised density is written as:

$$\pi(\mathbf{x}) \approx \xi + \tilde{g}(\mathbf{x})^2 \quad [30, \text{Eq. 19}], \quad (2.28)$$

where  $\xi$  is a positive constant added according to the ratio of the Lebesgue weighted L2-norm error and the Lebesgue weighting (see Eq. 2.24) such that

$$0 \leq \xi \leq \frac{1}{\lambda(\mathcal{X})} \|\tilde{g} - \sqrt{\pi}\|_{L_\lambda^2(\mathcal{X})}^2 \quad [30, \text{Eq. 35}]. \quad (2.29)$$

This leads to the normalised target function

$$f_X(\mathbf{x}) \approx \frac{1}{z} \left( \lambda(\mathbf{x}) \xi + \lambda(\mathbf{x}) \tilde{g}(\mathbf{x})^2 \right) \quad [30, \text{Eq. 19}], \quad (2.30)$$

which is the normalisation constant  $z = \int_{\mathcal{X}} f_X(\mathbf{x}) d\mathbf{x}$ . Given the tensor train approximation of  $\sqrt{\pi}$ , the marginal function  $f_{X_k}(x_k)$  can be expressed as

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

### 2.3.1 Marginal Functions

TT-approximations are handy when approximating integrals, as marginal functions can be easily computed which may simplify the integration significantly. We compute those by a procedure, to which Cui et al. [30] refer to as backwards marginalisation, see Prop. 2, and to which I add the forward marginalisation, see Prop. 1. This is similar to the left and right orthogonalisation of TT-cores [34, 35]. The backwards marginalisation provides us with the coefficient matrices  $\mathbf{B}_k$ , while the forward marginalisation gives the coefficient matrices  $\mathbf{B}_{\text{pre},k}$ . These matrices enable the efficient evaluation of marginal functions since they integrate over the coordinates either left or right of the  $k$ -th dimension, as in [30]. In doing so, we define the mass matrix  $\mathbf{M}_k \in \mathbb{R}^{n_k \times n_k}$  as

$$\mathbf{M}_k[i, j] = \int_{\mathcal{X}_k} \phi_k^{(i)}(x_k) \phi_k^{(j)}(x_k) \lambda(x_k) dx_k, \quad i, j = 1, \dots, n_k, \quad [30, \text{Eq. 22}], \quad (2.32)$$

where  $\{\phi_k^{(i)}(x_k)\}_{i=1}^{n_k}$  denotes the set of basis functions for the  $k$ -th coordinate. The proposition used to compute  $\mathbf{B}_k$ , stated in Prop. 1, is adapted directly from [30].

**Proposition 1** (Backwards Marginalisation as in [30]): 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  $\mathbf{C}_k \in \mathbb{R}^{r_{k-1} \times n_k \times r_k}$ :

$$\mathbf{C}_k[\alpha_{k-1}, \tau, l_k] = \sum_{i=1}^{n_k} \mathbf{B}_k[\alpha_{k-1}, i, l_k] \mathbf{L}_k[i, \tau] \quad [30, \text{Eq. (27)}]. \quad (2.33)$$

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

$$\mathbf{Q}_k \mathbf{R}_k = (\mathbf{C}_k^{(R)})^\top \quad [30, \text{Eq. 28}]. \quad (2.34)$$

3. Compute the new coefficient tensor:

$$\mathbf{B}_{k-1}[\alpha_{k-2}, i, l_{k-1}] = \sum_{\alpha_{k-1}=1}^{r_{k-1}} \mathbf{A}_{k-1}[\alpha_{k-2}, i, \alpha_{k-1}] \mathbf{R}_k[l_{k-1}, \alpha_{k-1}] \quad [30, \text{Eq. 29}]. \quad (2.35)$$

**Proposition 2** (Forward Marginalisation): Starting with the first coordinate  $k = 1$ , we set  $\mathbf{B}_{\text{pre},1} = \mathbf{A}_1$ . The following procedure can be used to obtain the coefficient tensor  $\mathbf{B}_{\text{pre},k+1} \in \mathbb{R}^{r_k \times n_{k+1} \times r_{k+1}}$  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  $\mathbf{C}_k \in \mathbb{R}^{r_{k-1} \times n_k \times r_k}$ :

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

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

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

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

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

After computing the coefficient tensors  $\mathbf{B}_{\text{pre},k+1}$  as in Prop. 2 and  $\mathbf{B}_{k+1}$  from Prop. 1,

the marginal PDF of  $k$ -th dimension can be expressed as

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

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

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

For the first dimension,  $f_{X_1}(x_1)$  can be expressed as

$$f_{X_1}(x_1) \approx \frac{1}{z} \left( \xi \prod_{i=2}^d \lambda_i(\mathcal{X}_i) + \sum_{l_1=1}^{r_1} \left( \sum_{i=1}^n \phi_1^{(i)}(x_1) \mathbf{D}_1[i, l_1] \right)^2 \right) \lambda_1(x_1) \quad [30, \text{Eq. 30}], \quad (2.41)$$

where  $\mathbf{D}_1[i, l_1] = \mathbf{B}_1[\alpha_0, i, l_1]$  and  $\alpha_0 = 1$ , and similarly in the last dimension

$$f_{X_d}(x_d) \approx \frac{1}{z} \left( \xi \prod_{i=1}^{d-1} \lambda_i(\mathcal{X}_i) + \sum_{l_{n-1}=1}^{r_{d-1}} \left( \sum_{i=1}^n \phi_d^{(i)}(x_d) \mathbf{D}_d[l_{n-1}, i] \right)^2 \right) \lambda_d(x_d), \quad (2.42)$$

where  $\mathbf{D}_d[l_{n-1}, i] = \mathbf{B}_{\text{pre},d}[l_{n-1}, i, \alpha_{n+1}]$  and  $\alpha_{d+1} = 1$ . Note that we calculate the normalisation numerically within the process of finding the marginals so that  $\sum f_{X_k}(x_k) = 1$ .

### 2.3.2 Sampling from a TT Approximation

If instead of evaluating integrals we like to draw samples from the approximated function we do this via the inverse Rosenblatt transform (IRT), as in [32], to preserve the correlation structure. Since we approximate the square root of the target function, Cui et. al. [30] call that the squared inverse Rosenblatt transform (SIRT).

#### Algorithm 2: Squared Inverse Rosenblatt Transform (SIRT)

```

1: Input: seeds  $\{\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(N)}\} \sim \mathcal{U}(0, 1)^d$  and  $\mathbf{B}_1, \dots, \mathbf{B}_d$  from Prop. 1
2: for  $s = 1, \dots, N$  do
3:   for  $k = 1, \dots, d$  do
4:     compute normalised PDF  $f_{X_k|X_{<k}}(x_k|x_{k-1}^{(s)}, \dots, x_1^{(s)})$ , Eq. 2.45
5:     compute cumulative distribution function  $F_{X_k|X_{<k}}(x_k)$ , Eq. 2.43,
6:     project sample  $x_k^{(s)} = F_{X_k|X_{<k}}^{-1}(u_k^{(s)})$ 
7:     interpolate  $\mathbf{G}_k(x_k^{(s)})$ , Eq. 2.44
8:     update  $\mathbf{G}_{\leq k}(x_{\leq k}^{(s)}) = \mathbf{G}_{<k}(x_{<k}^{(s)}) \mathbf{G}_k(x_k^{(s)})$ 
9:   end for
10: end for
11: Output: samples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$ , where each  $\mathbf{x}^{(s)} \in \mathbb{R}^d$  for  $s = 1, \dots, N$ 

```

We start by calculating the Backward marginals  $\mathbf{B}_1, \dots, \mathbf{B}_d$  as in Prop. 1 and draw  $N$  uniformly distributed seeds  $\{\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(N)}\} \sim \mathcal{U}(0, 1)^d$ , where each  $\mathbf{u}^{(s)}$  is  $d$ -dimensional

for  $s = 1, \dots, N$ . Then we calculate the first marginal  $f_{X_1}(x_1)$  as in Eq. 2.41 and normalise with  $z = \int_{\mathcal{X}_1} f_{X_1}(x_1) dx_1$ . Next, we compute the cumulative distribution function (CDF)

$$F_{X_k|X_{<k}}(x_k) \approx \int_{-\infty}^{x_k} f_{X_k|X_{<k}}(\hat{x}_k|x_{k-1}, \dots, x_1) d\hat{x}_k [30, \text{Eq. 17}] \quad (2.43)$$

for the first dimension  $k = 1$  and then project the seed on the parameter space  $x_k^{(s)} = F_{X_k|X_{<k}}^{-1}(u_k^{(s)})$ . Once that is done, we use a piecewise polynomial interpolation

$$\mathbf{G}_k(x_k^{(s)}) \approx \frac{x_k^{(s)} - x_k^{(i)}}{x_k^{(i+1)} - x_k^{(i)}} \mathbf{G}_k(x_k^{(i+1)}) + \frac{x_k^{(i+1)} - x_k^{(s)}}{x_k^{(i+1)} - x_k^{(i)}} \mathbf{G}_k(x_k^{(i)}), \quad (2.44)$$

for  $x_k^{(i)} \leq x_k^{(s)} \leq x_k^{(i+1)}$  in between two grid points  $i$  and  $i + 1$  as in [32]. Through  $\mathbf{G}_k(x_k^{(s)}) \in \mathbb{R}^{1 \times r_{k-1}}$  we condition on the previous samples, which denotes the product of all approximated tensors of the previous  $k - 1$  samples to preserve the correlation structure. Then we marginalise over the dimensions  $k + 1, \dots, d$  via  $\mathbf{B}_k$  so that the next "conditional marginal" is given as:

$$\begin{aligned} f_{X_k|X_{<k}}(x_k|x_{k-1}^{(s)}, \dots, x_1^{(s)}) &\approx \frac{1}{z} \left( \xi \prod_{i=k+1}^d \lambda_i(X_i) + \right. \\ &\quad \left. \sum_{l_k=1}^{r_k} \left( \sum_{i=1}^n \phi_k^{(i)}(x_k^{(s)}) \left( \sum_{\alpha_{k-1}=1}^{r_{k-1}} \mathbf{G}_{<k}^{(\alpha_{k-1})}(x_{<k}^{(s)}) \mathbf{B}_k[\alpha_{k-1}, i, l_k] \right) \right)^2 \right) \lambda_k(x_k) [30, \text{Eq. 31}]. \end{aligned} \quad (2.45)$$

We repeat the procedure for each  $u_k^{(s)} \in \mathbf{u}^{(s)}$  to gain samples  $\mathbf{x}^{(s)} \sim f_X(\mathbf{x})$ , see algorithmic box 3 for a summarised version.

### MH - correction step

Since the samples using the SIRT scheme are samples from an approximation it is sensible to correct those using a Metropolis-Hastings (MH) importance sampler. In doing so we compute the acceptance probability  $\alpha = \min(w^{(s+1)}/w^{(s)}, 1)$ , where

$$w(\mathbf{x}) = \frac{\pi(\mathbf{x})}{f_X(\mathbf{x})} = \frac{\pi(\mathbf{x})}{\gamma + \tilde{g}(\mathbf{x})^2} \quad (2.46)$$

is the importance ratio. In practise we calculate the importance ratio in the log space, where  $\log f_X(\mathbf{x}) = \log f_{X_1}(x_1) + \log f_{X_2|X_1}(x_2|x_1) + \dots + \log f_{X_k|X_{<k}}(x_k|x_{k-1}, \dots, x_1)$  is given as in Eq. 2.45, see [32]. We refer to this as the SIRT-MH scheme, which provides the corrected chain  $\{\mathbf{x}_{\text{MH}}^{(1)}, \dots, \mathbf{x}_{\text{MH}}^{(N)}\} \sim \pi(\mathbf{x})$ .

**Algorithm 3:** MH correction step

```

1: Input: samples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N+1)}\}$ , where each  $\mathbf{x}^{(s)} \in \mathbb{R}^d$  for  $s = 1, \dots, N + 1$ 
2: for  $s = 1, \dots, N$  do
3:   compute MH ratio  $\frac{w^{(s+1)}}{w^{(s)}} = \frac{\pi(\mathbf{x}^{(s+1)})}{\pi(\mathbf{x}^{(s)})} \frac{f_X(\mathbf{x}^{(s)})}{f_X(\mathbf{x}^{(s+1)})}$ 
4:   compute acceptance probability  $\alpha = \min(w^{(s+1)}/w^{(s)}, 1)$ 
5:   Draw  $u \sim \mathcal{U}(0, 1)$ 
6:   if  $\alpha \geq u$  then
7:     Accept and set  $\mathbf{x}_{\text{MH}}^{(s+1)} = \mathbf{x}^{(s+1)}$ 
8:   else
9:     Reject and keep  $\mathbf{x}_{\text{MH}}^{(s+1)} = \mathbf{x}^{(s)}$ 
10:  end if
11: end for
12: Output: corrected sample chain  $\{\mathbf{x}_{\text{MH}}^{(1)}, \dots, \mathbf{x}_{\text{MH}}^{(N)}\}$ , where each  $\mathbf{x}_{\text{MH}}^{(s)} \in \mathbb{R}^d$  for  $s = 1, \dots, N$ 
```

### 2.3.3 On the Error of a TT Approximation

A straightforward way to asses the error from the TT approximation is to calculate the relative root mean squared error (RMS)

$$\left( \frac{\int_{\mathcal{X}} (\pi(\mathbf{x}) - (\gamma + \tilde{g}(\mathbf{x})^2))^2 \lambda(\mathbf{x}) d\mathbf{x}}{\int_{\mathcal{X}} \pi(\mathbf{x})^2 \lambda(\mathbf{x}) d\mathbf{x}} \right)^{1/2} = \frac{\|\pi(\mathbf{x}) - (\gamma + \tilde{g}(\mathbf{x})^2)\|_{L^2_\lambda(\mathcal{X})}}{\|\pi(\mathbf{x})\|_{L^2_\lambda(\mathcal{X})}}. \quad (2.47)$$

We can approximate the this integral as

$$\left( \frac{1}{N} \sum_{i=1}^N (\pi(\mathbf{x}^{(i)}) - (\gamma + \tilde{g}(\mathbf{x}^{(i)})^2))^2 \lambda(\mathbf{x}) \right)^{1/2} \approx \left( \int_{\mathcal{X}} (\pi(\mathbf{x}) - (\gamma + \tilde{g}(\mathbf{x})^2))^2 \lambda(\mathbf{x}) d\mathbf{x} \right)^{1/2} \quad (2.48)$$

and similarly  $\int_{\mathcal{X}} \pi(\mathbf{x})^2 \lambda(\mathbf{x}) d\mathbf{x}$ .

#### Absolute Error Bound

One way to assess the error between two distributions is to calculate the Wasserstein distance, because the Kantorovich-Rubinstein duality, as in [36, 37], says that the 1-Wasserstein distance is equal to the upper bound of differences in expectations of a function  $h$  between two probability distributions.

We define the 1-Wasserstein distance as

$$W_1(\pi, \tilde{\pi}) = \inf_{\nu \in \Pi(\pi, \tilde{\pi})} \int_{\mathcal{X} \times \mathcal{X}} c(\mathbf{x}, \tilde{\mathbf{x}}) \nu(\mathbf{x}, \tilde{\mathbf{x}}) d\mathbf{x} d\tilde{\mathbf{x}}, \quad (2.49)$$

where  $\nu$  couples  $\mathbf{x}$  and  $\tilde{\mathbf{x}}$  so that the integral over the distance  $c(\mathbf{x}, \tilde{\mathbf{x}})$  weighted by the probability measures  $\pi$  and  $\tilde{\pi}$  is the greatest lower bound of all integrals with respect to a  $\nu$  in the set of all couplings  $\Pi(\pi, \tilde{\pi})$ . Often  $\nu$  is called a transport plan, where

$c(\mathbf{x}, \tilde{\mathbf{x}})$  is the (ground) cost function, and  $\nu(\mathbf{x}, \tilde{\mathbf{x}})$  is related to the mass which has to be transported and the 1-Wasserstein distance is the earth mover distance. On the other hand (Kantorovich-Rubinstein duality), we can describe the 1-Wasserstein distance

$$W_1(\pi, \tilde{\pi}) = \sup_{\|h(\mathbf{x}) - h(\tilde{\mathbf{x}})\|_{L^2} \leq \|\mathbf{x} - \tilde{\mathbf{x}}\|_{L^2}} \left\{ \int_{\mathcal{X}} h(\mathbf{x}) d\pi(\mathbf{x}) - \int_{\mathcal{X}} h(\tilde{\mathbf{x}}) d\tilde{\pi}(\tilde{\mathbf{x}}) \right\} \quad (2.50)$$

$$= \sup_{\|h(\mathbf{x}) - h(\tilde{\mathbf{x}})\|_{L^2} \leq \|\mathbf{x} - \tilde{\mathbf{x}}\|_{L^2}} \left\{ \mathbb{E}_{\mathbf{x} \sim \pi}[h(\mathbf{x})] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{\pi}}[h(\tilde{\mathbf{x}})] \right\}. \quad (2.51)$$

as the lowest upper bound of differences in expectations of the 1-Lipschitz function  $h$  in between the two distributions  $\pi$  and  $\tilde{\pi}$ , with distance measure  $c(\mathbf{x}, \tilde{\mathbf{x}}) = \|\mathbf{x} - \tilde{\mathbf{x}}\|_{L^2}$  for  $\mathcal{X} \in \mathbb{R}^d$ . For two sample sets  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} \sim \pi$  and  $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(M)}\} \sim \tilde{\pi}$  the calculation of the Wasserstein distance becomes an optimisation problem, that is to find the best coupling of samples weighted by their distribution value according to an appropriate distance measure [38]. More specifically the 1-Wasserstein distance becomes

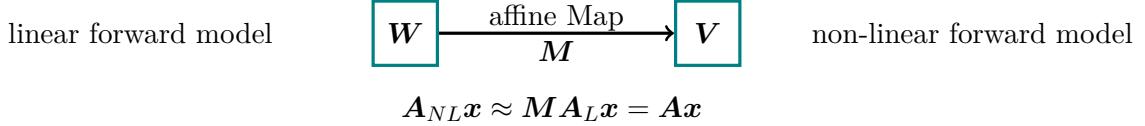
$$W_1(\pi, \tilde{\pi}) = \min_{\nu \in \Pi(\pi, \tilde{\pi})} \sum_{j=1}^M \sum_{i=1}^N \nu_{ij} \|\mathbf{x}^{(i)} - \tilde{\mathbf{x}}^{(j)}\|_{L^2}, \quad (2.52)$$

where the transport plan  $\nu \in \mathbb{R}_{\geq 0}^{N \times M}$  defines the coupling  $\nu_{ij} \in \nu$  as  $\nu_{ij} := \pi(\mathbf{x}^{(i)}) \tilde{\pi}(\tilde{\mathbf{x}}^{(j)})$  similar to [38, Eq. 3.166]. Additionally we require that  $\sum_{i=1}^N \pi(\mathbf{x}^{(i)}) = \sum_{j=1}^M \tilde{\pi}(\tilde{\mathbf{x}}^{(j)}) = 1$ . This gives us an upper bound of the absolute error in between the expected value of any 1-Lipschitz function  $h$ , e.g the upper bound of absolute differences in means related to the probability measures  $\pi$  and  $\tilde{\pi}$ .

## 2.4 Affine Map

The forward map, which we introduce in Ch. 3, poses a weakly non-linear forward problem, which we could tackle by treating the problem as a linear problem and then iteratively updating the non-linear part after each parameter sample. Instead, we approximate the non-linear model using an affine map  $\mathbf{M} : \mathbf{A}_L \rightarrow \mathbf{A}_{NL}$ , which maps from the linear model to the non-linear model, so that we set  $\mathbf{A} = \mathbf{M}\mathbf{A}_{NL} \approx \mathbf{A}_{NL}$ . Here, we give a brief introduction to affine maps and present our approach to calculating the affine map deterministically. Alternatively, one can also determine this map using other methods, e.g. machine learning methods or matrix inversion.

An affine map is any linear map between two vector spaces or affine spaces, where an affine space does not need to preserve a zero origin, see [39, Def. 2.3.1]. In other words, an affine map does not need to map to the origin of the associated vector space or be a linear map on vector spaces, including a 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 the books [39, 40]



**Figure 2.3:** This Figure shows the schematic representation of the affine map  $\mathbf{M}$ , which approximates the non-linear forward model from the linear forward model. Here,  $\mathbf{V}$  contains values produced by the linear forward model, and  $\mathbf{W}$  contains the corresponding values from the non-linear forward model. Both  $\mathbf{V}$  and  $\mathbf{W}$  are affine subspaces over the same field. The affine map  $\mathbf{M}$  projects elements from the linear forward model space  $\mathbf{V}$  onto their counterparts in the non-linear forward model space  $\mathbf{W}$ .

## 2.5 Regularisation

As mentioned in the introduction, the currently most used method to analyse data in atmospheric physics is regularisation-based. Since we want to show that our methods are computationally comparable if not faster, and provide more information than regularisation, we choose a linear-Gaussian Bayesian framework closest to our regulariser, see section ??.

The Tikhonov regularisation approach provides one solution  $\mathbf{x}_\lambda$  that minimises both the data misfit norm

$$\|\mathbf{y} - \mathbf{Ax}\|_{L^2} \quad (2.53)$$

and a regularisation semi-norm

$$\lambda \|\mathbf{T}\mathbf{x}\|_{L^2}, \quad (2.54)$$

for a given regularisation parameter  $\lambda > 0$  as described in [14], with a linear forward model matrix  $\mathbf{A}$ , the data  $\mathbf{y}$ , a regularisation operator  $\mathbf{T}$ . The regularisation parameter weights the semi-norm and penalises  $\mathbf{x}$  according to that. If  $\lambda$  is large, then the effect of the data on the solution  $\mathbf{x}_\lambda$  is small or negligible. If  $\lambda$  is small, the solution  $\mathbf{x}_\lambda$  will be dominated by the noisy data, resulting in an overfitted  $\mathbf{x}_\lambda$ . We refer to [41] and [7] for a more comprehensive analysis of the effects of the regularisation parameter on the solution, e.g. due to small singular values of the forward model.

For a fixed  $\lambda$ , the regularised solution

$$\mathbf{x}_\lambda = \arg \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{Ax}\|_{L^2}^2 + \lambda \|\mathbf{T}\mathbf{x}\|_{L^2}^2 \quad (2.55)$$

is obtained by taking the derivative with respect to  $\mathbf{x}$ :

$$\nabla_{\mathbf{x}} \left\{ (\mathbf{y} - \mathbf{Ax})^T (\mathbf{y} - \mathbf{Ax}) + \lambda \mathbf{x}^T \mathbf{T}^T \mathbf{T} \mathbf{x} \right\} = 0 \quad (2.56)$$

$$\iff \nabla_{\mathbf{x}} \left\{ \mathbf{y}^T \mathbf{y} + \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} - 2 \mathbf{y}^T \mathbf{Ax} + \lambda \mathbf{x}^T \mathbf{T}^T \mathbf{T} \mathbf{x} \right\} = 0 \quad (2.57)$$

$$\iff 2 \mathbf{A}^T \mathbf{Ax} - 2 \mathbf{A}^T \mathbf{y} + 2 \lambda \mathbf{T}^T \mathbf{T} \mathbf{x} = 0, \quad (2.58)$$

also known as the "regularised normal equations"  $\mathbf{A}^T \mathbf{y} = \mathbf{A}^T \mathbf{A} \mathbf{x} + \lambda \mathbf{T}^T \mathbf{T} \mathbf{x}$  [42]. Solving this equation yields the regularised solution

$$\mathbf{x}_\lambda = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{L})^{-1} \mathbf{A}^T \mathbf{y}, \quad (2.59)$$

where we define  $\mathbf{L} := \mathbf{T}^T \mathbf{T}$ , which typically represents a discrete matrix approximation of a differential operator choice [7]. For example

$$\mathbf{T} = \frac{1}{h} \begin{bmatrix} -1 & 1 & & & \\ 0 & -1 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & 0 & -1 & 1 \\ & & & 0 & -1 \end{bmatrix}, \quad (2.60)$$

is the first order derivative with equal spacing  $h$  as in [7] then

$$\mathbf{L} = \frac{1}{h^2} \begin{bmatrix} 1 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix}, \quad (2.61)$$

is the second order derivative with Neumann boundary conditions, see [43].

In practice,  $\mathbf{x}_\lambda$  is computed for a range of  $\lambda$ -values and evaluated based on the trade-off between the data misfit and the regularisation semi-norm. The optimal value of  $\lambda$  corresponds to the point of maximum curvature of the so-called L-curve [44], where the data misfit norm versus the regularisation semi-norm is plotted, see Fig. 4.13.

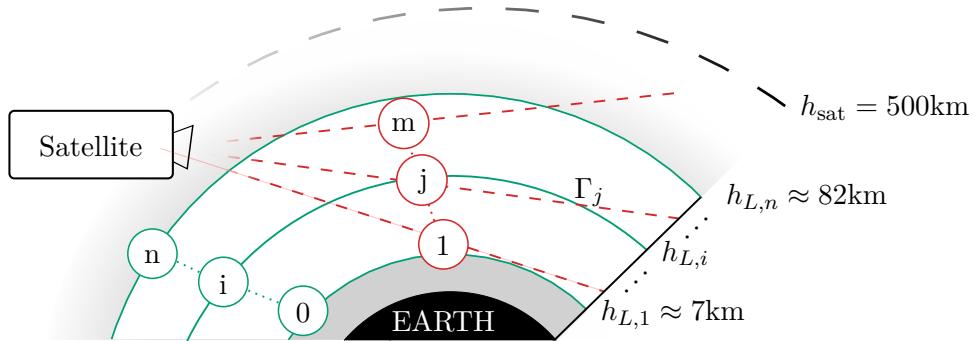
Additionally one can think about regularisation as a Lagrange multiplier  $\mathcal{L}(\mathbf{x}, \lambda) := \lambda \mathbf{x}^T \mathbf{L} \mathbf{x} + \|\mathbf{y} - \mathbf{A} \mathbf{x}\|_{L^2}$ , which minimises  $\mathbf{x}_\lambda = \arg \min_{\mathbf{x}} \mathbf{x}^T \mathbf{L} \mathbf{x}$  with respect to a constant  $\|\mathbf{y} - \mathbf{A} \mathbf{x}\|_{L^2}$ , see [14, fn. 6] and [45, Fig. 2.13]. So every solution  $\mathbf{x}_\lambda$  is an extremum (the most regularised solution for a given data misfit norm) and almost every sample of the posterior, which represents a feasible solution given the data, has a higher  $\mathbf{x}^T \mathbf{L} \mathbf{x}$  value and lays above the L-Curve.



# 3

## Forward Model

In this chapter, we present the forward model on which we apply all our methodology. We follow the Michelson Interferometer for Passive Atmospheric Sounding (MIPAS) handbook [46] and simulate data according to a cloud-free atmosphere in local thermodynamic equilibrium and assume a measurement instrument with infinite spectral resolution and no pointing errors. We do not include any other instrument specific details, such as sensor area or antenna response, as they are not available to us.



**Figure 3.1:** Schematic of measurement and analysis geometry, not to scale. The stationary satellite, at a constant height  $h_{\text{sat}}$  above Earth, takes  $m = 30$  measurements along its line-of-sight defined by the line  $\Gamma_j$ . Each measurement has a limb height  $\ell_j$ ,  $j = 1, 2, \dots, m$  defined as the closest distance of  $\Gamma_j$  to the Earth's surface. Between  $h_{L,0} \approx 7\text{km}$  and  $h_{L,n} \approx 82\text{km}$ , the stratosphere is discretised into  $n = 34$  layers as illustrated by the solid green lines.

A satellite at a constant height  $h_{\text{sat}}$  points through the atmosphere (limb-sounding) and measures thermal radiation of gas molecules along its straight line of sight  $\Gamma_j$ , see Figure 3.1. One measurement of the thermal radiation of one specific molecule, in our case ozone, denoted by the ozone volume mixing ratio (VMR)  $x(r)$  at distance  $r$  from

the satellite, at the wave number  $\nu$ , is given by the path integral

$$y_j = \int_{\Gamma_j} B(\nu, T) k(\nu, T) \frac{p(T)}{k_B T(r)} x(r) \tau(r) dr + \eta_j \quad (3.1)$$

$$\tau(r) = \exp \left\{ - \int_{r_{\text{obs}}}^r k(\nu, T) \frac{p(T)}{k_B T(r')} x(r') dr' \right\}, \quad (3.2)$$

which is the radiative transfer equation (RTE) [46]. For more information on the processes within the atmosphere for ozone, we refer to [47]. We define a tangent height  $h_{\ell_j}$  and  $\Gamma_j$  for each  $j = 1, 2, \dots, m$ , so that the data vector  $\mathbf{y} \in \mathbb{R}^m$  including some noise  $\eta_j$ . Within the atmosphere, 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  $k_B$ . 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 absorption constant

$$k(\nu, T) = L(\nu, T_{\text{ref}}) \frac{Q(T_{\text{ref}})}{Q(T)} \frac{\exp \{-c_2 E''/T\}}{\exp \{-c_2 E''/T_{\text{ref}}\}} \frac{1 - \exp \{-c_2 \nu/T\}}{1 - \exp \{-c_2 \nu/T_{\text{ref}}\}} \quad (3.3)$$

is dependent on the line intensity  $L(\nu, T_{\text{ref}})$  at reference temperature  $T_{\text{ref}} = 296K$ , the lower-state energy of the transition  $E''$ , the second radiation constant  $c_2 = 1.4387769\text{cmK}$  all provided by the HITRAN database [48]. Since we assume that the measurement device has negligible frequency window we neglect line broadening around  $\nu_0$  for the calculations of  $L(\nu, T_{\text{ref}})$ , which would normally be modelled as a convolution of the normalised Lorentz profile (collisional/pressure broadening) and the normalised Doppler (thermal broadening) profile [46]. Additionally, since we target one specific molecule, we simplify the calculation of  $k(\nu, T)$ , which usually involves summing the individual absorption constants for each targeted molecule weighted by the respective volume mixing ratio [46]. The total internal partition function for the lower-state energy is given as:

$$Q(T) = g'' \exp \left\{ -\frac{c_2 E''}{T} \right\}, \quad (3.4)$$

with the statistical weight  $g''$  (also called the degeneracy factor) accounting for the molecule's non-rotational and rotational energy states, see [49]. Under the assumption of local thermodynamic equilibrium (LTE), the black body radiation acts as a source function

$$B(\nu, T) = \frac{2hc^2\nu^3}{\exp \left\{ \frac{hc\nu}{k_B T} \right\} - 1}, \quad (3.5)$$

with Planck's constant  $h$  and speed of light  $c$ . For fundamentals on the Radiative transfer equation we recommend [50, Chapter 1], and for a more comprehensive model we refer to [51]

To enable matrix-vector multiplication, we discretise the atmosphere in  $n$  layers, where the  $i^{\text{th}}$  layer is defined by two spheres of radii  $h_{L,i-1} < h_{L,i}$ , for  $i = 1, \dots, n$ , with  $h_{L,0}$  and

$h_{L,n}$ . Then we can discretise the ozone, pressure and temperature profiles as a function of height; in between the heights  $h_{L,i-1}$  and  $h_{L,i}$ , each of the ozone concentration  $x_i$ , the pressure  $p_i$ , the temperature  $T_i$ , as well as all other height dependent parameters are assumed to be constant. Above  $h_{L,n}$  and below  $h_{L,0}$ , the ozone concentration is set to zero, so no signal can be obtained. We rewrite the integral in Eq. (3.1) for one noise free measurement, using the trapezoidal rule, as a vector-vector multiplication  $\mathbf{A}_j(\mathbf{x}, \mathbf{p}, \mathbf{T})$ , where the ozone volume mixing ratio  $\mathbf{x} = \{x_1, x_2, \dots, x_n\} \in \mathbb{R}^n$ . Note that since we aim to provide estimates over pressure  $\mathbf{p}$  and temperature  $\mathbf{T}$ , we explicitly include them as parameters in our forward model. Here, the non-linear absorption  $\tau(r)$  is another vector-vector multiplication and included in each entry of  $\mathbf{A}_j(\mathbf{x}, \mathbf{p}, \mathbf{T}) \in \mathbb{R}^n$  which is the  $j$ -th row of the matrix  $\mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T})$ . Then given a noise vector  $\boldsymbol{\eta} \in \mathbb{R}^m$  the data vector is given as:

$$\mathbf{y} = \mathbf{A}_{NL} + \boldsymbol{\eta}. \quad (3.6)$$

Here we define the non-linear forward model matrix as  $\mathbf{A}_{NL} := \mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T}) \in \mathbb{R}^{m \times n}$  for simplicity. Similarly we define  $\mathbf{A}_L$ , which denotes the linear forward model matrix and neglects absorption (e.g. set  $\tau = 1$  in Eq. (3.2)). Then we can compute noise-free linear data as matrix-vector multiplication  $\mathbf{A}_L \mathbf{x}$ . Further, we classify the inverse problem as weakly non-linear, see e.g. Fig. 4.9, as neglecting the absorption changes the measurements only slightly.

### 3.1 Singular Value Decomposition of the Forward Model

Before simulating some data, we provide a quick and intuitive way of assessing if the data collection is effective, how much information is passed through the forward model, depending on how we measure and how the signal to noise ratio (SNR) affects that information. One way of doing this is via a singular value decomposition (SVD) of the forward model matrix

$$\mathbf{A}_L = \sum_{i=1}^r \mathbf{u}_i \sigma_i \mathbf{v}_i^T = \mathbf{U} \Sigma \mathbf{V}^T \quad (3.7)$$

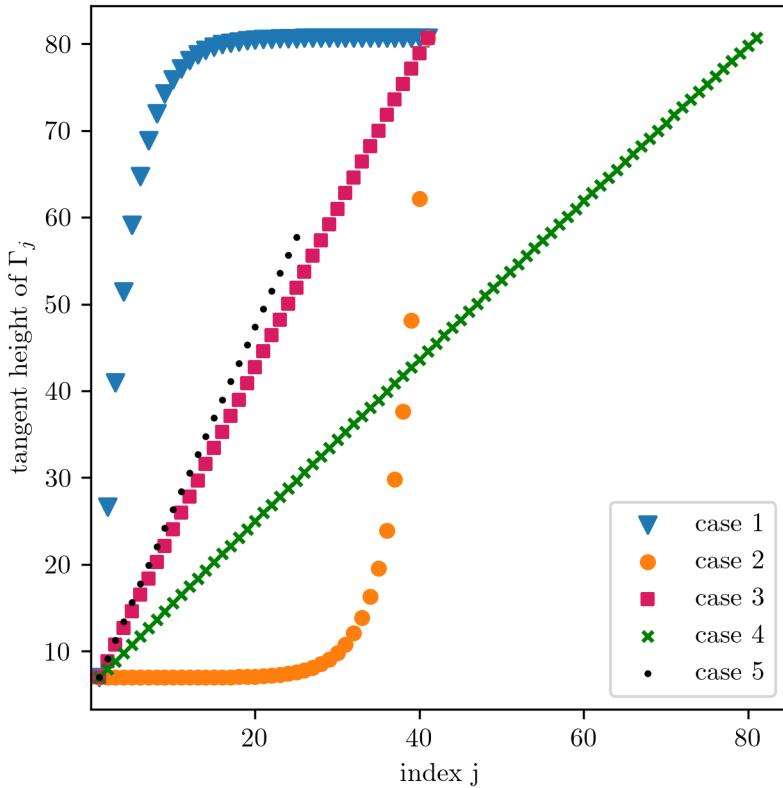
where  $r = \min\{m, n\}$  for a forward model  $\mathbf{A}_L \in \mathbb{R}^{m \times n}$ . Consider noise-free measurements  $\mathbf{A}_L \mathbf{x}$  for a satellite at a fixed height of  $h_{\text{sat}} = 500\text{km}$  above sea level, where  $\mathbf{x}$  is the ozone VMR. The SVD gives us information on how information is picked up from the parameter space by the forward model, described through the right singular values  $\mathbf{v}_i$ . The singular values  $\sigma_i$ , ordered in size from the largest  $\sigma_1$  to the smallest  $\sigma_r$ , weight that information from the right singular values to the left singular values  $\mathbf{u}_i$ , which project onto the data space. For a large singular value, we can say that the forward model is informative about structures in the corresponding right singular vector and vice versa. Note that we obtain

the same results using the non-linear forward model  $\mathbf{A}(\mathbf{x})$  to do this analysis, where we would rewrite to the matrix-vector multiplication  $\mathbf{A}_{NL}\mathbf{x}$ , where  $\mathbf{A}_{NL}$  is depend on  $\mathbf{x}$ .

Further, for very small singular values  $\sigma_i \ll \sigma_1/\text{SNR}$  below the RMS noise level or the noise standard deviation (STD), we can introduce an effective rank  $r_{\text{eff}} \leq r$ . Then information of parameter space spanned by  $\{\mathbf{v}_{r_{\text{eff}}+1}, \dots, \mathbf{v}_r\}$  is not passed through the forward model and the data is noise dominated in the corresponding data space, see Figure 3.6. This is based on the rough assumption that if we define the SNR as

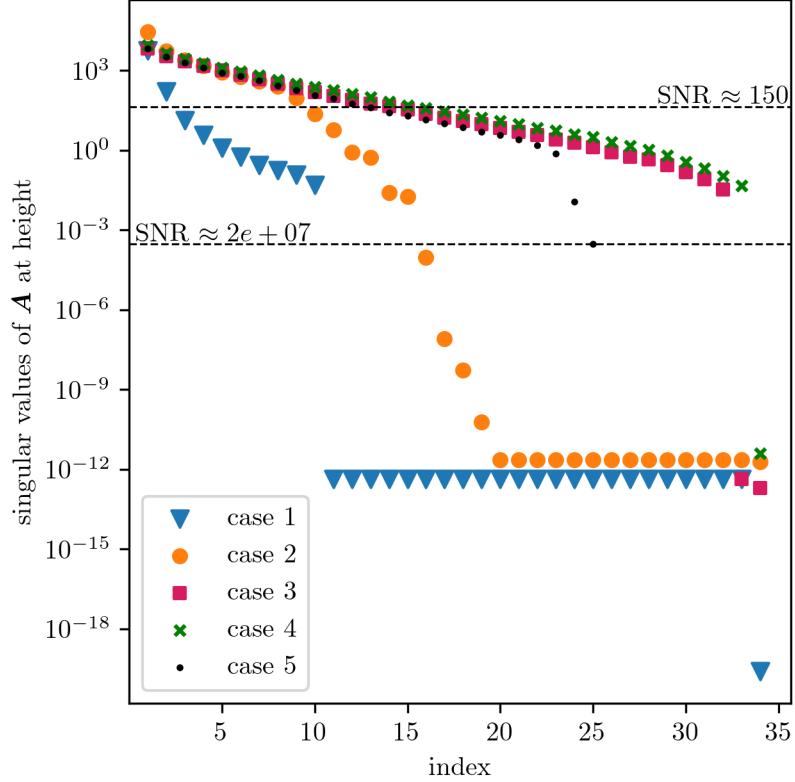
$$\text{SNR} := \frac{\max(y)}{\text{STD noise}} = \frac{\text{peak signal}}{\text{RMS noise}} [52] \quad (3.8)$$

then the maximum singular value  $\max(y) \approx \sigma_1$  and the information transmitted through the forward model corresponds roughly to the singular values  $\sigma_i \gtrsim \max(y)/\text{SNR}$ . See [7] for a more comprehensive analysis.



**Figure 3.2:** We plot the tangent heights for different cases of measurements.

Next, we plot the singular values for 5 different measurement scenarios, where we either measure at equidistance spaced tangent heights or collect more data from high signal regions at low altitudes, to see which of the tested cases is most effective. We assess the number of singular values above and below a certain SNR visually. Our objective is to measure ozone  $\mathbf{x}$  so our forward model  $\mathbf{A}$  includes temperature and pressure, the

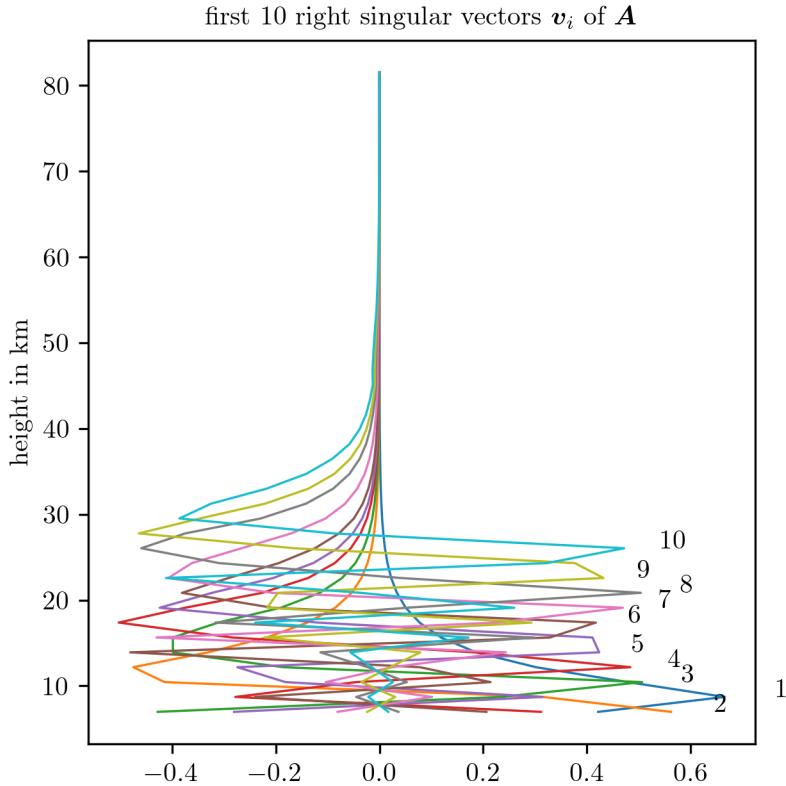


**Figure 3.3:** We plot the singular values of the forward model matrix for different sequences of measurements. The corresponding tangent heights of the test cases are plotted in Fig. 3.2. The dotted vertical line marks where the singular values are dominated by noise according to a SNR.

latter is dominant, see Fig. 4.17, and decreases exponentially in height and hence does affect the information passed through the model. If the pressure is high, the noise is low, and if the pressure is low, the data is noise-dominated. We start with case 3 in Fig. 3.2 where measurements are spaced according to a pointing accuracy of 150arc sec, given to us by the team of the University of New South Wales Canberra Space [53]. The pointing accuracy determines how well the satellite can point in a certain direction and, hence, roughly the spacing in between two measurements. The corresponding singular values are plotted in Fig. 3.3, of which the first 25 decrease linearly in log-space and about 10-15 singular values lie above the SNR. In comparison, if we measure a lot of times in regions where the data is noise-dominated (high altitude), case 1, we do obtain more information since the singular values decrease rapidly. Measuring lots of times at low altitudes, where the data is informative, and less at higher altitudes, case 2, does not seem optimal either, as we observe one larger singular value, but the other singular values decrease faster compared to case 3. Now, if we double the number of measurements compared to case 3, see case 4, we do get slightly larger singular values, but not significantly so that it would be worth the engineering effort required to achieve that. The measurements with

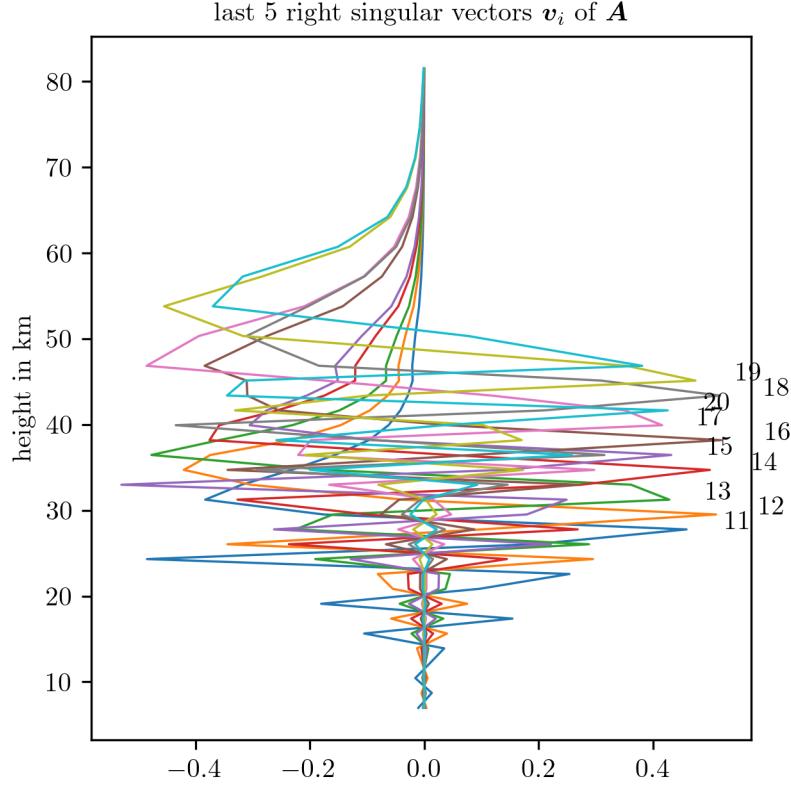
equidistance-spaced tangent height seem to be most informative. By exploratory analysis, we find that we can tolerate a slightly larger distance between tangent heights (pointing accuracy of 175arc sec) than required by [53], see case 5. In that case, we also stop measuring when the signal is too noisy and decrease the number of measurements taken without losing information. We note that if one wanted to obtain all information provided by the forward model, we would need a signal-to-noise ratio of roughly  $10^7$ .

In principle, we show that it does not depend on how one measures, one can not get more information by measuring more in regions where the information content is low or high. This contradicts the current measurement setup on the AURA MLS [8], which reports high noise in lower atmospheric regions, due to thermal radiation from the earth, and measures more in those regions.



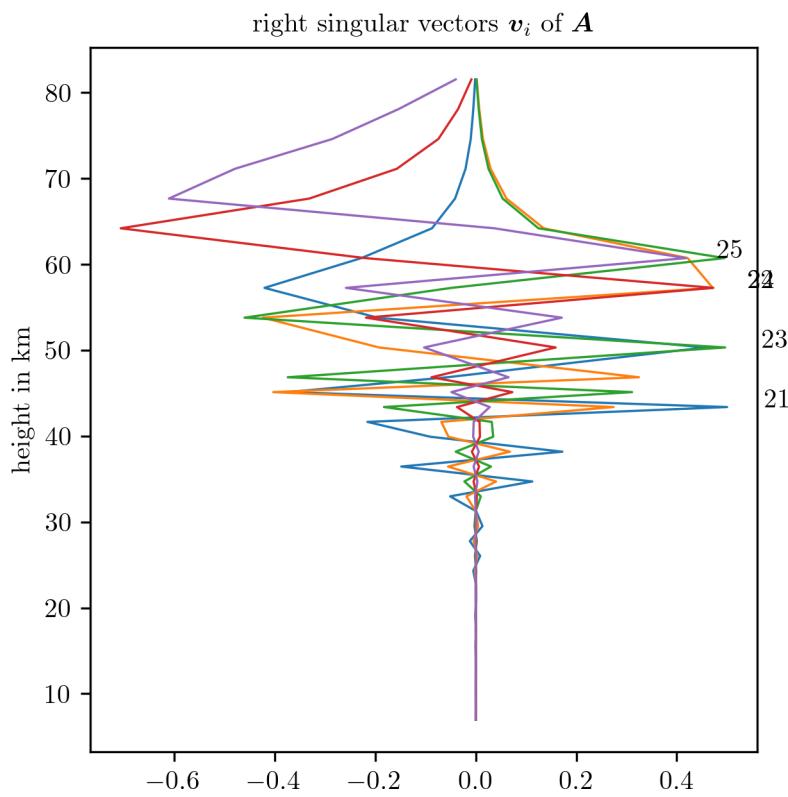
**Figure 3.4:** We plot the first 10 right singular vectors of the forward model matrix for case 5 sequence of measurements, see Fig. 3.3. These singular vectors correspond to high singular values of the forward model, see Fig. 3.3.

Consequently, we proceed with case 5 and plot the right singular vectors of the forward model versus height in the atmosphere to see where our model is most informative, or which structures of the parameter space are picked up by the model. The first 10 right singular vectors, in Fig. 3.4, corresponding to the 10 largest singular values, pick up structures in lower atmospheric regions. So we can assume that, given some data, we will



**Figure 3.5:** We plot the right singular vectors with index  $i = 11, \dots, 19$  of the forward model matrix for case 5 sequence of measurements, see Fig. 3.3. These singular vectors correspond to singular values around the noise level of the measurement, see Fig. 3.3.

be able to provide good reconstructions of the parameter in lower altitudes. Next, we plot the right singular vectors corresponding to the singular values  $\sigma_j$  for  $j = 11, \dots, 20$  in Fig. 3.5, where the noise starts to dominate the data. These singular values lie in regions around the SNR, see Fig. 3.3, and pick up values in the middle of the atmosphere. Consequently, we expect a higher uncertainty of reconstructed parameter values in the middle atmospheric regions. The singular vectors corresponding to the last 5 singular values pick up structures in higher altitudes, but since the singular values are very small, we will not be able to retrieve those structures. More specifically, the retrieved parameter values at higher altitudes will be fully determined by the prior or, in the case of regularisation, by the regulariser [7].



**Figure 3.6:** We plot the last 5 right singular vectors of the forward model matrix for the case 5 sequence of measurements, as displayed in Fig. 3.2. These singular vectors correspond to small singular values of the forward model, see Fig. 3.3.

# 4

## Results and Conclusions

In this chapter, we use the forward model to generate some data given an underlying ground truth and then guide the reader through the process of setting up a Bayesian framework and ultimately obtaining the posterior distributions of parameters of interest, such as ozone concentration or pressure and temperature profiles. We are using a directed acyclic graph (DAG) to visualise hierarchical and correlation structures of a Bayesian model. Next, we established a choice of prior distributions within our Bayesian model and formulated the posterior distributions. Based on the linear forward model  $\mathbf{A}_L$ , we characterise the marginal posterior for ozone and compare that to the TT approximation. Then we calculate the mean and the covariance matrix of the full posterior for ozone, which we use to find affine map to approximate the non-linear forward model  $\mathbf{A}_{NL} \approx \mathbf{M}\mathbf{A}_{NL}$  (see Sec. 4.3). In Sec. 4.4, we repeat the MTC scheme to provide a posterior distribution of ozone based on the approximate forward model and compare to a regularisation approach and a ground truth. Additionally, in sec 4.5, we extend the hierarchical Bayesian model and MTC scheme to jointly provide conditional posterior distributions of ozone, temperature and pressure. Here we elaborate on some aspects of prior modelling and on our findings when using a TT to approximate the higher dimensional marginal posterior. Lastly, we evaluate some errors occurring during the process. All of programming and analysis is done in Python and the reported computation times correspond to a MacBook Pro from 2019 with a 2.4 GHz quad-core Intel Core i5 processor.

## 4.1 Simulate Data based on a Ground Truth

We take a ground truth ozone VMR at distinct pressure values generated from some data [9] of the MLS on the Aura satellite within the Antarctic region and with a peak in high altitude, see Fig. 4.6, which seems to be a typical nighttime profile [47].

We target Ozone at a frequency of 235.71 GHz, which lies within the region where the MLS observes ozone [10, 54]. The corresponding wave number is  $\nu = 7.86\text{cm}^{-1}$ . We recursively calculate the geometric height with the hydrostatic equilibrium equation

$$\frac{dp}{p} = \frac{-gM}{R^*T} dh, \quad (4.1)$$

with the acceleration due to gravity

$$g = g_0 \left( \frac{r_0}{r_0 + h} \right), \quad (4.2)$$

where the polar radius pf the earth is  $r_0 \approx 6356\text{km}$ , the gravitation at sea level is  $g_0 \approx 9.81\text{m/s}^2$ ,  $R^* = 8.31432 \times 10^{-3}\text{Nm/kmol/K}$  and the mean molecular weight of the air is  $M = 28.97\text{kg/kmol}$  [55]. This holds up to a geometric height of 86km, where we ignore a 0.04% non-linear change in  $M$  from 80km to 86km in geometric altitude.

Following [55] we form the temperature function

$$T(h) = \begin{cases} T_0 & , h = 0 \\ T_0 + a_0 h & , 0 \leq h < h_{T,1} \\ T_0 + a_0 h_{T,1} & , h_{T,1} \leq h < h_{T,2} \\ T_0 + a_0 h_{T,1} + a_1(h_{T,2} - h_{T,1}) + a_2(h - h_{T,2}) & , h_{T,2} \leq h < h_{T,3} \\ T_0 + a_0 h_{T,1} + a_1(h_{T,2} - h_{T,1}) \\ \quad + a_2(h_{T,3} - h_{T,2}) + a_3(h - h_{T,3}) & , h_{T,3} \leq h < h_{T,4} \\ T_0 + a_0 h_{T,1} + a_1(h_{T,2} - h_{T,1}) \\ \quad + a_2(h_{T,3} - h_{T,2}) + a_3(h_{T,4} - h_{T,3}) + a_4(h - h_{T,4}) & , h_{T,4} \leq h < h_{T,5} \\ T_0 + a_0 h_{T,1} + a_1(h_{T,2} - h_{T,1}) \\ \quad + a_2(h_{T,3} - h_{T,2}) + a_3(h_{T,4} - h_{T,3}) + a_4(h_{T,5} - h_{T,4}) \\ \quad + a_5(h - h_{T,5}) & , h_{T,5} \leq h < h_{T,6} \\ T_0 + a_0 h_{T,1} + a_1(h_{T,2} - h_{T,1}) \\ \quad + a_2(h_{T,3} - h_{T,2}) + a_3(h_{T,4} - h_{T,3}) + a_4(h_{T,5} - h_{T,4}) \\ \quad + a_5(h_{T,6} - h_{T,5}) + a_6(h - h_{T,6}) & , h_{T,6} \leq h \lesssim 86 \end{cases} \quad (4.3)$$

with gradient and height values provided by [55], see Tab. 4.1. This acts as the ground truth temperature profile, see Fig. 4.15.

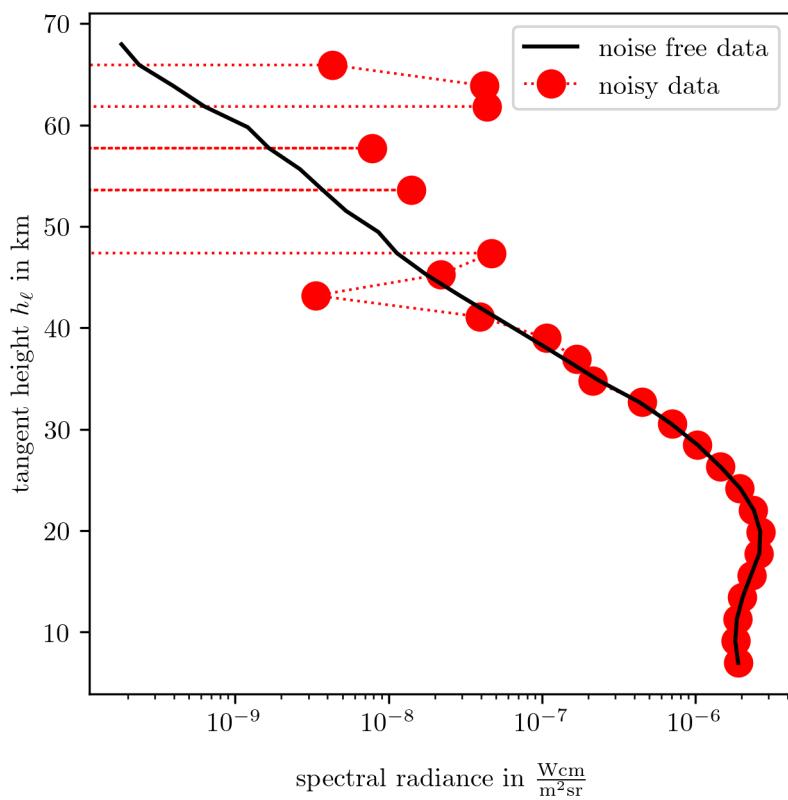
subscript $i$	geometric height $h_{T,i}$ in km	gradient $a_i$
0	0	-6.5
1	11	0
2	20.1	1
3	32.2	2.8
4	47.4	0
5	51.4	-2.8
6	71.8	-2

**Table 4.1:** Definition of height depending temperature gradients.

Then we can compute a data vector  $\mathbf{y} = \mathbf{A}_{NL} + \boldsymbol{\eta}$ , with  $m = 30$  measurements according to the radiative transfer equation (RTE), see Eq. 3.1 which we solve using the trapezoidal integration rule, determined by the satellite pointing accuracy of 175arc sec, see Fig. 3.2. We assume an atmosphere between  $h_{L,1} = 7\text{km}$  and  $h_{L,n} = 83.3\text{km}$  with  $n = 45$  equidistant layers. The height value  $h_{L,i}$  for each layer  $i = 1, \dots, n$  is defined by the pressure values from [9] and the hydrostatic equilibrium equation, see Eq. 4.1. We calculate the absorption constant  $k(\nu, T)$  as in Eq. 3.2, following the HITRAN database [48], which provides the line intensity  $L(\nu, T_{\text{ref}})$  for the isotopologue  $^{16}\text{O}_3$  with the AFGL Code 666. This gives us a non-linear forward model matrix  $\mathbf{A}_{NL} = \mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T}) \in \mathbb{R}^{m \times n}$ , where  $\mathbf{x} \in \mathbb{R}^n$  is vector related to the ozone VMR,  $\mathbf{p} \in \mathbb{R}^n$  is the vector describing the pressure values and  $\mathbf{T} \in \mathbb{R}^n$  the temperature values.

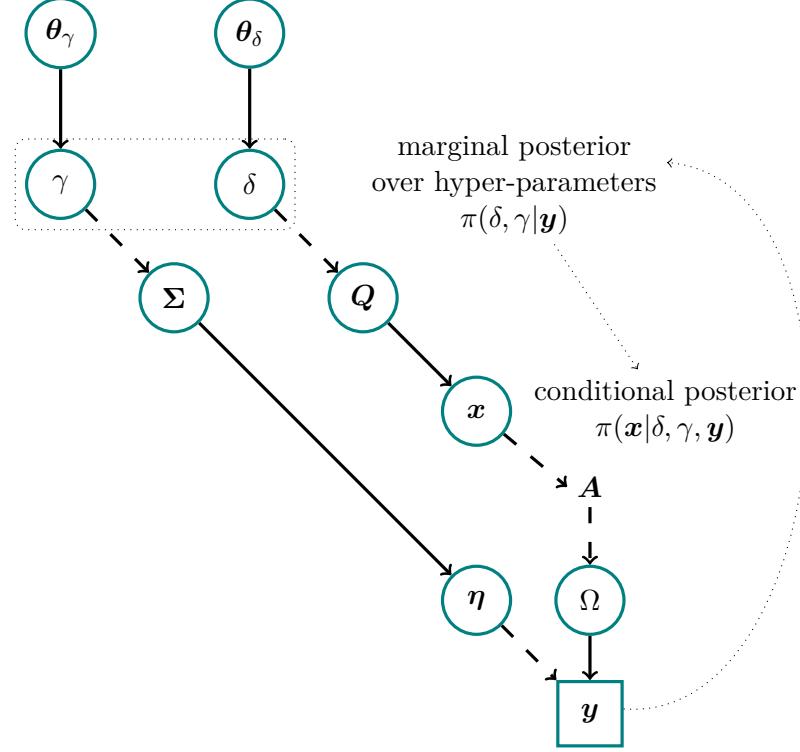
Lastly we add normally distributed  $\boldsymbol{\eta} \sim \mathcal{N}(0, \gamma^{-1} \mathbf{I})$  noise so that the SNR is 150, see Eq. 3.8, similar to [11], where a signal with a maximal spectral intensity of around 100K and a noise range of 0.4 to 1.6K is reported. We note that the methods used in this thesis will work with different SNRs or other frequencies.

When we plot the data in Fig. 4.1, we see that, as mentioned in Section 3.1, the data is noise-dominated in higher altitudes. Now, given the data, we like to determine the posterior distributions over ozone  $\mathbf{x}$ , pressure  $\mathbf{p}$  and temperature  $\mathbf{T}$  at different heights.



**Figure 4.1:** Logarithmic plot of data points at different tangent height. Note that negative values are not appearing, and we see that the noise is dominating at high altitudes.

## 4.2 Hierarchical Bayesian Framework for Ozone



**Figure 4.2:** DAG for visualisation of hierarchical modelling and measuring process of ozone including the MTC scheme. The hyper-parameter  $\gamma$  deterministically (dotted line) sets the noise covariance  $\Sigma = \gamma^{-1} \mathbf{I}$  and hence the random (solid line) noise vector  $\eta \sim \mathcal{N}(0, \gamma^{-1} \mathbf{I})$ . The hyper-parameter  $\delta$  determines (dotted line) the prior precision matrix  $Q = \delta \mathbf{L}$  for the normally distributed (solid line) prior  $x|\delta \sim \mathcal{N}(0, \delta \mathbf{L})$ , where  $\mathbf{L}$  is a graph Laplacian, see Eq. 4.5. The hyper-prior distributions (solid line)  $\pi(\delta, \gamma)$  are defined by  $\theta_\gamma$  and  $\theta_\delta$ . Through the linear forward model  $\mathbf{A}$  we generate a space of all measurable noise free data  $\mathbf{Ax}$  from which we randomly observe a data set  $\mathbf{y}$  including some added noise  $\eta$ . Within the MTC scheme we evaluate the marginal posterior over the hyper-parameters  $\pi(\gamma, \delta | \mathbf{y})$  first and then the conditional posterior  $\pi(x | \delta, \gamma, \mathbf{y})$ . This breaks the correlation structure of  $x$  and  $\delta$  and  $\gamma$ , and allows to evaluate the marginal posterior independent of  $x$ .

In this section, we setup the hierarchically-ordered linear-Gaussian Bayesian framework to determine the ozone posterior distribution, conditioned on ground truth temperature and pressure. Where for now we define the forward model matrix  $\mathbf{A} = \mathbf{A}_L$  and define the distributions of that Bayesian model, similarly to a regularisation approach, as:

$$\mathbf{y} | \mathbf{x}, \gamma, \delta \sim \mathcal{N}(\mathbf{A}\mathbf{x}, \gamma^{-1} \mathbf{I}) \quad (4.4a)$$

$$\mathbf{x} | \delta \sim \mathcal{N}(\mathbf{0}, (\delta \mathbf{L})^{-1}) \quad (4.4b)$$

$$\delta \sim \Gamma(\alpha_\delta = 1, \beta_\delta = 10^{-35}) \quad (4.4c)$$

$$\gamma \sim \Gamma(\alpha_\gamma = 1, \beta_\gamma = 10^{-35}). \quad (4.4d)$$

Assuming Gaussian noise  $\eta \sim \mathcal{N}(0, \gamma^{-1} \mathbf{I})$ , the likelihood function is a normal distributions with mean  $\mathbf{Ax}$  and covariance matrix  $\gamma^{-1} \mathbf{I}$ . We define the normal prior-distribution

$\pi(\mathbf{x}|\delta)$ , with zero mean and precision matrix  $\delta\mathbf{L}$ , where  $\delta$  is a smoothness hyper-parameter and  $\mathbf{L}$  is the second order discrete derivate operator (see Eq. 4.5). Here the hyper-prior distributions  $\pi(\delta)$  and  $\pi(\gamma)$  are gamma distributions with shape  $\alpha$  and rate  $\beta$ .

We can visualise this hierarchical structure and the correlations in between different hyper-parameters and parameters through a DAG. The hyper-parameter  $\gamma$  sets the noise covariance deterministically (dotted line), but is it self statistically (solid line) defined by the hyper-prior distribution  $\pi(\gamma)$ . In this case  $\theta_\gamma$  determines the hyper-prior distribution  $\pi(\gamma)$ , and similarly  $\theta_\delta$  for  $\pi(\delta)$ , which then deterministically sets the prior precision  $\mathbf{Q}$ . In our case,  $\delta$  accounts for smoothness of the ozone profile. Then  $\mathbf{Ax}$  determines the space of all measurable noise-free data sets  $\Omega$ , through the linear forward model, from which we observe a data set  $\mathbf{y}$  including some noise  $\eta$ . From that data we then "reverse the arrows" to determine the posterior distribution over the parameter  $\mathbf{x}$ . Since noise is a random process with a defined distribution, the posterior distribution  $\pi(\mathbf{x}|\mathbf{y})$  is well defined. Usually, due to underlying correlation structures, evaluating this posterior poses a significant challenge, here the MTC scheme provides the marginal posterior  $\pi(\delta, \gamma|\mathbf{y})$  first and then the conditional posterior  $\pi(\mathbf{x}|\delta, \gamma, \mathbf{y})$ .

### 4.2.1 Prior Modelling

To complete the Bayesian framework, we have to define prior distributions over the hyper-parameters and parameters. Ideally, we define the prior distributions as uninformative as possible, and include functional dependencies and physical properties.

By choosing a normally distributed prior  $\pi(\mathbf{x}|\delta)$  with zero mean and no other restrictions, we can already see that our model is not taking into account that ozone values can not be negative. As already mentioned, we set the precision matrix of that prior distribution to

$$\delta\mathbf{L} = \delta \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{bmatrix} \quad (4.5)$$

which is a 1-dimensional Graph Laplacian as in [14, 43] with Dirichlet boundary condition. This matrix will also act as the regulariser later in the Regularisation section, see Sec. 2.5. We reduce the dimension of  $\mathbf{x}$  from 45 to 34 by discarding every second ozone VMR over a height of  $\approx 47\text{km}$ . Doing that, while not changing  $\mathbf{L}$ , we effectively induce a larger correlation between points at higher altitude. We plot the corresponding prior ozone profiles according to  $\mathbf{x} \sim \mathcal{N}(0, (\delta\mathbf{L})^{-1})$  in Fig. B.1.

For  $\delta$  and  $\gamma$  we pick relatively uninformative gamma distributions so that  $\gamma \sim \mathcal{T}(\boldsymbol{\theta}_\gamma) \propto \gamma^{\alpha_\gamma - 1} \exp(-\beta_\gamma \gamma)$  and  $\delta \sim \mathcal{T}(\boldsymbol{\theta}_\delta)$ , where  $\boldsymbol{\theta}_\gamma = \{\alpha_\gamma, \beta_\gamma\} = \{\alpha_\delta, \beta_\delta\} = \boldsymbol{\theta}_\delta = (1, 10^{-35})$ , see Fig. 4.10, similar to [14]. Those gamma distributions have another advantage when using a MWG algorithm to sample from the marginal posterior distribution  $\pi(\delta, \gamma | \mathbf{y})$ . In doing so, we introduce the regularisation parameter  $\lambda = \delta/\gamma$  so that  $\pi(\gamma | \lambda, \mathbf{y}) \sim \mathcal{T}(\cdot)$  is a gamma distribution and easy to sample from.

### 4.2.2 Posterior Distribution – linear Model

As explained in Sec. 2.1 we factorise the posterior

$$\pi(\mathbf{x}, \delta, \gamma | \mathbf{y}) \propto \pi(\mathbf{y} | \mathbf{x}, \delta, \gamma) \pi(\mathbf{x}, \delta, \gamma) \quad (4.6)$$

into

$$\pi(\mathbf{x}, \delta, \gamma | \mathbf{y}) = \pi(\mathbf{x} | \delta, \gamma, \mathbf{y}) \pi(\delta, \gamma | \mathbf{y}) \quad (4.7)$$

the marginal posterior  $\pi(\delta, \gamma | \mathbf{y})$  and conditional posterior  $\pi(\mathbf{x} | \delta, \gamma, \mathbf{y})$  (see Eq. 2.5). As discussed in sec. 2.1, for the linear-Gaussian case,  $\mathbf{x}$  cancels in the marginal posterior over the hyper-parameters. Following the MTC scheme, we characterise the marginal posterior first and *then* the conditional posterior of  $\pi(\mathbf{x} | \delta, \gamma)$ .

#### Marginal Posterior

Consequently, for the hierarchical model specified in Eq. 4.4, the marginal posterior distribution over the hyper-parameters is given by

$$\pi(\lambda, \gamma | \mathbf{y}) \propto \lambda^{n/2 + \alpha_\delta - 1} \gamma^{m/2 + \alpha_\delta + \alpha_\gamma - 1} \exp\left\{-\frac{1}{2}g(\lambda) - \frac{\gamma}{2}f(\lambda) - \beta_\delta \lambda \gamma - \beta_\gamma \gamma\right\}, \quad (4.8)$$

with the introduced regularisation parameter  $\lambda = \delta/\gamma$ , and

$$f(\lambda) = \mathbf{y}^T \mathbf{y} - (\mathbf{A}^T \mathbf{y})^T (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{L})^{-1} (\mathbf{A}^T \mathbf{y}), \quad (4.9a)$$

$$\text{and } g(\lambda) = \log \det(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{L}). \quad (4.9b)$$

Note that, when changing variables from  $\delta = \lambda \gamma$  to  $\lambda$  the hyper-prior distribution changes to  $\pi(\lambda) \propto \lambda^{\alpha_\delta - 1} \gamma^{\alpha_\delta} \exp(-\beta_\delta \lambda \gamma)$ , due to  $d\delta/d\lambda = \gamma$ . Most of the computational effort, for each function evaluation of the marginal posterior in Eq. 4.8, lies in the calculation of  $f(\lambda)$  in Eq. 4.30a and  $g(\lambda)$  in Eq. 4.30b. In Fig. 4.3 we see that  $f(\lambda)$  and  $g(\lambda)$  are well behaved within the region of interest and approximate  $f(\lambda) \approx \tilde{f}(\lambda)$  with a 3rd order Taylor series around the mode  $\lambda_0$  of  $\pi(\lambda, \gamma | \mathbf{y})$ . We also note that  $\tilde{g}(\lambda) \approx g(\lambda)$  behaves linearly around  $\lambda_0$  in the log-space. The approximations are implicitly given by

$$f^{(r)}(\lambda_0) = (-1)^{r+1} (\mathbf{A}^T \mathbf{y})^T (\mathbf{B}_0^{-1} \mathbf{L})^r \mathbf{B}_0^{-1} \mathbf{A}_L^T \mathbf{y} \quad (4.10)$$

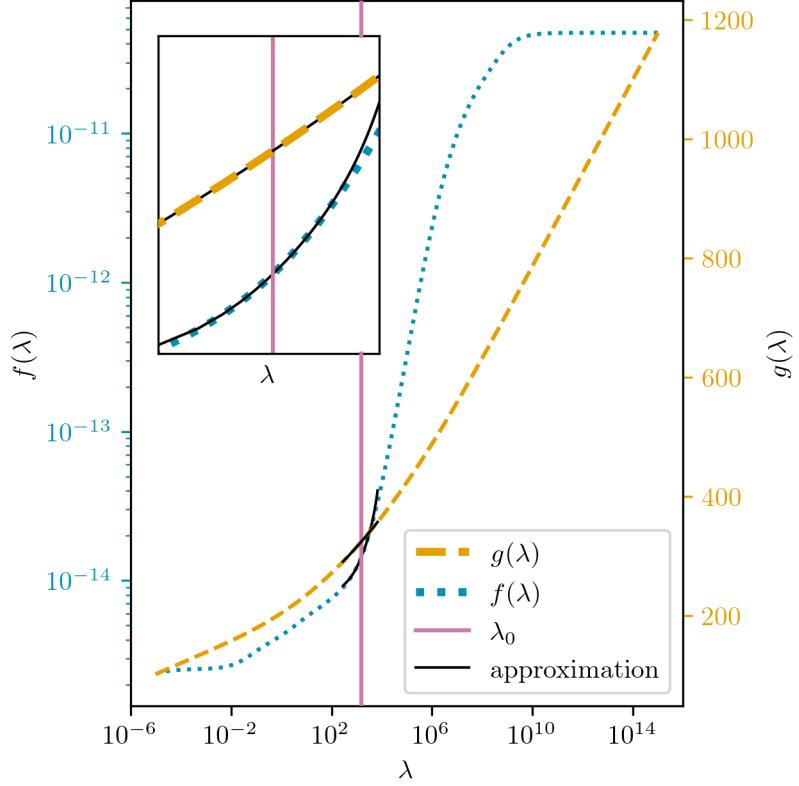
$$\text{and } \log \tilde{g}(\lambda) = (\log \lambda - \log \lambda_0) \frac{\log g(\lambda_{\max}) - \log g(\lambda_0)}{\log \lambda_{\max} - \log \lambda_0} + \log g(\lambda_0) \quad (4.11)$$

with  $\mathbf{B}_0 = \mathbf{A}^T \mathbf{A} + \lambda_0 \mathbf{L}$ . We plot the approximations

$$\tilde{f}(\lambda) = \sum_{r=0}^3 f^{(r)}(\lambda_0) (\lambda - \lambda_0)^r, \quad (4.12a)$$

$$\text{and } \tilde{g}(\lambda) = \exp \log \tilde{g}(\lambda), \quad (4.12b)$$

in Fig. 4.3 and elaborate on the approximation errors in Sec 4.2.2. Note that usually a Taylor series includes a factor  $(r!)^{-1}$ , in this case it cancels in  $f^{(r)}(\lambda_0)$ , see [14].



**Figure 4.3:** Plot of the functions  $f(\lambda)$  and  $g(\lambda)$  from the marginal posterior for a wide range of  $\lambda = \delta/\gamma$ . We plot the third-order Taylor series in black around the mode of the marginal posterior (vertical line) for the sampling range of  $\lambda$  within the MWG algorithm.

**Error due to Approximation of  $f$  and  $g$**  We report an maximum approximation error in between  $f$  and  $g$  For 2nd order taylor 3rd tayle and 4th order taylor

then we neglect ther erro in  $g$  beacus ethabsolute in  $f$  and  $g$  are and propagation error in marginal is due to  $f$  is realtive RMS and the maximum error is ...

When approximating the functions  $f(\lambda)$  and  $g(\lambda)$ , we find that the 3rd-order Taylor series of  $f(\lambda)$  and a linear approximation of  $g(\lambda)$  in log-space give the smallest error. The Taylor series truncation error of  $f(\lambda)$  is bounded by the fourth order Taylor series  $E_f = \arg \max_{\lambda} f^{(4)}(\lambda_0)/4! (\lambda - \lambda_0)^4$  and corresponds to an relative error bounded by 20%. Since the maximum absolute error of the approximation  $\arg \max_{\lambda} |\tilde{g}(\lambda) - g(\lambda)| \approx 1$  corresponds to an relative error of approximately 0.3% and is small compared to  $E_f \approx 1e8$  we ignore the approximation error of  $g(\lambda)$ . Then the maximum relative propagation error  $\arg \max_{\lambda, \gamma} 0.5\gamma E_f / \log \pi(\lambda, \gamma | \mathbf{y})$  is bound by approximately 5%.

**Sample from Marginal Posterior** Using these approximation we can either utilise a TT approximation of the marginal posterior, see Sec. 2.3.1, over a predefined grid and calculate the marginals  $\pi(\gamma | \mathbf{y})$  and  $\pi(\lambda | \mathbf{y})$ , or employ a Metropolis within Gibbs (MWG)

sampler to sample from  $\pi(\lambda, \gamma | \mathbf{y})$ , see sec. 2.2.1. More specifically, we implement a Metropolis random walk on

$$\pi(\lambda | \gamma, \mathbf{y}) \propto \lambda^{n/2 + \alpha_\delta - 1} \exp \left\{ -\frac{1}{2} g(\lambda) - \frac{\gamma}{2} f(\lambda) - \beta_\delta \gamma \lambda \right\}. \quad (4.13)$$

We accept or reject a proposal  $\lambda' \sim \mathcal{N}(0, \sigma_\lambda)$  according to the acceptance ratio in log space

$$\log \left\{ \frac{\pi(\lambda' | \gamma^{(k)}, \mathbf{y})}{\pi(\lambda^{(k)} | \gamma^{(k)}, \mathbf{y})} \right\} = \log \{ \pi(\lambda' | \gamma^{(k)}, \mathbf{y}) \} - \log \{ \pi(\lambda^{(k)} | \gamma^{(k)}, \mathbf{y}) \} \quad (4.14)$$

$$= \frac{n}{2} (\log \{ \lambda' \} - \log \{ \lambda^{(k)} \}) + \frac{1}{2} \Delta g + \frac{\gamma^{(k)}}{2} \Delta f + \beta_\delta \gamma^{(k)} \Delta \lambda, \quad (4.15)$$

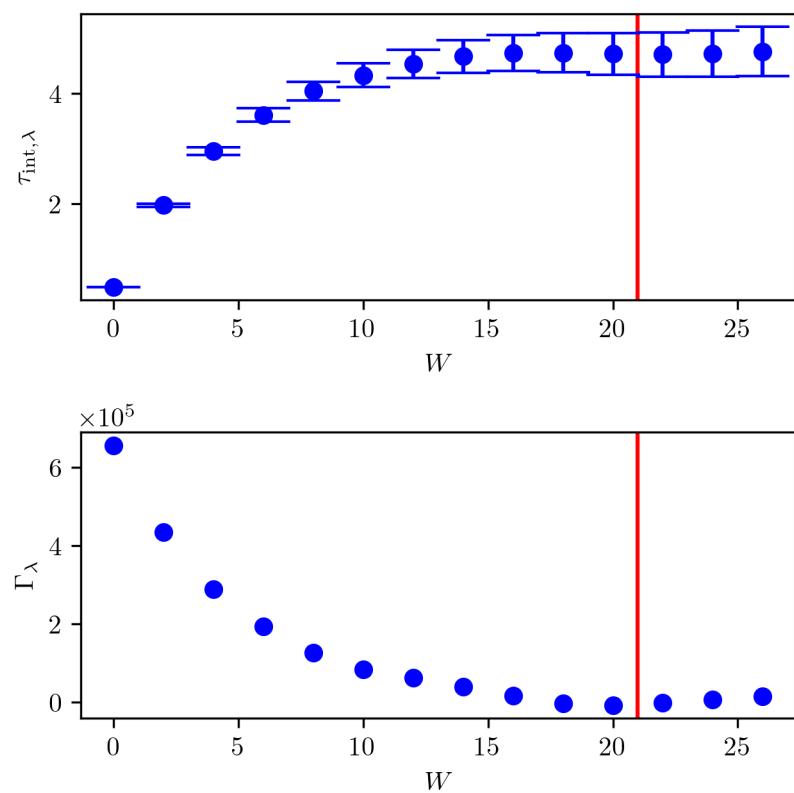
where  $\Delta \lambda = \lambda' - \lambda^{(k)}$  and  $\Delta f \approx \tilde{f}(\lambda') - \tilde{f}(\lambda^{(k)}) = \sum_{r=1}^3 f^{(r)}(\lambda_0) (\Delta \lambda' - \Delta \lambda^{(k)})^r$ , with  $\Delta \lambda' = \lambda' - \lambda_0$  and  $\Delta \lambda^{(k)} = \lambda^{(k)} - \lambda_0$ . Similarly we approximate  $\Delta g \approx \tilde{g}(\lambda') - \tilde{g}(\lambda^{(k)})$ .

Lastly, we do a Gibbs step on

$$\gamma^{(k+1)} | \lambda^{(k+1)}, \mathbf{y} \sim \Gamma \left( \frac{m}{2} + \alpha_\delta + \alpha_\gamma, \frac{1}{2} f(\lambda^{(k+1)}) + \beta_\gamma + \beta_\delta \lambda^{(k+1)} \right) \quad (4.16)$$

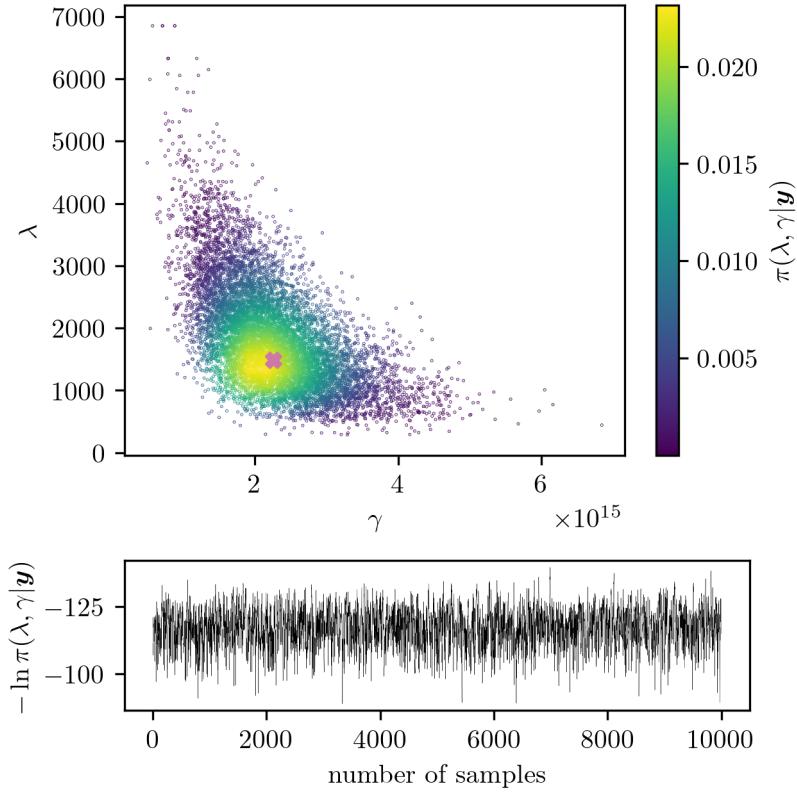
to generate marginal posterior samples  $(\lambda, \gamma)^{(1)}, \dots, (\lambda, \gamma)^{(N)} \sim \pi(\lambda, \gamma | \mathbf{y})$ .

The mode  $(\lambda_0, \gamma_0)$  of  $\pi(\lambda, \gamma | \mathbf{y})$  is provided by the `scipy.optimize.fmin` function and we approximate  $f(\lambda)$  and  $g(\lambda)$  accordingly. In doing so, we compute the vector  $\mathbf{B}_0^{-1} \mathbf{A}^T \mathbf{y} = (\mathbf{A}^T \mathbf{A} + \lambda_0 \mathbf{L}) \mathbf{A}^T \mathbf{y}$ , the matrix  $\mathbf{B}_0^{-1} \mathbf{L}$  and the determinant in  $g(\lambda)$  using Cholesky decomposition. Then we approximate  $f(\lambda)$  with a 3rd order Taylor series and  $g(\lambda)$  with a linear approximation in the log-space, where for the approximation in  $g(\lambda)$  we set  $\lambda_{\max}$  to the maximum value of  $\lambda$  on the TT-grid (see next Paragraph). To sample from  $\pi(\lambda, \gamma | \mathbf{y})$  we employ the MWG algorithm, see Alg. Box 1, initialised at the mode  $(\lambda^{(0)}, \gamma^{(0)}) = (\lambda_0, \gamma_0)$  and take  $N = 10000$  plus  $N_{\text{burn-in}} = 100$  steps in approximately 0.3s. The standard deviation of the normal proposal distribution is set to  $\sigma_\lambda = 0.8 \lambda_0$  so that the acceptance rate is  $\approx 0.5$  as suggested in [56]. The samples are plotted in Fig. 4.5 as a 2D scatter plot, as well as the trace of the MWG to show ergodicity. We calculate the integrated autocorrelation time (IACT) with the Python implementation of [23], provided by [57], which gives us  $\tau_{\text{int}, \gamma} =$  and  $\tau_{\text{int}, \delta} =$ , see Fig. 4.4 and Fig. B.2.



**Figure 4.4:** Here the autocorrelation function  $\Gamma_\lambda$  at different lags  $W$  is plotted as well as the IATC  $\tau_{\text{int},\lambda}$  for the samples from  $\pi(\gamma, \lambda | \mathbf{y})$  based on the linear forward model.

**TT Approximation of Marginal Posterior** We approximate the square root of marginal posterior on a predefined univariate grid, where  $\gamma = [0.25 \times 10^{15}, 5.5 \times 10^{15}]$  and  $\lambda = [100, 5000]$ . We set the number of grid points to  $n = 20$  and the number of ranks  $r = 5$ , which we keep constant. Since we do not approximate  $\sqrt{\pi(\lambda, \gamma|\mathbf{y})}$  in the log-space we introduce a "normalisation constant"  $c = 340$ . This avoids underflow so that the values  $\sqrt{\pi(\lambda, \gamma|\mathbf{y})} = \exp\{0.5 \log \pi(\lambda, \gamma|\mathbf{y}) + c\}$  are within computer precision. Then we initialise the `tt.cross.rectcross.rect_cross.cross` function, based on the TT cross algorithm in [58, 59], from the Python package `ttypy` [60], with a random tensor and do 1 sweep to obtain a TT approximation of  $\pi(\lambda, \gamma|\mathbf{y})$ . This takes about 0.1s for 400 function evaluations. Then we compute the marginals  $\pi(\lambda|\mathbf{y})$  and  $\pi(\gamma|\mathbf{y})$  as described in Sec. 2.3.1. In doing so we calculate the coefficient tensor  $\mathbf{B}$  and  $\mathbf{B}_{\text{pre}}$  as in Prop. 1 and 2, where we set  $\xi = 1/\lambda(\mathcal{X})$  and  $\lambda(x) = 1$ , so that for Cartesian basis  $\mathbf{M}_k = \text{diag}(\lambda_k(\mathcal{X}_k))$  in Eq. 12. We plot the TT approximation as a colour map on top of the obtained samples in the scatter plot in Fig. 4.5. We report a RMS error over the whole grid of ...



**Figure 4.5:** We scatter plot the samples of  $\lambda = \delta/\gamma$  and  $\gamma$  from the marginal posterior  $\pi(\lambda, \gamma|\mathbf{y})$  and colour code the samples using the TT approximation of  $\pi(\lambda, \gamma|\mathbf{y})$ . The mode of  $(\lambda_0, \gamma_0)$  of  $\pi(\lambda, \gamma|\mathbf{y})$  is marked by the pink cross. To show ergodicity we plot the trace of the samples of the MWG sampler below.

## Full Posterior Ozone Mean and Variance

Then we evaluate the normally distributed conditional posterior distribution

$$\mathbf{x}|\delta, \gamma, \mathbf{y} \sim \mathcal{N}\left(\underbrace{\mathbf{A}^T \mathbf{A} + \delta/\gamma \mathbf{L}}_{\mathbf{x}_\lambda}^{-1} \mathbf{A}^T \mathbf{y}, \underbrace{(\gamma \mathbf{A}^T \mathbf{A} + \delta \mathbf{L})^{-1}}_{\gamma \mathbf{B}_\lambda}\right), \quad (4.17)$$

as in Eq. 2.15, with  $\lambda = \delta/\gamma$ . In this thesis, we compute the mean

$$\mu_{\mathbf{x}|\mathbf{y}} = \int \mathbf{x}_\lambda \pi(\lambda|\mathbf{y}) d\lambda \approx \sum \mathbf{x}_{\lambda_i} \pi(\lambda_i|\mathbf{y}), \quad (4.18)$$

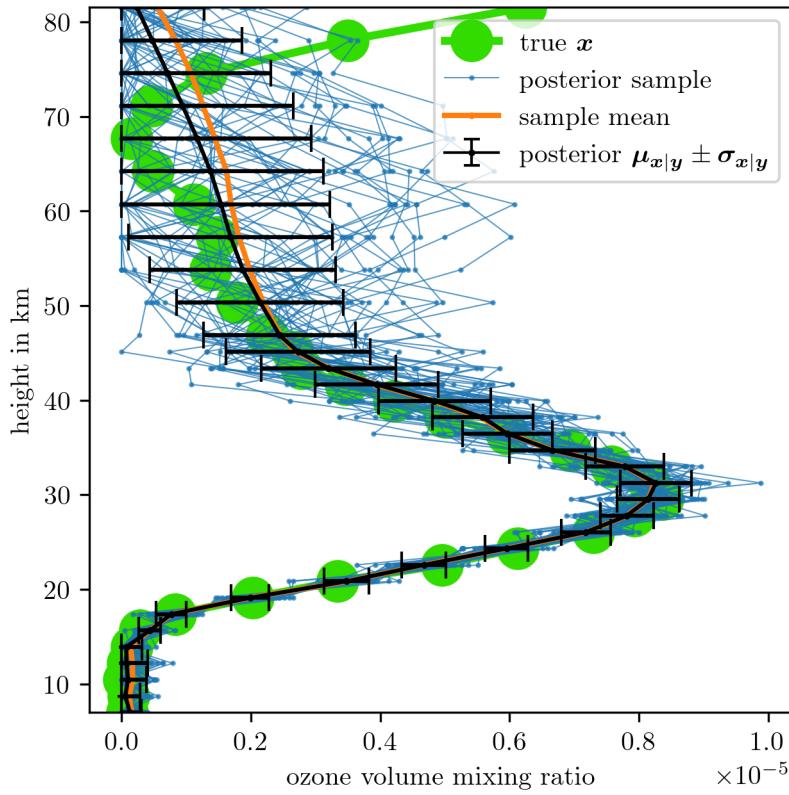
and covariance

$$\Sigma_{\mathbf{x}|\mathbf{y}} = \int \gamma^{-1} \pi(\gamma|\mathbf{y}) d\gamma \int \mathbf{B}_\lambda^{-1} \pi(\lambda|\mathbf{y}) d\lambda \approx \sum \gamma_i^{-1} \pi(\gamma_i|\mathbf{y}) \sum \mathbf{B}_{\lambda_i}^{-1} \pi(\lambda_i|\mathbf{y}) \quad (4.19)$$

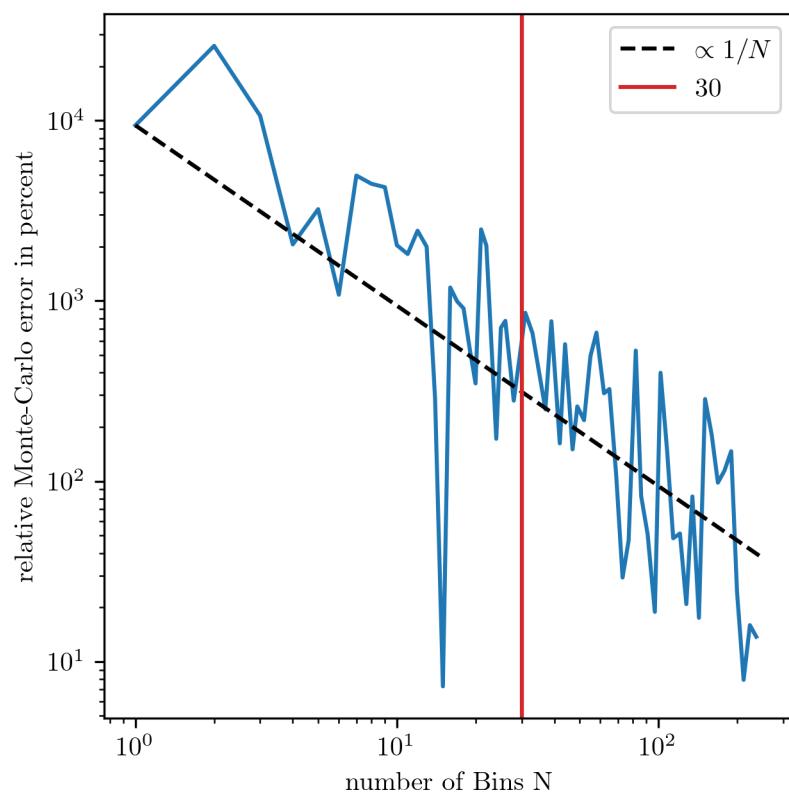
of  $\pi(\mathbf{x}|\delta, \gamma, \mathbf{y})$  as weighted expectations, by quadrature [61, Sec. 2.1], with  $\sum \pi(\lambda_i|\mathbf{y}) = \sum \pi(\gamma_i|\mathbf{y}) = 1$ . The weights  $\pi(\lambda_i|\mathbf{y})$  and  $\pi(\gamma_i|\mathbf{y})$  are either given by the TT approximation or by the bins for the sample-based histograms. If calculating the variance is too costly, the RTO method (see Sec. 2.2.3) may be a feasible alternative to draw a sample from Eq. 4.17.

Based on the marginal posterior distribution  $\pi(\lambda, \gamma|\mathbf{y})$  we calculate the mean and covariance of the conditional posterior  $\pi(\mathbf{x}|\lambda, \gamma, \mathbf{y})$  by quadrature as in Eq. 4.18 and Eq. 4.19. We can either use the sample based histogram bins as weights or the TT approximation to integrate over marginal approximations  $\pi(\lambda|\mathbf{y})$  and  $\pi(\gamma|\mathbf{y})$ . More precisely, for the sample based evaluation, the height of the histogram bars act as quadrature weights, e.g.  $\pi(\lambda_i|\mathbf{y})$  at the centre  $\lambda_i$  of each bin. Then, we obtain the full conditional mean  $\mu_{\mathbf{x}|\mathbf{y}}$  and covariance matrix  $\Sigma_{\mathbf{x}|\mathbf{y}}$  as weighted expectations. Again, we use Cholesky decomposition to invert  $\mathbf{B}_\lambda = \mathbf{A}^T \mathbf{A} + \lambda \mathbf{L}$  and to calculate  $\mathbf{x}_\lambda = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{L})^{-1} \mathbf{A}^T \mathbf{y}$ . In total we have to evaluate  $\mathbf{x}_\lambda$  and invert  $\mathbf{B}_\lambda$  20 times to obtain mean and covariance of  $\pi(\mathbf{x}|\mathbf{y})$ , Fig. 4.6, which takes less than 0.2s. We plot posterior samples of  $\pi(\mathbf{x}|\mathbf{y})$  in Fig. 4.6, where we set negative ozone VMR to zero. The fact that we have to deal with negative ozone values is due to the poor prior choice of  $\pi(\mathbf{x}|\delta)$ . Note, that the sample mean is slightly larger than posterior mean at heights where the data is noise dominated, and the ozone values are determined by the prior, or where the ground truth is close to zero. This indicates that we should use a different more physical based prior or model to parametrise the ozone profile. Note, that the posterior samples do not represent the ozone peak at around 80km.

**Errors due to number of grid points** we plot RMS for mean and variance compared to a solutop with 200 grid points. The relative error behaves proportionally to  $1/N$ , see Fig. 4.7 and Eq. A.11, and we consider a relative error less than 0.1% good enough. This happens roughly at a bin size of 25, which is our TT grid size. Note that we exclude the error due to  $\tau_{\text{int}}$  the IACT and that we choose the grid according to the sampled values so that the sampling region is the same as the region in which we approximate the posterior distributions. . The error of is influenced by rank and gridsize interpolation error



**Figure 4.6:** We draw ozone samples from the full posterior distribution  $\pi(\mathbf{x}|\mathbf{y})$  after characterising mean and covariance of  $\pi(\mathbf{x}|\mathbf{y})$  by weighted expectations over the marginal posterior  $\pi(\lambda, \gamma|\mathbf{y})$ . we determine  $\pi(\lambda, \gamma|\mathbf{y})$  either through sampling or via TT approximation based on the linear forward map  $\mathbf{A}_L$ . Note that we set negative values ozone VMR values to zero. We will use those samples to find the affine map  $\mathbf{M}$ , see section 4.3



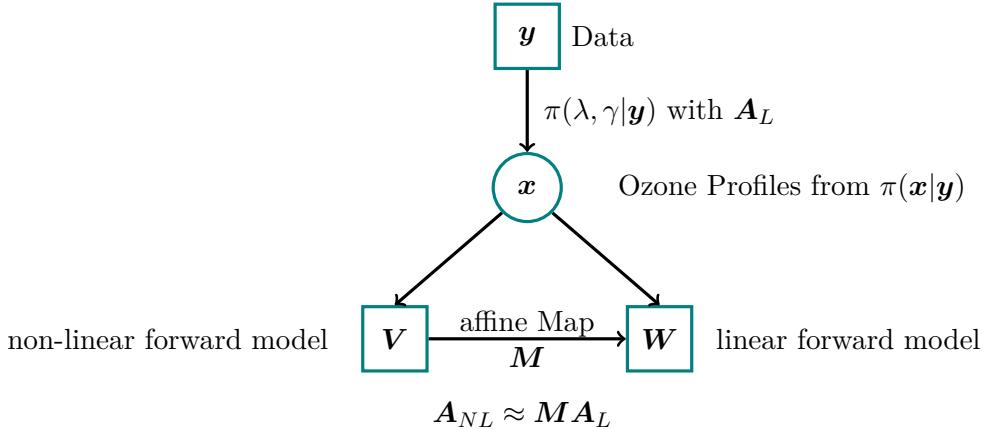
**Figure 4.7:** Assessment of Monte-Carlo error, where we calculate the relative error of the mean due to binning up the samples compared to the sample mean  $\|\mu_{\text{samp}} - \mu_{\text{distr}}\| / \|\mu_{\text{samp}}\|$ .

### 4.3 Approximate non-linear Forward Model with an Affine Map

Given the posterior distribution for ozone  $\pi(\mathbf{x}|\mathbf{y})$ , we can now approximate the non-linear forward model

$$\mathbf{A}_{NL} \approx \mathbf{M}\mathbf{A}_L, \quad (4.20)$$

with an affine map  $\mathbf{M}$ , see Fig. 4.8 for the summarised strategy. Here we write  $\mathbf{A}_{NL}\mathbf{x}$ , which implies that we construct the non-linear forward model and compute non-linear noise-free data  $\mathbf{A}(\mathbf{x}, \mathbf{p}, \mathbf{T})$  based on ground truth pressure and temperature. Using



**Figure 4.8:** The strategy to find the affine map consist of evaluating the marginal posterior for ozone using the linear forward model. Then we draw ozone samples from the conditional posterior and calculate noise free data based on the linear and non-linear forward model. Next we find a mapping in between those two space so that we can approximate the non-linear forward model using an affine map  $\mathbf{M}$ .

posterior ozone samples we generate two affine subspaces and then find the mapping between those. The subspace  $\mathbf{W}$  is created by noise free data based on the linear model  $\mathbf{A}_L$  and  $\mathbf{V}$  by noise free data based on the non-linear model  $\mathbf{A}_{NL}$ , given  $m$  samples  $\mathbf{x}^{(j)} \sim \pi(\mathbf{x}|\mathbf{y})$  for  $j = 1, \dots, m$ . We report a relative RMS difference between  $\mathbf{W}$  and  $\mathbf{V}$  of about 1%, which we aim to reduce through the affine map  $\mathbf{M}$ . More specifically, the affine subspace associated with the linear forward model is

$$\mathbf{W} = \begin{bmatrix} | & | & | \\ \mathbf{A}_L \mathbf{x}^{(1)} & \dots & \mathbf{A}_L \mathbf{x}^{(j)} & \dots & \mathbf{A}_L \mathbf{x}^{(m)} \\ | & | & | \end{bmatrix} \in \mathbb{R}^{m \times m} \quad (4.21)$$

and with the non-linear forward model is

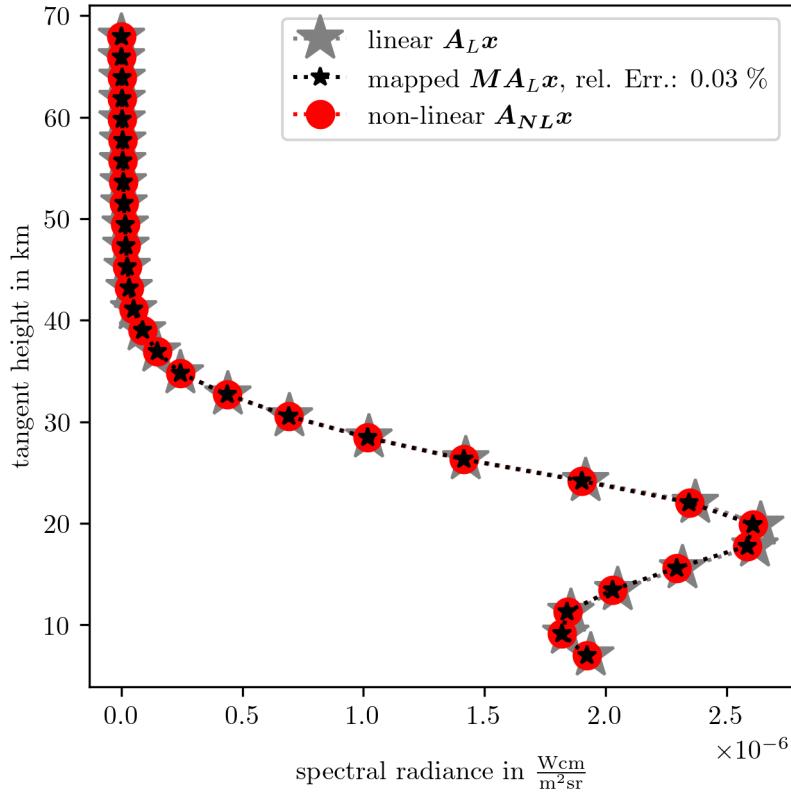
$$\mathbf{V} = \begin{bmatrix} | & | & | \\ \mathbf{A}_{NL}\mathbf{x}^{(1)} & \dots & \mathbf{A}_{NL}\mathbf{x}^{(j)} & \dots & \mathbf{A}_{NL}\mathbf{x}^{(m)} \\ | & | & | \end{bmatrix} = \begin{bmatrix} — & v_1 & — \\ & \vdots & \\ — & v_j & — \\ & \vdots & \\ — & v_m & — \end{bmatrix} \in \mathbb{R}^{m \times m}. \quad (4.22)$$

Then we calculate affine map

$$\mathbf{V}\mathbf{W}^{-1} = \mathbf{M} = \begin{bmatrix} — & r_1 & — \\ & \vdots & \\ — & r_j & — \\ & \vdots & \\ — & r_m & — \end{bmatrix} \in \mathbb{R}^{m \times m}. \quad (4.23)$$

by solving  $v_j = r_j\mathbf{W}$  for each row  $r_j$  in  $\mathbf{M}$ , where  $j = 1, \dots, m$ , using the Python function `numpy.linalg.solve`. We can do that because every measurement in the data vector  $\mathbf{y}$  is independent of each other, and then every row  $v_j$  of  $\mathbf{V} \in \mathbb{R}^{m \times m}$  is independent of each other as well.

We assess the affine map by calculating the relative RMS difference  $\|\mathbf{MW} - \mathbf{V}\|_{L^2}/\|\mathbf{MW}\|_{L^2}$  between the mapped linear noise free data and the non-linear noise free data, which is approximately 0.001%. In Fig. 4.9, we show the mapping for one posterior ozone sample which has not been used to create this mapping. In other word this is an unseen event not in the training data. The relative RMS error for this approximation is roughly 0.03% and much smaller than the relative difference between noise free linear data and non-linear data. Consequently, from here onwards, we use the approximated forward map.



**Figure 4.9:** We asses how good we can map a new ozone sample  $\mathbf{x} \sim \pi(\mathbf{x}|\mathbf{y})$  from the linear forward model onto the non-linear forward model using the previous calculated affine map  $\mathbf{M}$ . The sample has not been used to create this affine map. The gray stars represent noise free linear data, where as the red circles present noise free non-linear data. Then we map the linear noise free data onto the non-linear noise free data, black start, and provide the relative RMS error in between the mapped noise free data and the non-linear data.

## 4.4 Regularisation Solution vs. Bayesian Approach – approximated Model

With the affine approximation we define the forward model matrix

$$\mathbf{A} := \mathbf{M}\mathbf{A}_L \quad (4.24)$$

using the affine map  $\mathbf{M}$ . Here we compare the posterior distribution of ozone to a regularisation approach.

### 4.4.1 Posterior Distribution for Ozone

Again we us the MTC scheme and the exact same setup as in Sec. 4.2.2 to evaluate the marginal posterior and then the conditional posterior of ozone.

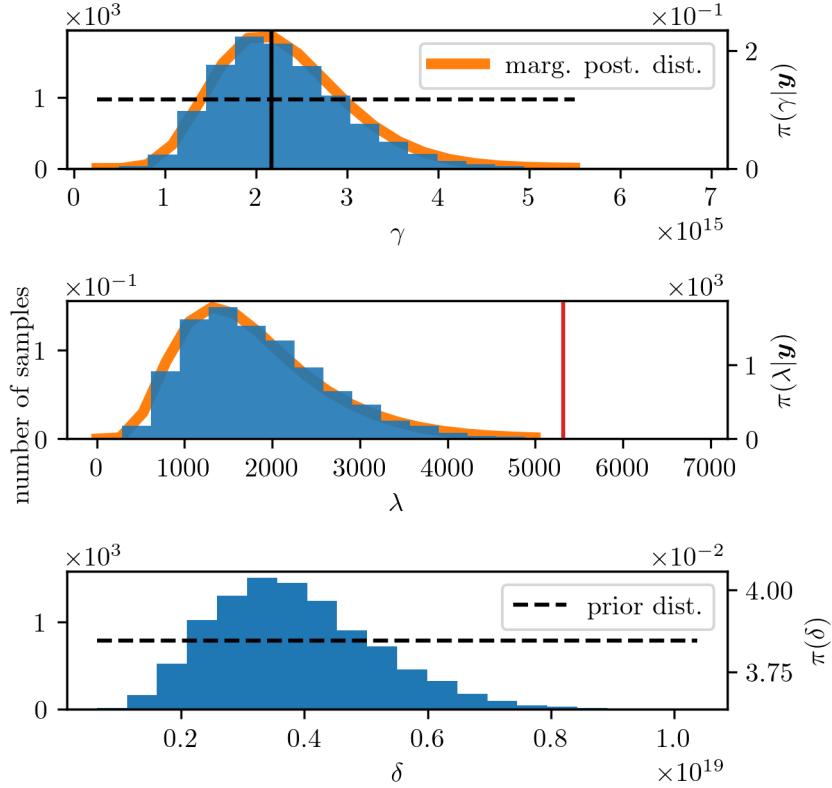
#### Marginal Posterior

The marginal posterior is defined as in Eq. 4.8, but with updated forward model. We initialise the MWG at the mode of  $\pi(\lambda, \gamma|\mathbf{y})$  and approximate  $f(\lambda)$  and  $g(\lambda)$  around the mode as in Eq. 4.12a and Eq. 4.12b. Then we run the MWG algorithm for  $N = 10000$  plus  $N_{\text{burn-in}} = 100$  steps and plot the samples in Fig. 4.10 as well as the marginal approximations provided by the TT decomposition on the same grid as previously defined with same number of ranks (see Sec. 4.2.2).

#### Full Posterior Ozone Mean and Variance

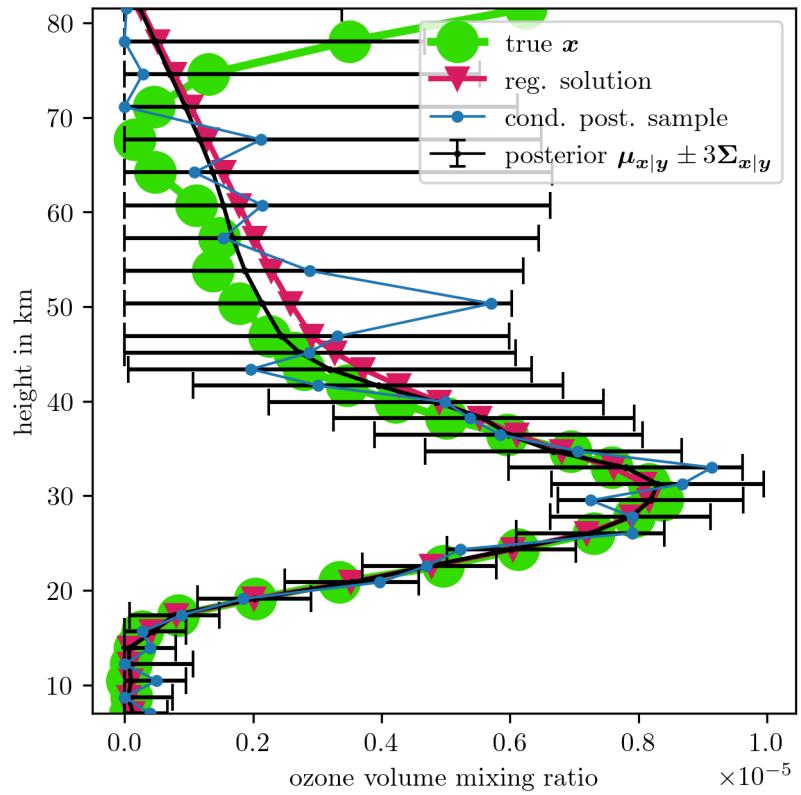
Next, we characterise the conditional posterior  $\pi(\mathbf{x}|\mathbf{y})$  as in Eq. 4.17. Again, we calculate the full posterior mean  $\mu_{\mathbf{x}|\mathbf{y}}$ , see Eq. 4.18, and covariance matrix  $\Sigma_{\mathbf{x}|\mathbf{y}}$  4.19 as weighted expectation over a 20-point grid provided by either the marginal TT approximations of  $\pi(\gamma|\mathbf{y})$  and  $\pi(\lambda|\mathbf{y})$  or by the bins of the sample histogram as quadrature weights. We plot the conditional mean and variance in Fig. 4.11, the regularised solution (see next section), and one sample from the posterior, which represents a feasible solution to this inverse problem. We can see that the ground truth lays within 3 times of the STD (accounting for  $\approx 99\%$  of all posterior samples) around the mean except for the peak at around 80km. We also note that compared to the previously calculated mean and variance based on the linear forward model, see 4.6, the approximated based posterior distribution does not differ significantly. This is expected since the difference between the linear and non-linear forward map of  $\approx 1\%$  is small.

Additionally in Fig. 4.12, we plot the singular values of the covariance matrix  $\Sigma_{\mathbf{x}|\mathbf{y}}$ , to visualise how many ozone values are informative. We observe that the last 10 singular values are very small and correspond to ozone values at the high altitudes

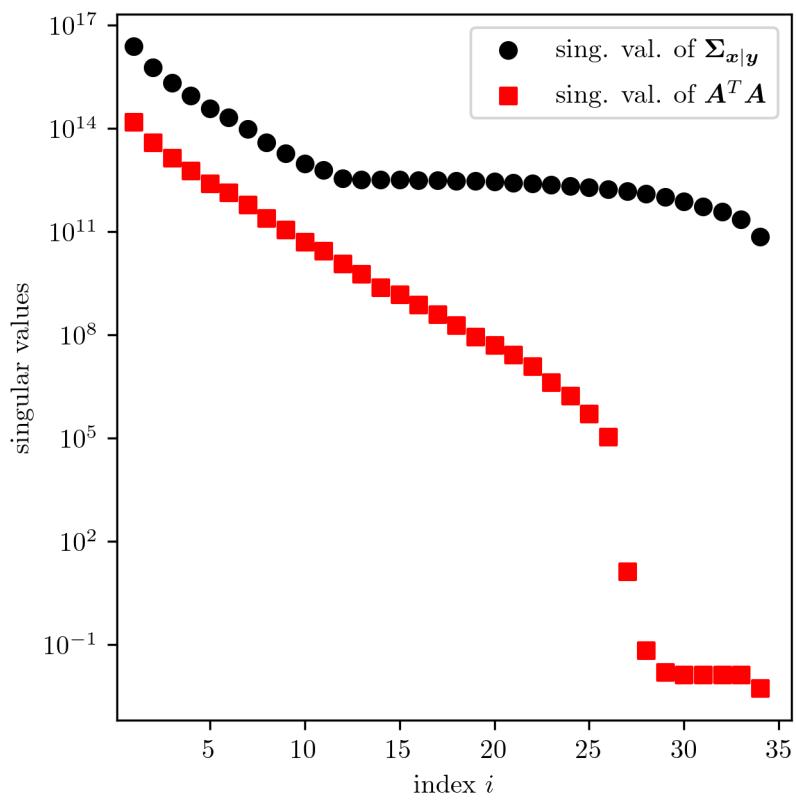


**Figure 4.10:** We plot the TT approximation of marginal posterior in orange and the samples as a histogram as well as the prior distribution with a dotted line. Note that we sample  $\lambda$  and  $\gamma$  using the Metropolis-within-Gibbs sampler and can calculate  $\delta$  for every sample of the marginal posterior, we can not do this for the TT approximation. The regularised parameter corresponding to the regularised solution is marked thought the red vertical line at  $\lambda_{\text{reg}} =$ .

$\gtrsim 45\text{km}$  with a large variance, see Fig. 4.11. This is also roughly the region where we introduce a higher correlation structure due (see Eq. 4.5 Sec. 4.2.1) and the samples are determined by the prior.



**Figure 4.11:** We plot the conditional posterior mean and variance in black and the regularised solution on top of the ground truth ozone profile in green. We use the updated forward map  $MA_L$



**Figure 4.12:** Singular values of the covariance matrix of  $\Sigma_{\mathbf{x}|\mathbf{y}}^{-1} = (\sum(\mathbf{A}^T \mathbf{A} + \mathbf{L})^{-1})^{-1}$  of the posterior distribution  $\pi(\mathbf{x}|\mathbf{y})$  for ozone.

#### 4.4.2 Solution by Regularisation

Since we claim that Bayesian analysis is superior to regularisation methods we compare the MTC method to a Tikhonov regularisation solution, see Sec. 2.5 and [14]. This is most similar to our chosen linear-Gaussian Bayesian framework. The Tikhonov regularised solution is defined as in [14, 63]

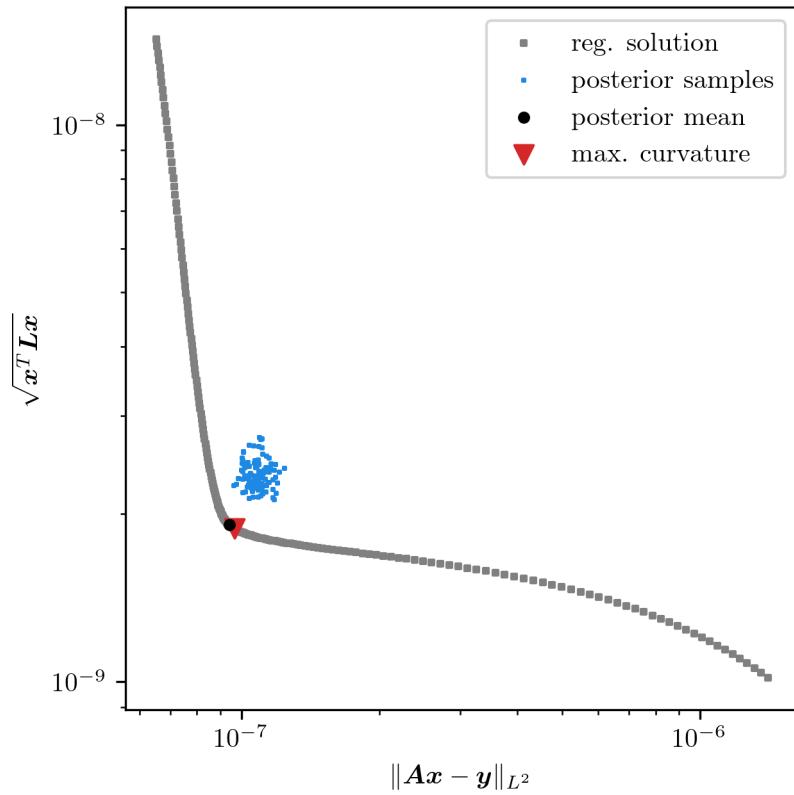
$$\mathbf{x}_\lambda = \arg \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \lambda \mathbf{x}^T \mathbf{Lx}, \quad (4.25)$$

with the regularisation parameter  $\lambda$ . The regularised solution is typically calculated by solving the normal equations, see Sec. 2.5,

$$\mathbf{x}_\lambda = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{L})^{-1} \mathbf{A}^T \mathbf{y}. \quad (4.26)$$

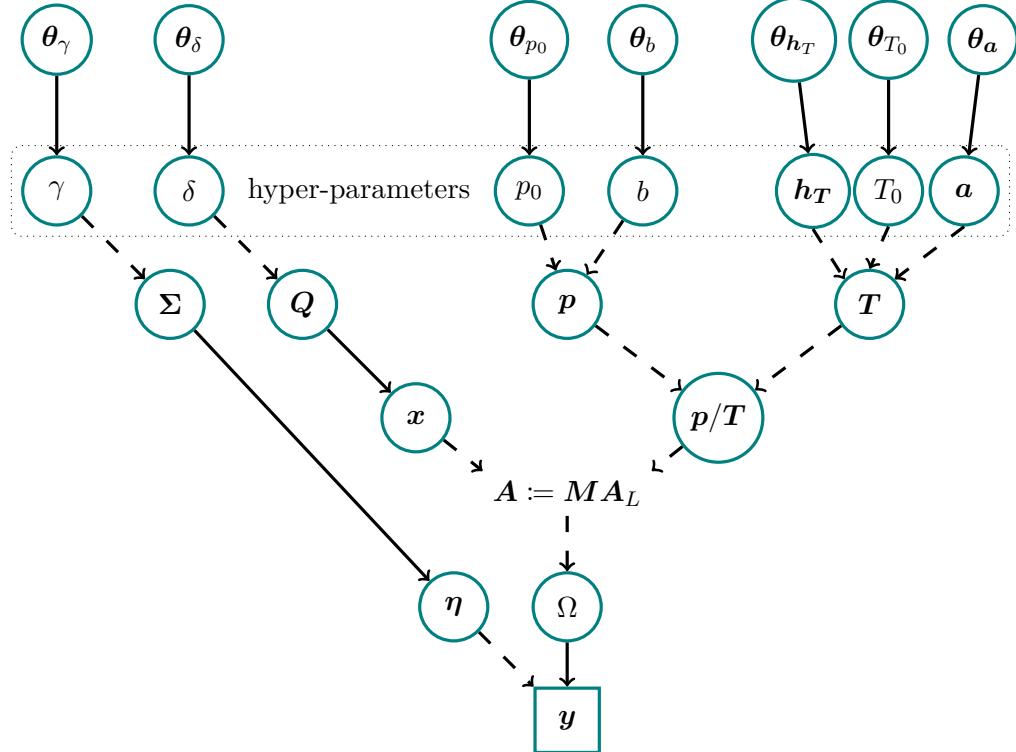
To find the best regularised solution, we use the L-curve method [44]. Within this method we compute  $\mathbf{x}_\lambda$ , for 200 different  $\lambda$  values in between  $10^{0.5}$  to  $10^{6.5}$  and plot the solution semi norm  $\sqrt{\mathbf{x}_\lambda^T \mathbf{L} \mathbf{x}_\lambda}$  against the data misfit norm  $\|\mathbf{Ax}_\lambda - \mathbf{y}\|$ , see Figure 4.13. The best regularised solution corresponding to the corner of the L-curve is located at the point of maximum curvature, see triangle in Fig. 4.13, which we find with the kneedle algorithm [64] using the python function `kneed.KneeLocator` in less 0.1s.

We plot the regularised solution in Fig. 4.11 and observe that it is very similar to the posterior mean. It is pretty clear that the regularised solution accounts for only one possible solution and does not provide uncertainties. The regularised solution is not similar to the samples drawn from the posterior  $\pi(\mathbf{x}|\mathbf{y})$ , see also Fig. 4.6. In Fig. 4.13, the samples of  $\pi(\mathbf{x}|\mathbf{y})$  lie above the L-Curve where as the posterior mean and the regularised solution lie on the L-Curve. This does make sense, if one thinks about the mean (average over non-smooth samples) and the regularised solution as extremely smooth ozone profiles, see also Sec. 2.5. In contrast the samples are less regularised and hence lie above the L-Curve, but have a similar data misfit norm and as already mentioned are all feasible solution to the data. Neither the regularisation solution nor the posterior ozone profiles capture the ozone peak of the ground truth at high altitudes.



**Figure 4.13:** We calculate regularised solution as in Eq. ?? and plot the regularised semi norm  $\sqrt{\mathbf{x}^T \mathbf{L} \mathbf{x}}$  against the data misfit norm  $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|$  to find the regularised solution at the point of maximum curvature of the so-called L-Curve. Additionally we calculate the data misfit norm and the regularised norm for the ozone posterior and for samples of the conditional posterior distribution. [make box around Kneedle reagion](#)

## 4.5 Hierarchical Bayesian Framework for Ozone, Pressure and Temperature



**Figure 4.14:** DAG of Bayesian model for ozone, pressure and temperature. The hyper-parameters  $\mathbf{h}_T = \{h_{T,1}, h_{T,2}, h_{T,3}, h_{T,4}, h_{T,5}, h_{T,6}\}$ ,  $\mathbf{a} = \{a_0, a_1, a_2, a_3, a_4, a_5, a_6\}$ ,  $T_0$ ,  $b$  and  $p_0$  deterministically (dotted line) describe pressure parameter  $\mathbf{p}$  through the function in Eq. 4.28, and temperature parameter  $\mathbf{T}$  through the function in Eq. 4.3. In this case, we decide the hyper-prior distribution  $\pi(\mathbf{h}_T, \mathbf{a}, b, p_0)$  to be a normal distribution, determined by  $\theta_{h_T}, \theta_a, \theta_{T_0}, \theta_b, \theta_{p_0}$ , which represent mean and variances. The ozone parameter  $\mathbf{x}$  is statistically (solid line) described by the prior distribution  $\mathbf{x}|\delta \sim \mathcal{N}(0, \delta \mathbf{L})$ . Here the hyper-parameter  $\delta$  accounts for smoothness in the ozone profile and defines the precision matrix  $\mathbf{Q} = \delta \mathbf{L}$ , where  $\mathbf{L}$  is graph Laplacian as in Eq. 4.5. The noise covariance  $\Sigma = \gamma^{-1} \mathbf{I}$  of the random noise vector  $\boldsymbol{\eta} \sim \mathcal{N}(0, \gamma^{-1} \mathbf{I})$  is defined by the noise hyper-parameter  $\gamma$ . As in Sec. 4.2 described,  $\theta_\delta$  and  $\theta_\gamma$  define the hyper-priors  $\pi(\delta, \gamma)$ . Then, we randomly observe a data set  $\mathbf{y}$  from the space of all measurables  $\Omega$  through the approximated forward model  $\mathbf{MA}_L$  including some added noise. Given the data we like to determine the posterior distribution over the hyper-parameters  $\pi(\mathbf{h}_T, \mathbf{a}, b, p_0, \delta, \gamma | \mathbf{y})$  first and then  $\pi(\mathbf{x} | \mathbf{h}_T, \mathbf{a}, b, p_0, \delta, \gamma, \mathbf{y})$  utilising the MTC scheme.

As in Sec. 4.2, we use the DAG in Fig. 4.14 to visualise the measurement process and correlations in between parameters. We already see that the parameters pressure  $\mathbf{p}$ , temperature  $\mathbf{T}$  and ozone  $\mathbf{x}$  are correlated, and progress deterministically (dashed line) into the forward model, via  $\mathbf{x} \times \mathbf{p}/\mathbf{T}$ . Through their respective prior distributions they generate a space of all possible noise free data  $\Omega$ , from which we observe some data included some added normally distributed noise  $\boldsymbol{\eta}$ . This hierarchical Bayesian framework includes the hyper-parameters  $p_0, b, \mathbf{h}_T, \mathbf{T}_0, \mathbf{a}$ , for pressure and temperature,

model parameters	priors	TT bounds		$\tau_{\text{int}}$	Context
		lower	upper		
$\gamma$	$\mathcal{T}(1, 10^{-10})$	$5 \cdot 10^{-8}$	$4.5 \cdot 10^{-7}$	$9 \pm 0.1$	$\mathbf{y}$
$\delta$	$\mathcal{T}(1, 10^{-10})$	-	-	$1.5 \pm 0.1$	$\mathbf{x}$
$\lambda = \delta/\gamma$	-	500	$10^4$	$3.5 \pm 0.3$	$\mathbf{x}$
$\mathbf{x}$	$\mathcal{N}(0, \delta \mathbf{L})$	-	-	-	$\mathbf{x}$
$p_0$	$\mathcal{N}(1243, 5)$	1229	1259	$550 \pm 9$	$\mathbf{p}/\mathbf{T}$
$T_0$	$\mathcal{N}(288.15, 4.5)$	275	302	$2446 \pm 76$	$\mathbf{p}/\mathbf{T}$
$h_{T,1}$	$\mathcal{N}(11, 0.5)$	9.5	12.5	$1820 \pm 49$	$\mathbf{p}/\mathbf{T}$
$b$	$\mathcal{N}(0.167, 5 \cdot 10^{-4})$	0.165	0.171	$2813 \pm 92$	$\mathbf{p}/\mathbf{T}$
$h_{T,3}$	$\mathcal{N}(32.3, 2.5)$	25.2	39.8	$394 \pm 5$	$\mathbf{p}/\mathbf{T}$
$a_0$	$\mathcal{N}(-6.5, 0.01)$	-6.53	-6.47	$330 \pm 4$	$\mathbf{p}/\mathbf{T}$
$h_{T,2}$	$\mathcal{N}(20.1, 1.6)$	17.7	22.3	$454 \pm 7$	$\mathbf{p}/\mathbf{T}$
$a_1$	$\mathcal{N}(0, 0.1)$	-0.3	0.3	$508 \pm 8$	$\mathbf{p}/\mathbf{T}$
$a_2$	$\mathcal{N}(1, 0.01)$	0.97	1.03	$341 \pm 5$	$\mathbf{p}/\mathbf{T}$
$a_3$	$\mathcal{N}(2.8, 0.1)$	2.5	3.1	$316 \pm 4$	$\mathbf{p}/\mathbf{T}$
$h_{T,4}$	$\mathcal{N}(47.4, 5)$	45.9	48.9	$324 \pm 4$	$\mathbf{p}/\mathbf{T}$
$a_4$	$\mathcal{N}(0, 0.1)$	-0.3	0.3	$335 \pm 4$	$\mathbf{p}/\mathbf{T}$
$h_{T,5}$	$\mathcal{N}(51.4, 5)$	49.9	52.9	$319 \pm 4$	$\mathbf{p}/\mathbf{T}$
$a_5$	$\mathcal{N}(-2.8, 0.1)$	-3.1	-2.5	$335 \pm 4$	$\mathbf{p}/\mathbf{T}$
$h_{T,6}$	$\mathcal{N}(71.8, 3)$	62.5	80.8	$347 \pm 5$	$\mathbf{p}/\mathbf{T}$
$a_6$	$\mathcal{N}(-2, 0.01)$	-2.03	-1.97	$320 \pm 4$	$\mathbf{p}/\mathbf{T}$

**Table 4.2:** Summary of relevant parameter characteristics, bounds and sampling statistics. We denote  $\mathcal{N}(\mu, \sigma)$  as the Gaussian and  $\mathcal{T}(\alpha = \text{scale}, \beta = \text{rate})$  as the gamma distribution. The IACT  $\tau_{\text{int}}$  is estimated as in [65] from posterior samples based on the approximated forward map.

$\delta$  for ozone smoothness and  $\gamma$  for noise precision. Here pressure  $\mathbf{p}$  and temperature  $\mathbf{T}$  are functionally dependent on  $p_0, b, \mathbf{h}_T, \mathbf{T}_0, \mathbf{a}$ , see Eq. 4.3 and Eq. 4.28. Each of those hyper-parameters is described by the hyper-prior distribution  $\pi(\gamma, \delta, p_0, b, \mathbf{h}_T, \mathbf{T}_0, \mathbf{a})$  defined by us. Here  $\boldsymbol{\theta}_\gamma, \boldsymbol{\theta}_\delta, \boldsymbol{\theta}_{p_0}, \boldsymbol{\theta}_b, \boldsymbol{\theta}_{\mathbf{h}_T}, \boldsymbol{\theta}_{\mathbf{T}_0}, \boldsymbol{\theta}_{\mathbf{a}}$  determine gamma distributions  $\gamma, \delta \sim \pi(\boldsymbol{\theta}_\gamma, \boldsymbol{\theta}_\delta)$ , so that e.g.  $\gamma \sim \Gamma(\alpha_\gamma, \beta_\gamma)$  with  $\boldsymbol{\theta}_\gamma = \{\alpha_\gamma, \beta_\gamma\}$ , and a normal distribution  $p_0, b, \mathbf{h}_T, \mathbf{T}_0, \mathbf{a} \sim \pi(\boldsymbol{\theta}_{p_0}, \boldsymbol{\theta}_b, \boldsymbol{\theta}_{\mathbf{h}_T}, \boldsymbol{\theta}_{\mathbf{T}_0}, \boldsymbol{\theta}_{\mathbf{a}})$ , so that e.g.  $b \sim \mathcal{N}(\mu_b, \sigma_b)$  and  $\boldsymbol{\theta}_b = \{\mu_b, \sigma_b\}$ . We denote the forward model as  $\mathbf{A} := \mathbf{M}\mathbf{A}_L$ .

Then, we set up the hierarchical Bayesian framework

$$\mathbf{y}|\mathbf{x}, p_0, b, \mathbf{h}_T, T_0, \mathbf{a}, \delta, \gamma \sim \mathcal{N}(\mathbf{A}(p_0, b, \mathbf{h}_T, T_0, \mathbf{a}) \mathbf{x}, \gamma^{-1} \mathbf{I}) \quad (4.27a)$$

$$\mathbf{x}|\delta \sim \mathcal{N}(\mathbf{0}, (\delta \mathbf{L})^{-1}) \quad (4.27b)$$

$$\delta \sim \Gamma(\alpha_\delta, \beta_\delta) \quad (4.27c)$$

$$\gamma \sim \Gamma(\alpha_\gamma, \beta_\gamma) \quad (4.27d)$$

$$\mathbf{a} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{a}}, \boldsymbol{\Sigma}_{\mathbf{a}}) \quad (4.27e)$$

$$\mathbf{h}_T \sim \mathcal{N}(\boldsymbol{\mu}_T, \boldsymbol{\Sigma}_{\mathbf{h}_T}) \quad (4.27f)$$

$$T_0 \sim \mathcal{N}(\mu_{T_0}, \sigma_{T_0}) \quad (4.27g)$$

$$p_0 \sim \mathcal{N}(\mu_{p_0}, \sigma_{p_0}) \quad (4.27h)$$

$$b \sim \mathcal{N}(\mu_b, \sigma_b) \quad (4.27i)$$

and define a normally distributed likelihood (due to Gaussian noise) and normally distributed priors, where the hyper-prior means and variances are described through  $\boldsymbol{\theta}_{\mathbf{a}} = (\boldsymbol{\mu}_{\mathbf{a}}, \boldsymbol{\Sigma}_{\mathbf{a}})$ ,  $\boldsymbol{\theta}_{\mathbf{h}_T} = (\boldsymbol{\mu}_T, \boldsymbol{\Sigma}_{\mathbf{h}_T})$ ,  $\boldsymbol{\theta}_{T_0} = (\mu_{T_0}, \sigma_{T_0})$ ,  $\boldsymbol{\theta}_{p_0} = (\mu_{p_0}, \sigma_{p_0})$ , and  $\boldsymbol{\theta}_b = (\mu_b, \sigma_b)$ . which we summarise in Tab. 4.2. Before we formulate the posterior distribution we carefully define  $\boldsymbol{\theta}_\gamma, \boldsymbol{\theta}_\delta, \boldsymbol{\theta}_{p_0}, \boldsymbol{\theta}_b, \boldsymbol{\theta}_{\mathbf{h}}, \boldsymbol{\theta}_{T_0}, \boldsymbol{\theta}_{\mathbf{a}}$ .

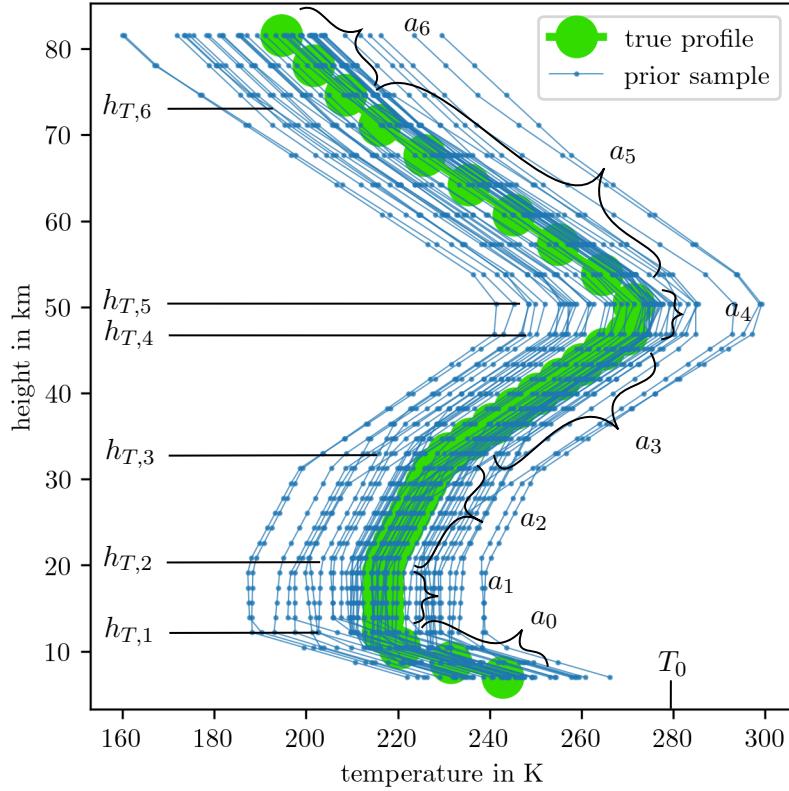
#### 4.5.1 Prior Modelling

The first we do, is to describe the pressure  $\mathbf{p}$  in between  $h_{L,0} \approx 7\text{km}$  and  $h_{L,n} \approx 82\text{km}$  with an exponential function

$$p(h) = \exp(-b h) p_0 \quad , h_{L,0} \leq h \leq h_{L,n} \quad (4.28)$$

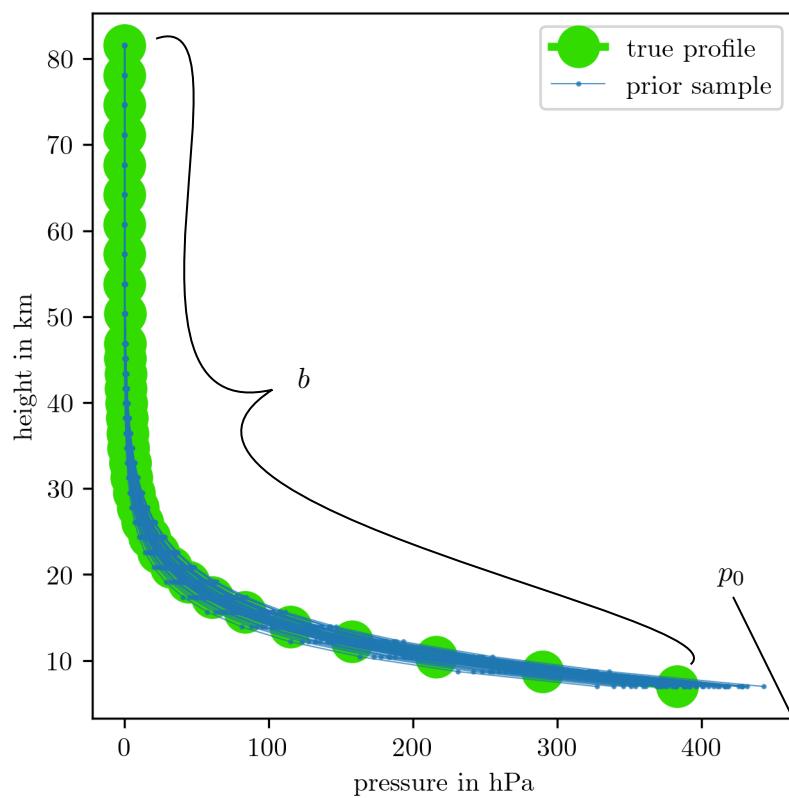
dependent on two hyper-parameters  $p_0, b$ , see Fig. 4.16. Similarly, the temperature as described in Eq. 4.3 can be parametrised with 14 hyper-parameters  $\mathbf{h}_T = \{h_{T,1}, h_{T,2}, h_{T,3}, h_{T,4}, h_{T,5}, h_{T,6}\}$ ,  $\mathbf{a} = \{a_0, a_1, a_2, a_3, a_4, a_5, a_6\}$  and  $T_0$  (see Fig. 4.15 and Eq. 4.3). To complete the model, we have to define a sensible hyper-prior distribution  $\pi(\mathbf{h}_T, \mathbf{a}, T_0)$ . In doing so, we choose the variance and mean of the normally distributed hyper-prior distribution  $\pi(\mathbf{h}_T)$  so that the temperature profile maintains its structure,  $h_{T,i} < h_{T,i+1}$  for  $i = 1, \dots, 5$  (see Fig. B.5). Further, we define  $\pi(\mathbf{a})$  as normally distributed, because we find (see Fig. 4.17) that the data is uninformative about the temperature profile. Similarly, we set  $\pi(T_0)$  to a normal distribution so that it mimics a daily temperature variability of roughly 30K. The hyper-prior distribution  $\pi(p_0, b)$  for pressure-related hyper-parameters is also normally distributed. We choose the variance for  $\pi(p_0)$  so that  $p_0$  has a variability of around 80hPa, close to what we can observe when looking at weather data. Note that we fit one exponential function to ground truth pressure values between  $h_{L,0} \approx 7\text{km}$  and

$h_{L,n} \approx 82$ , so that the pressure values  $p_0$  at sea level may be different to true observed pressure values due to that approximation. We set means of the normal distribution of  $\pi(\mathbf{h}_T, \mathbf{a}, T_0, p_0, b)$  to the ground truth values of  $\mathbf{p}$  and  $\mathbf{T}$ , and summarise the hyper-prior means and variances in Tab. 4.2.

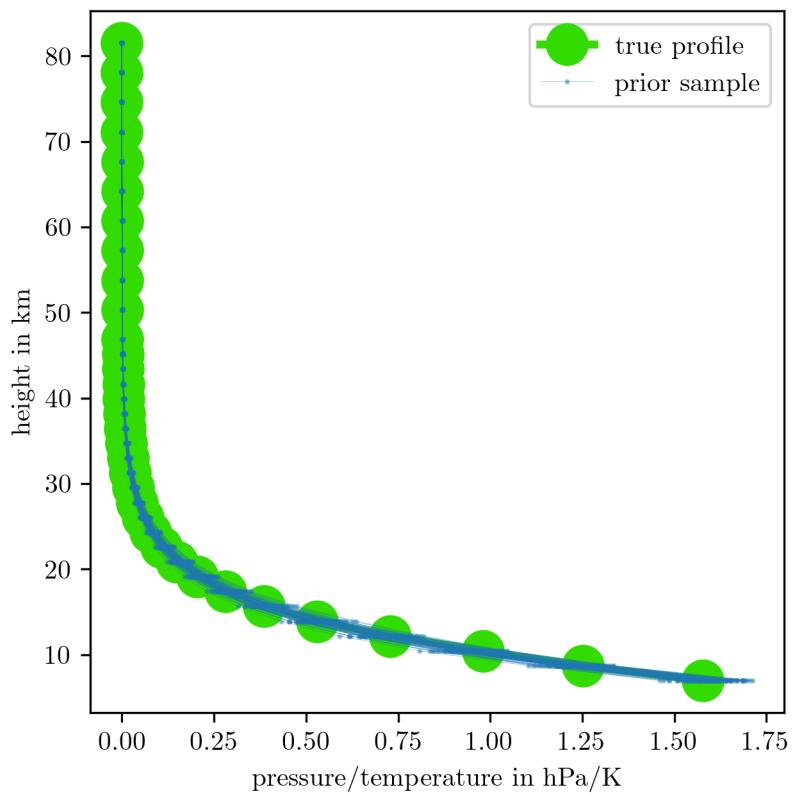


**Figure 4.15:** We draw samples from the hyper-prior distribution of  $h_1, h_2, h_3, h_4, h_5, h_6, a_0, a_1, a_2, a_3, a_4$  and  $T_0$  as defined in table 4.2 and then calculate  $\mathbf{T}$  according to the function in Eq. 4.3.

We plot prior samples of the pressure  $\mathbf{p}$  in Fig. 4.16, the temperature  $\mathbf{T}$  in Fig. 4.15 and the ratio  $\mathbf{p}/\mathbf{T}$  in Fig. 4.17 against the ground truth profiles. Additionally, we plot prior samples of  $1/\mathbf{T}$  in Fig. B.6. Here we already observe that  $\mathbf{p}/\mathbf{T}$  inherits the structure of the pressure function and hence the data is uninformative about the temperature, and that is one of the reasons why we chose those normally, rather informative, hyper-prior distributions.



**Figure 4.16:** We draw samples from the hyper-prior distribution of  $h_0$ ,  $b$  and  $p_0$  as defined in table 4.2 and then calculate  $\mathbf{p}$  according to the function in Eq. 4.28.



**Figure 4.17:** We draw samples from the hyper-prior distribution of  $h_0, b, p_0, h_1, h_2, h_3, h_4, h_5, h_6, a_0, a_1, a_2, a_3, a_4$  and  $T_0$  as defined in table 4.2 and then calculate  $p/T$  according to the functions in Eq. 4.28 and 4.3.

### 4.5.2 Posterior Distribution

Here, we define the marginal and then conditional posterior distribution for the described Bayesian model. We either use the `t-walk` algorithm [26] to draw samples from  $\pi(p_0, b, \mathbf{h}_T, \mathbf{c}_T, \mathbf{a}_T, \delta, \gamma | \mathbf{y})$  or we utilise a TT approximation on a predefined grid to draw samples via the SIRT method with an MH correction step from that marginal posterior. We define  $\boldsymbol{\theta}_{p,T} := \{p_0, b, \mathbf{h}_T, \mathbf{c}_T, \mathbf{a}_T\}$ , which includes all hyper-parameters related to pressure and temperature, for brevity. Lastly, we use the RTO Method to draw samples from the conditional  $\pi(\mathbf{x} | p_0, b, \mathbf{h}_T, \mathbf{c}_T, \mathbf{a}_T, \delta, \gamma, \mathbf{y})$ .

#### Marginal Posterior Distribution

The marginal posterior is given as

$$\pi(\lambda, \boldsymbol{\theta}_{p,T}, \gamma | \mathbf{y}) \propto \lambda^{n/2} \gamma^{m/2} \exp \left\{ -\frac{1}{2} g(\lambda, \boldsymbol{\theta}_{p,T}) - \frac{\gamma}{2} f(\lambda, \boldsymbol{\theta}_{p,T}) \right\} \pi(\lambda, \boldsymbol{\theta}_{p,T}, \gamma), \quad (4.29)$$

with  $\lambda = \delta/\gamma$ ,

$$f(\lambda, \boldsymbol{\theta}_{p,T}) = \mathbf{y}^T \mathbf{y} - (\mathbf{A}(\boldsymbol{\theta}_{p,T})^T \mathbf{y})^T (\mathbf{A}(\boldsymbol{\theta}_{p,T})^T \mathbf{A}(\boldsymbol{\theta}_{p,T}) + \lambda \mathbf{L})^{-1} (\mathbf{A}(\boldsymbol{\theta}_{p,T})^T \mathbf{y}), \quad (4.30a)$$

$$\text{and } g(\lambda, \boldsymbol{\theta}_{p,T}) = \log \det (\mathbf{A}(\boldsymbol{\theta}_{p,T})^T \mathbf{A}(\boldsymbol{\theta}_{p,T}) + \lambda \mathbf{L}). \quad (4.30b)$$

Here we use the approximated forward model  $\mathbf{A}(\boldsymbol{\theta}_{p,T}) := \mathbf{M}\mathbf{A}(\boldsymbol{\theta}_{p,T})_L$  and do not approximate  $f$  and  $g$  but calculate  $\mathbf{A}(\boldsymbol{\theta}_{p,T})$  for each function evaluation directly.

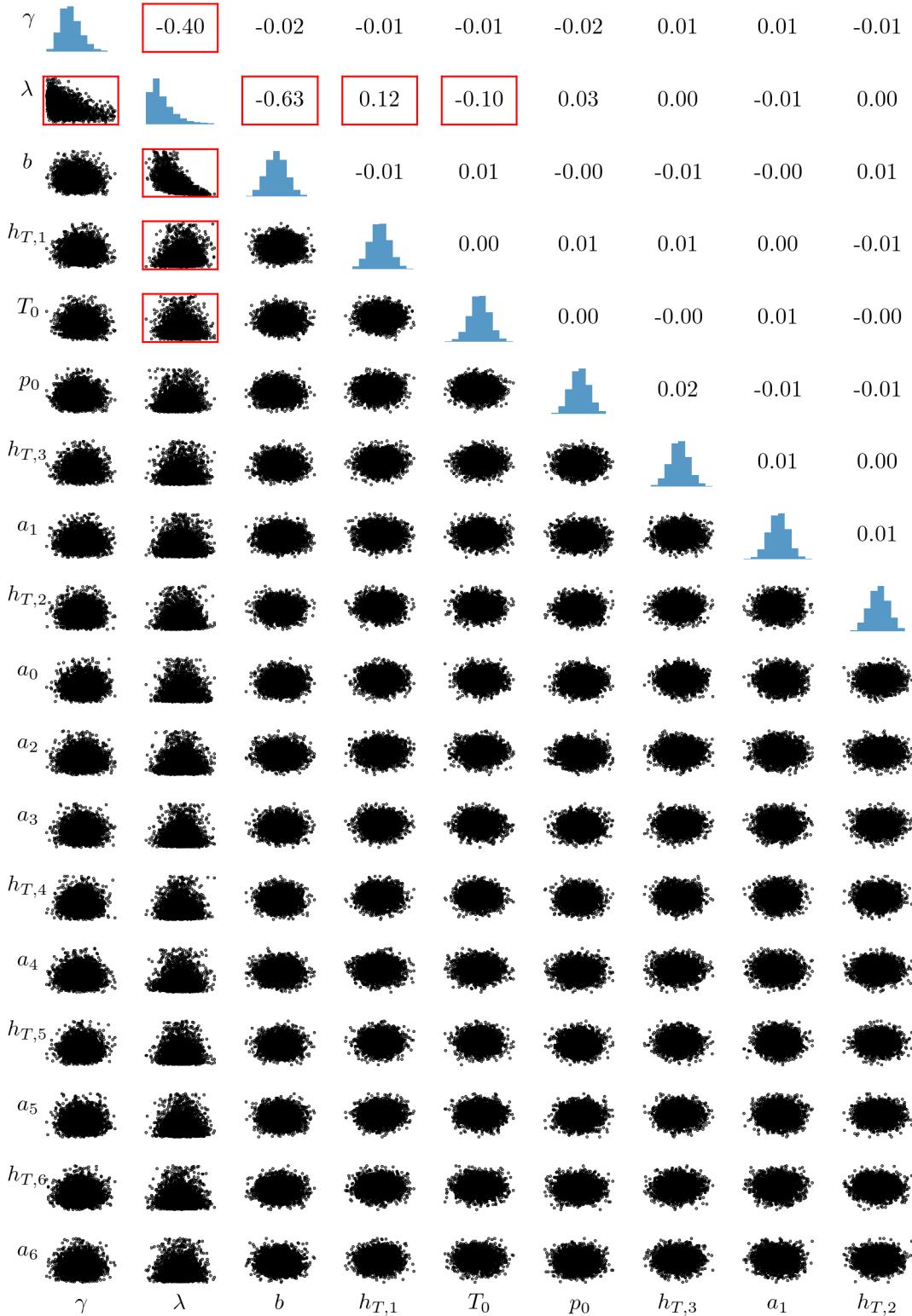
**Sampling from Marginal Posterior** For a ground truth, we run the `t-walk` [26] algorithm on  $\pi(\lambda, \boldsymbol{\theta}_{p,T}, \gamma | \mathbf{y})$  for  $N = 1000 \times 1000$  steps with a burn-in period of  $N_{\text{burn-in}} = 100 \times 1000$ , since we expect a maximal IATC, through some pre-analysis, of around  $\tau_{\text{int,max}} = 500$ , see Tab. 4.2, to generate around 1000 independent samples from the posterior. We initialise the Python implementation of the `t-walk`[66] around the hyperprior mean values and take a total number of  $N' = N + N_{\text{burn-in}} = 1100000$  steps around 5 mins. We iteratively define the TT grid according to the results of the `t-walk`, and also use these boundaries for the hyper-parameters when running the `t-walk` (see Tab. 4.2). We plot the resulting histograms in Fig. 4.20 to 4.24 and the trace of the samples in Fig. B.7.

**TT Approximation of Marginal Posterior** The aim now is to approximate the square root of the marginal posterior

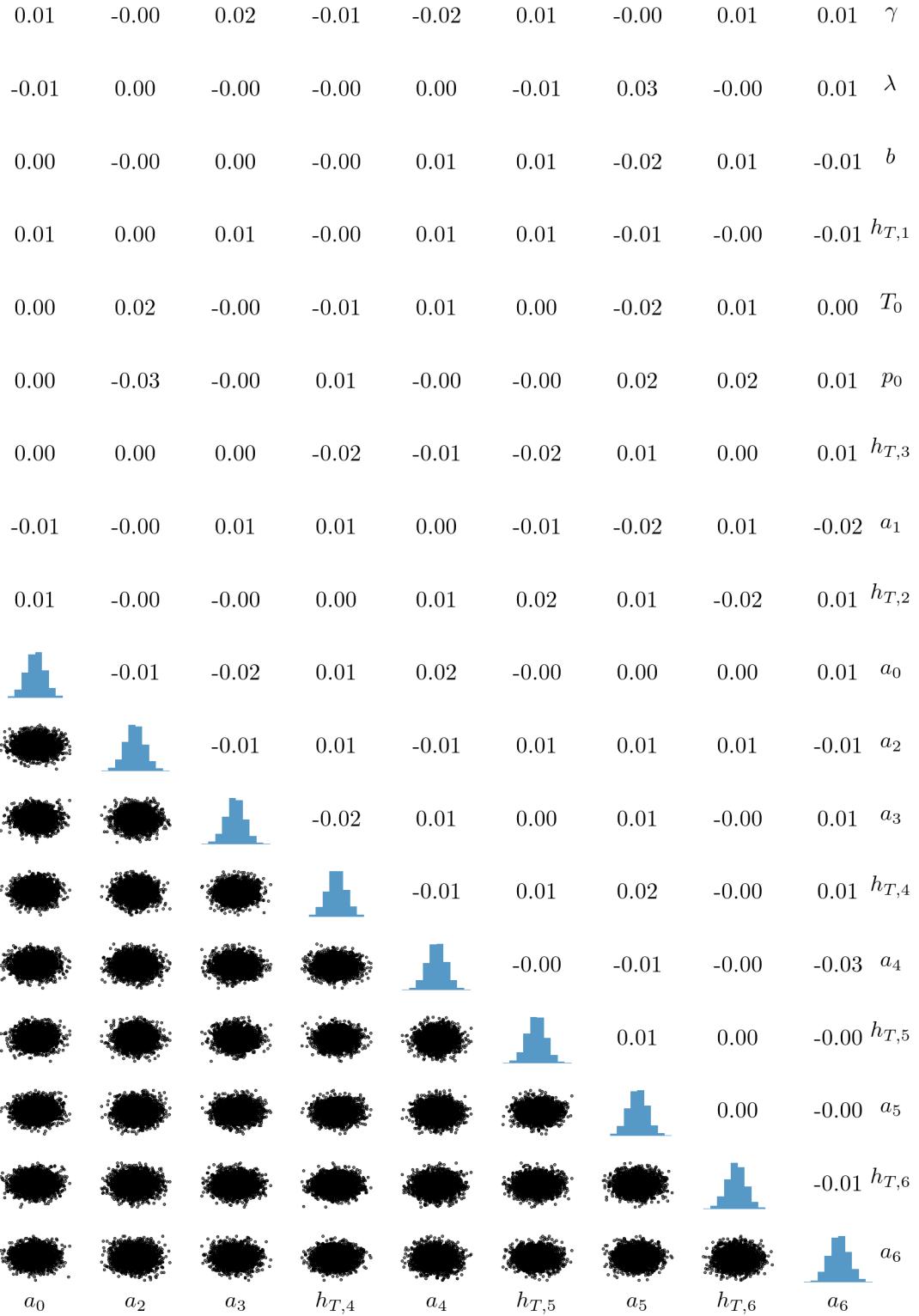
$$\pi(\lambda, \boldsymbol{\theta}_{p,T}, \gamma | \mathbf{y}) \propto \exp \left\{ \log \pi(\lambda, \boldsymbol{\theta}_{p,T}, \gamma | \mathbf{y}) + c \right\}, \quad (4.31)$$

where we introduce a "normalisation constant"  $c = -100$ , similar to Sec. 4.2.2, to avoid under or overflow. In doing so we run the `tt.cross.rectcross.rect_cross.cross` function from the `ttipy` python package [60]. We compute the marginal as in Sec. 2.3, where we set  $\xi = 1/\lambda(\mathcal{X})$  and  $\lambda(x) = 1$  so that for Cartesian basis  $\mathbf{M}_k = \text{diag}(\lambda_k(\mathcal{X}_k))$ . To draw samples from this TT approximation we use the SIRT-MH scheme as in Sec. ??

**Correlation structure** First, we arrange the order of hyper-parameters according to their correlation structure to make the TT approximation more efficient. In doing so, highly correlated hyper-parameter pairs are next each other, so that their TT cores are direct neighbours and linked through their shared rank. See Fig.4.18, where we scatter plot the samples from the TT approximation of  $\sqrt{\pi(p_0, b, T_0, \mathbf{h}_T, \mathbf{a}_T | \mathbf{y}, \gamma, \mathbf{x})}$  via SIRT-MH scheme and plot the Pearson correlation coefficient. A coefficient close to 1 or  $-1$  corresponds to high correlation, whereas low correlation between hyper-parameters has a coefficient close to zero. We observe that the  $\lambda$  hyper-parameters is highly correlated with  $b$  the hyper-parameter describing the pressure gradient and with  $\gamma$ . Additionally,  $h_{T,1}$  and  $T_0$  describing the temperature at low altitudes are correlated to  $b$  as well, since those do very slightly influence "the smoothness" of  $\mathbf{p}/\mathbf{T}$ , hard to see in Fig. 4.17. The IACTs in Tab. 4.2 agree with these results. Note that the correlation for parameters describing temperature in higher altitudes and surprisingly  $p_0$  are very much uncorrelated. Next we aim to find the optimal rank and grid size to approximate  $\sqrt{\pi(p_0, b, T_0, \mathbf{h}_T, \mathbf{a}_T | \mathbf{y}, \gamma, \mathbf{x})}$ , so that the number of function evaluation is as low as possible without too much sacrifice in accuracy.



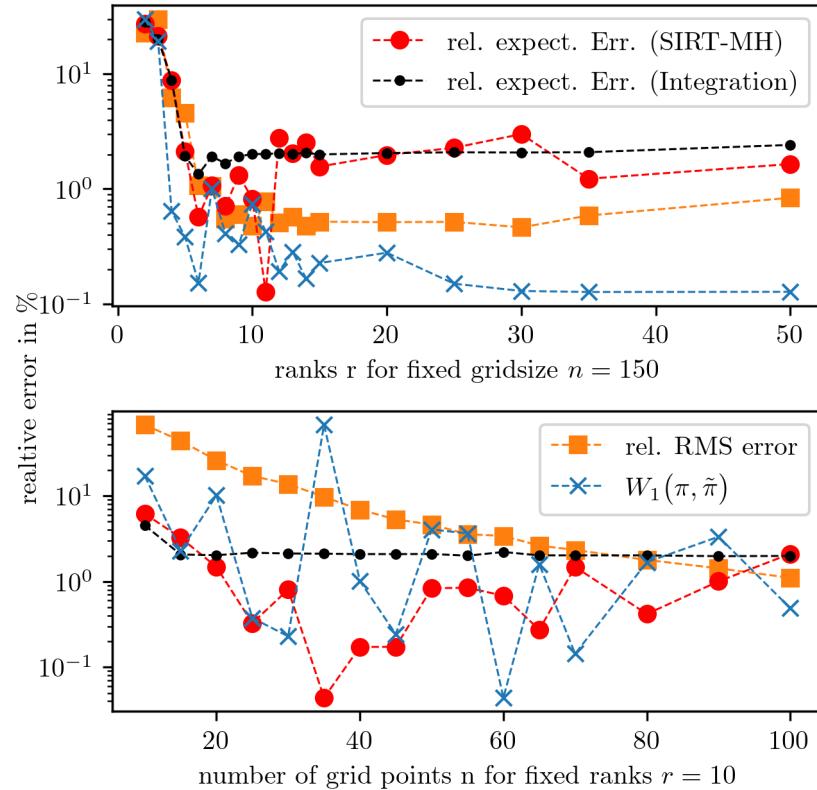
**Figure 4.18:** Scatter plot of samples from TT-approximation of  $\sqrt{\pi(p_0, b, T_0, \mathbf{h}_T, \mathbf{a}_T | \mathbf{y}, \gamma, \mathbf{x})}$  via SIRT scheme. We plot the Pearson correlation coefficient ranging from  $-1$  to  $1$  for each hyper-parameter pair.



Correlation plot of samples from TT-approximation of  $\sqrt{\pi(p_0, b, T_0, \mathbf{h}_T, \mathbf{a}_T | \mathbf{y}, \gamma, \mathbf{x})}$  via SIRT scheme.

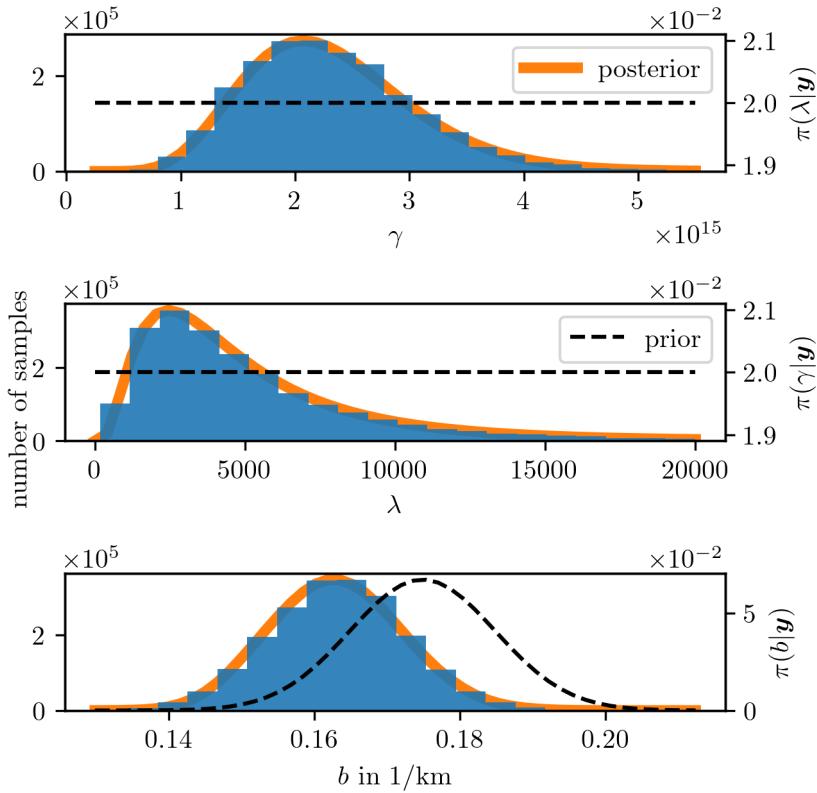
**Find optimal Rank and Grid Size** Next, we choose a relatively large number of grid points of  $n = 150$  and calculate different error measures for decreasing number of ranks to find the optimal number of ranks. Then we fix a small but tolerable rank and decrease the number of grid points until sufficient accuracy.

We calculate the 1-Wasserstein distance, see Sec. 2.3.3, between 5000 independent samples from the SIRT-MH scheme, weighted with the TT approximation of marginal posterior, and 4101 independent samples from the `t-walk`, weighted by the true posterior value. To calculate the 1-Wasserstein distance, as in Eq. 2.52, we use the `SamplesLoss("sinkhorn", p=1, blur=0.05, scaling=0.8)` function with default settings from the python package `geomloss` [67]. This provides the unbiased Sinkhorn divergence which converges towards the Wasserstein distance and can be understood as the generalised Quicksort algorithm [38]. Here,  $p = 1$  defines the distance measure  $\|\mathbf{x} - \tilde{\mathbf{x}}\|_{L^2}$ , the blur parameter can be understood as an entropic penalty and the scaling parameter specifies the trade-off between speed ( $\text{scaling} < .4$ ) and accuracy ( $\text{scaling} > .9$ ) [67]. Additionally we use the marginal functions of each TT approximation to calculate the weighted means  $\boldsymbol{\mu}_{\text{TT}} \in \mathbb{R}^{18}$  of each hyper-parameter and then the relative RMS difference  $\|\boldsymbol{\mu}_{\text{TT}} - \boldsymbol{\mu}_{\text{t-walk}}\|_{L^2}/\|\boldsymbol{\mu}_{\text{t-walk}}\|_{L^2}$  compared to "true means"  $\boldsymbol{\mu}_{\text{t-walk}}$  from the `t-walk`. Additionally, we compare the mean  $\boldsymbol{\mu}_{\text{SIRT-MH}}$  based on the samples from the SIRT-MH scheme with the "true" `t-walk` mean. We plot all of these measures in Fig. 4.19 and decide that a rank  $r = 10$  sufficient because the 1-Wasserstein distance is relatively stable for  $r \geq 10$ . We decide that a grid size  $n = 40$  is large enough based in the relative differences of the samples based and integrated mean compared to the `t-walk` mean, which also seem to be relatively stable for  $n \geq 40$ .



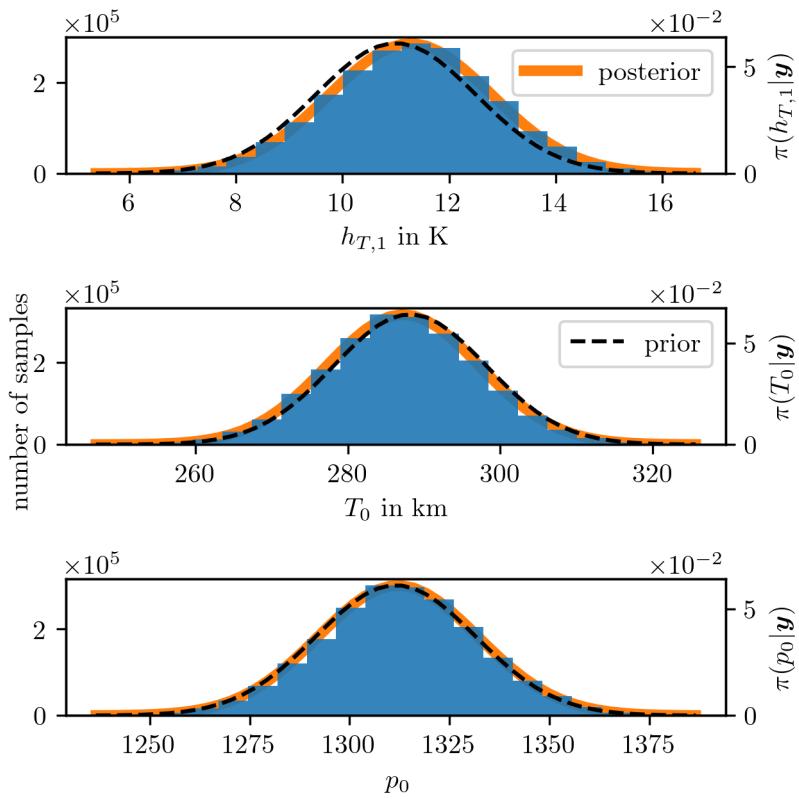
**Figure 4.19:** Given a TT approximation of  $\sqrt{\pi(\lambda, \theta_{p,T}, \gamma | \mathbf{y})}$  we calculate the realtive RMS error and the 1-Wasserstein distance between approximated values at sample points provided by the SIRT-MH and the true function values. Further we use the marginal function from the TT approximation to calculate mean values of each hyper-parameter as weighted expectations and compare to the sample mean provided by the **t-walk**. Additionally we compare sample based mean from the SIRT-MH and the **t-walk**. We plot the relative RMS difference between those calculations. Note that the y-axis represents the recreative errors related to all measures but  $W_1$ . We do so since the each hyper-parameters has a different scale we are only interested in the trend of the  $W_1$ .

To decrease the number of functions evaluation and define ranks  $r = [1, 10, 10, 10, 10, 10, 10, 5, 5, 5, 5, 5, 5, 3, 2]$ , according to the correlation structure of  $\pi(\lambda, \theta_{p,T}, \gamma | \mathbf{y})$  (see Fig. 4.18). We need one sweep in `tt.cross.rectcross.rect_cross.cross`, reducing the computation time to  $\approx 0.5\text{min}$  and the number of functions evaluations to 48240. Note that we can initialise the `tt.cross.rectcross.rect_cross.cross` algorithm at a previously calculated approximation. We plot the marginals for each hyper-parameter in Fig. 4.20 to Fig. 4.25 and samples in Fig. 4.18 .

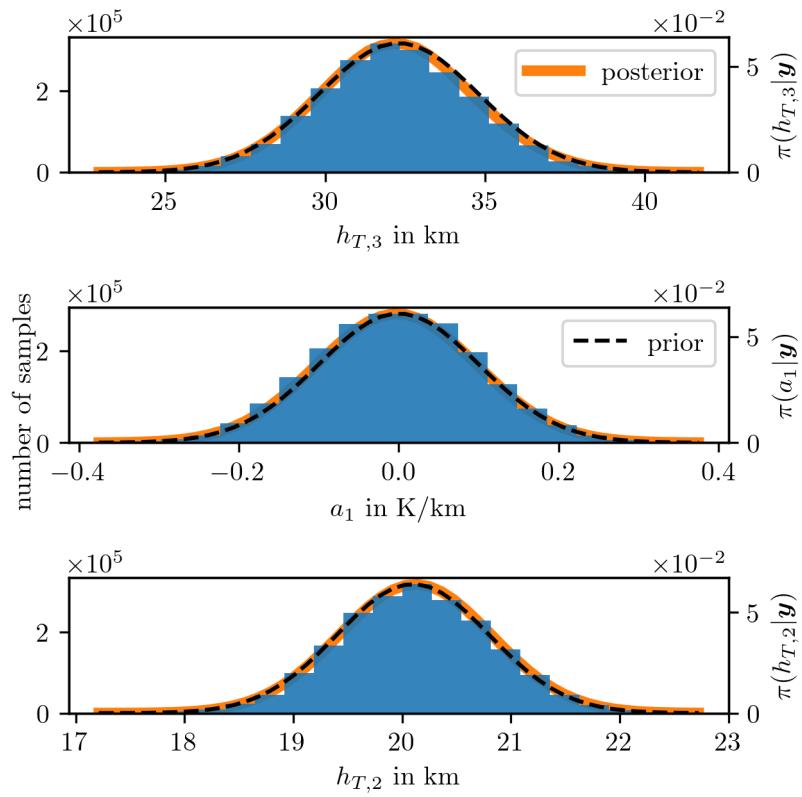


**Figure 4.20:** We plot the TT approximation of marginal posterior in orange and the samples as a histogram as well as the prior distribution with a dotted line.

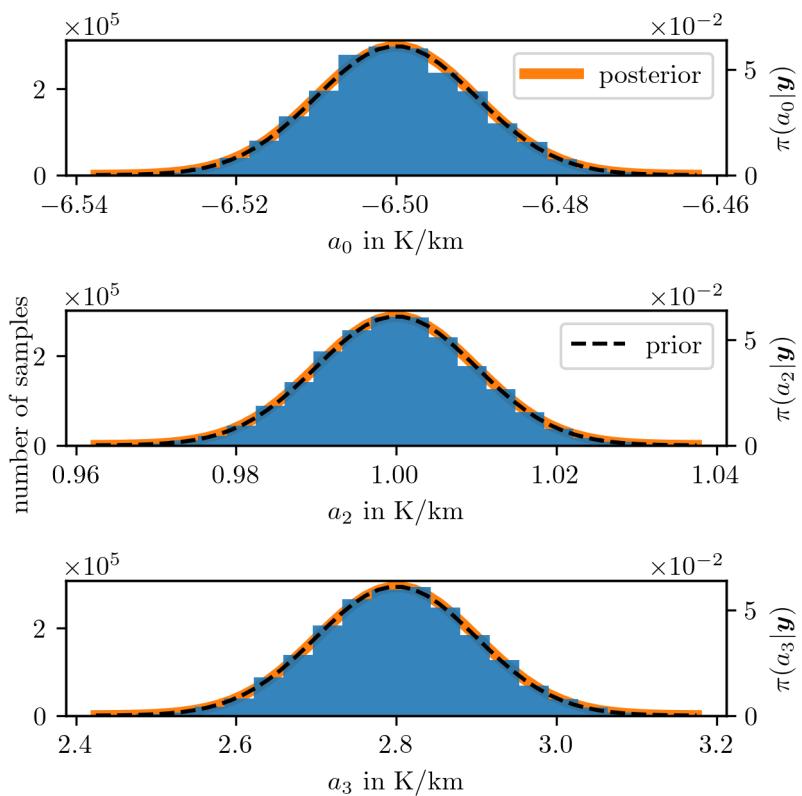
We observe that the marginal do not differ significantly from the t-walk samples and report a relative RMS error between true and TT approximated function values at SIRT-MH samples of  $\approx 7\%$ . We report IATC for the SIRT-MH around 2 – 4 and the IATC given in Tab. 4.2. Additionally, we see when comparing the marginal posterior distributions to the prior distributions of the hyper-parameters that the hyper-parameter  $b$  is the only hyper-parameter related to pressure and temperature which is affected by the data.



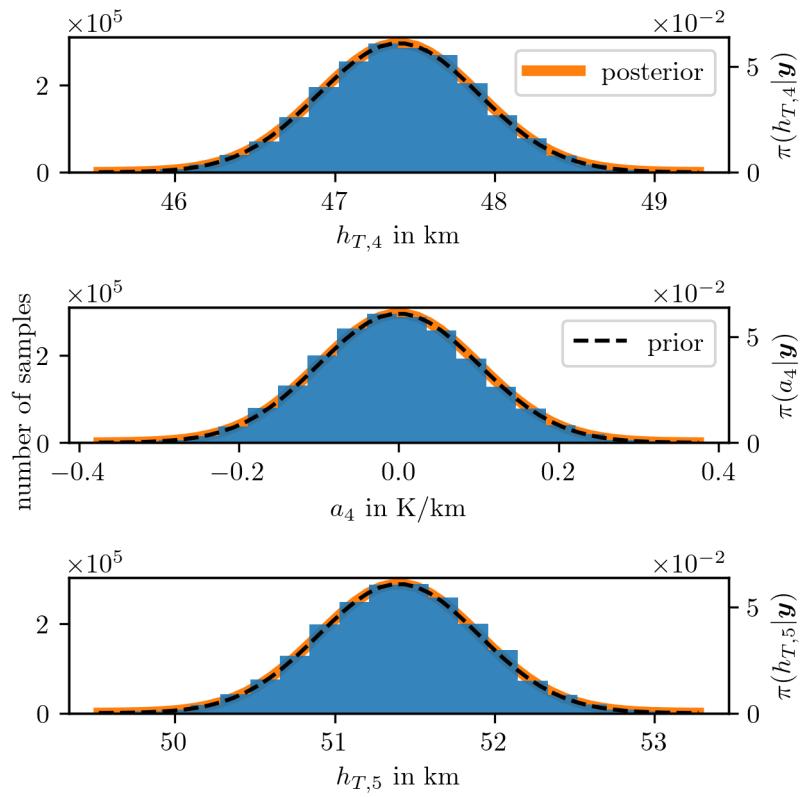
**Figure 4.21:** We plot the TT approximation of marginal posterior in orange and the samples as a histogram as well as the prior distribution with a dotted line.



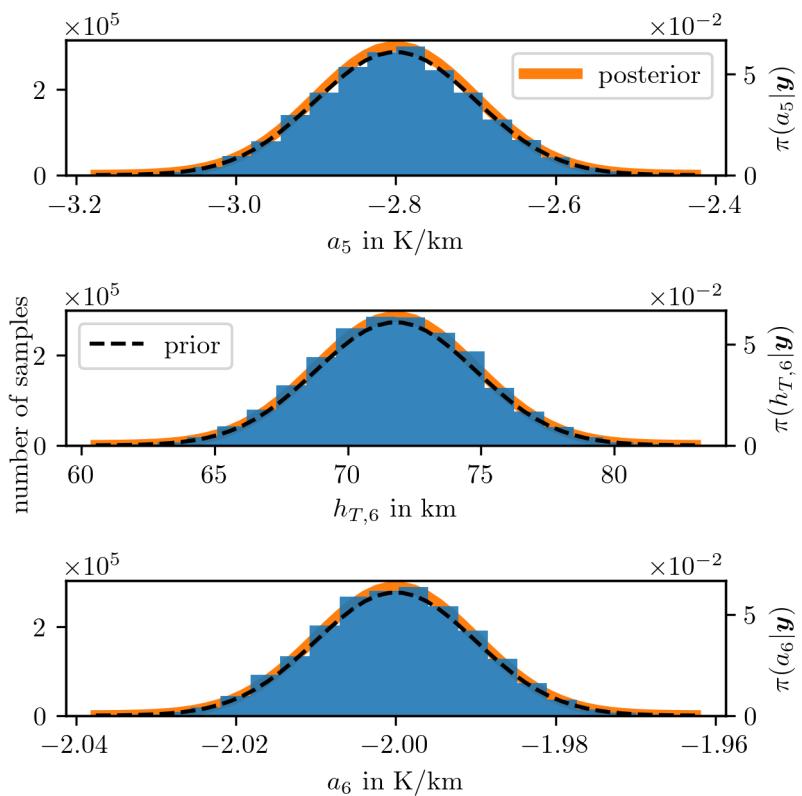
**Figure 4.22:** We plot the TT approximation of marginal posterior in orange and the samples as a histogram as well as the prior distribution with a dotted line.



**Figure 4.23:** We plot the TT approximation of marginal posterior in orange and the samples as a histogram as well as the prior distribution with a dotted line.



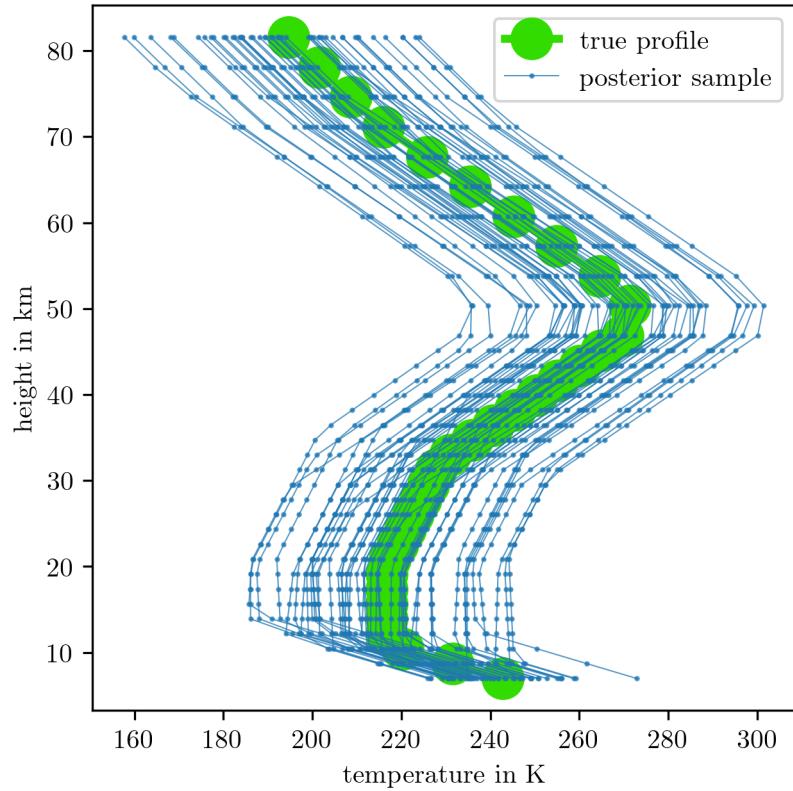
**Figure 4.24:** We plot the TT approximation of marginal posterior in orange and the samples as a histogram as well as the prior distribution with a dotted line.



**Figure 4.25:** We plot the TT approximation of marginal posterior in orange and the samples as a histogram as well as the prior distribution with a dotted line.

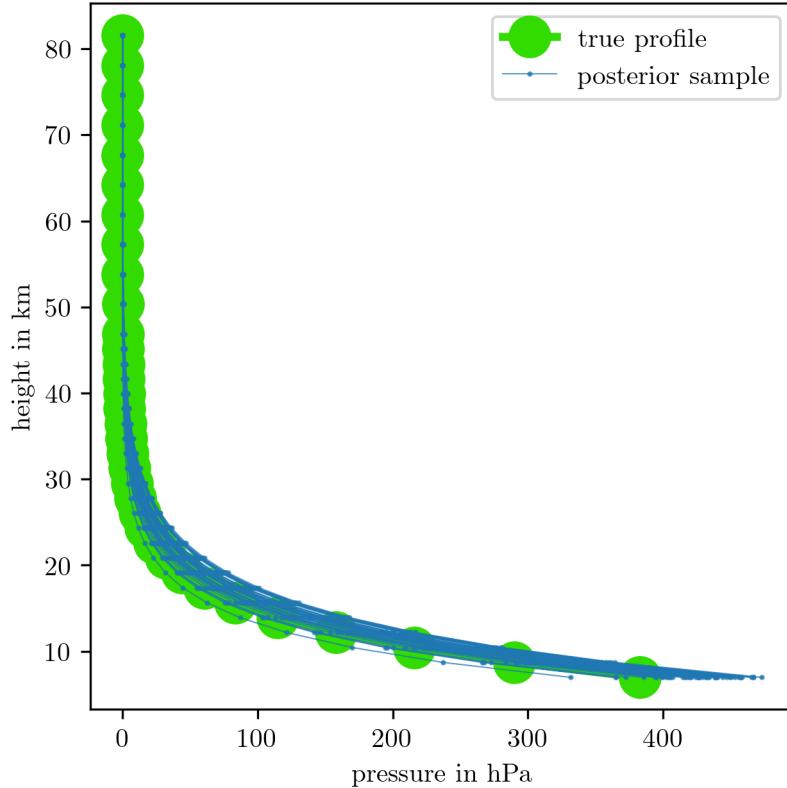
### Conditional Posterior Distribution

We use the RTO method (see Sec. 2.2.3) to obtain ozone samples from the conditional posterior. To obtain posterior pressure and temperature samples take hyper-parameter samples directly from the marginal posterior. Then we calculate pressure and temperature profile according to their respective function (see Eq. 4.28 and Eq. 4.3). We plot the posterior temperature and pressure profiles in Fig. 4.26 and Fig. 4.27 and posterior ozone in 4.28.



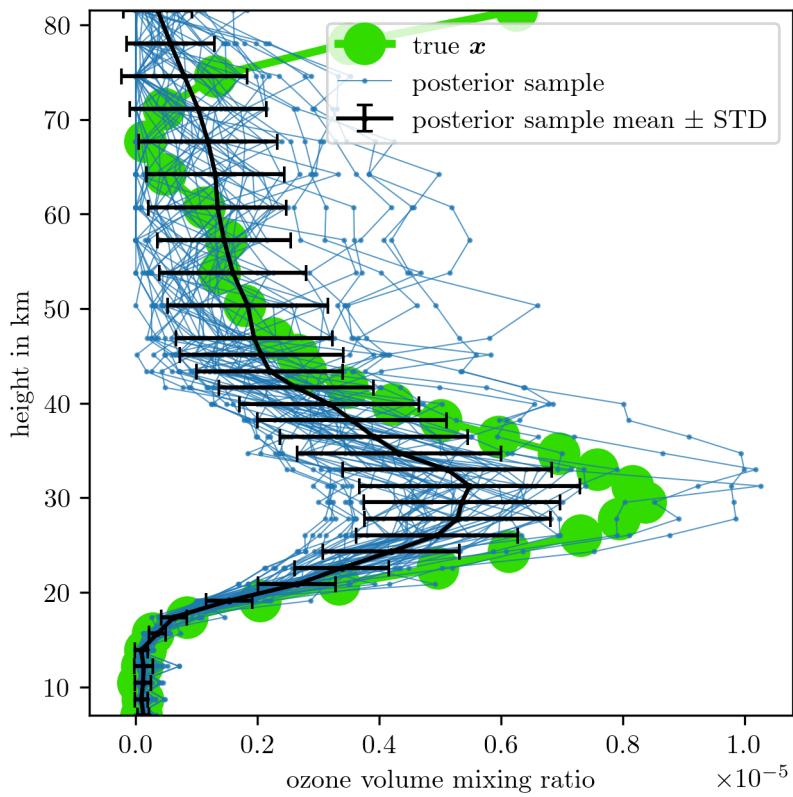
**Figure 4.26:** We take samples from the posterior distribution, as plotted in Figures 4.20 to 4.23 and plot the corresponding temperature function, see Eq: 4.3.





**Figure 4.27:** We take samples from the posterior distribution, as plotted in Fig. 4.24 and plot the corresponding pressure function, see Eq: 4.28.

We observe that pressure and ozone are highly correlated. Since the hyper-parameter  $b$  is smaller than the ground truth value the posterior pressure profile does not exponentially decrease as much compared the ground truth. This results in posterior pressure values which are slightly larger than the ground truth. In comparison, the corresponding posterior ozone profiles have much smaller individual ozone VMRs compared to the ground truth, but still a similar structure compared to the ground truth. Again we are not able to recover the peak at large altitude.



**Figure 4.28:** We take samples from the posterior distribution, as plotted in Fig. 4.24 and plot the corresponding pressure function, see Eq: 4.28.

# 5

## Summary and Outlook

In this chapter we draw conclusion based on the results from the previous chapter. We compare the regularised solution to the mean and to the samples from the full posterior. We elaborate on the occurring approximation errors. We compare the marginal posterior distributions based on the drawn samples and from the TT-decomposition. While elaboration about the different methods, we also elaborate on how informative the data and what the means in terms of ozone, pressure and temperature profile.

### 5.1 Regularisation Solution vs. Hierarchical Bayesian Approach

200 functino ecalution of  $\mathbf{x}_\lambda$  vs 10000 samples vs 600 funciton evalution of marg plus 20  $\mathbf{x}_\lambda$  and 20  $\mathbf{B}_\lambda^{-1}$  time similar within a second but no varicne As already mentioned the regularisation approach only provides on solution, see Fig. 4.11. In Fig. 4.13 we plot samples from the full conditional, which lay above L-Curve and make sense in terms of the Lagrange multipliers as the point on the L-Curve can be seen as extreme values. So the regularised estimate does not correlate to posterior solutions of the inverse problem. We note that the mean of full conditional is very similar to the regularised solution but is also some sort of an extreme value.

In comparison to the regularisation solution we can provide the mean and variance of the full conditional posterior distribution, as well as the sample mean.

### 5.2 Sampling Methods vs. TT Approximation

10000 samples vs 600 funciton evalution of marg similair time due approxiamtions of f and g the recto corss semms to be optimised than my code

We can conclude that the TT approximation is faster or as fast as sampling methods. For the marginal posterior  $\pi(\gamma, \lambda | \mathbf{y})$  the calculation of the TT-cores takes 0.1s, which we consider similar to the sampling time of 0.5s. But the TT approximation needs less function evaluations than the MWG sampler. More precisely, the TT needs  $n_{\text{tot}} = 2n_{\text{sweeps}}((d-2)r^2n + 2nr) = 400$  function evaluations, with number of sweeps  $n_{\text{sweeps}} = 2$  and rank  $r = 10$  and grid size  $n = 25$  compared to 10000 samples. Error due to approcimation of f and g but also due TT

for PT once trained signifincaly faster but t walk is not the optimal sampler since it doesn consiger correaltopij structure Times

When approximating the posterior distribution of the temperature pressure ratio we are much faster compared to sampling methods. Since the parameter space is 16-dimensional we have to run the t-walk for about 2 million steps. In addition of checking the trace of the samples, we also estimate the IATC with [65] see Tab. 4.2. Since for shorter chains with a sample size of  $10^6$  the error for the IATC estimate is much larger we decide to a sample size of  $4 \cdot 10^6$  is sufficient. This comes with a sampling time of 20mins, much larger compare to the 2.5min. Which makes sense as we need  $n_{\text{tot}} = 2n_{\text{sweeps}}((d-2)r^2n + 2nr) = 384838438$  function evaluations. We also note that we do run into problems especially in higher dimensional functions as we have a large range of values and hence introduce the constant c as already mentioned. The t-walk is more robust but the TT approximation is faster.

But both the samples and the TT approximation point towards the same results.

Reduce correlation structure by rotating coordinate system Errors

Within the TT approximation we run into numerical problems. One way of solving this issue could be to use a different basis set such as Lagrange polynomials as these exactly fit to a Gaussian or Chebychev polynomial as basis functions. Another idea is to use different reference measure for integration, such as a Gaussian measure instead of the current Lebesgue measure. Or that the TT finds normalisation constants automatically.

The t-walk is a robust easy to implement sampling method, of course one could employ a more efficient sampler such as a gibbs sampler or something similar [].

We consider the approximation errors of the functions  $f(\lambda)$ ,  $g(\lambda)$  and propagation error into the marginal posterior for sampling of about 10% good enough. The TT approximation error from the marginal posterior is with about 10% also good enough since we do not believe that our model is accurate enough to capture those differences.

When approximation the affine map we get an relative error of about 0.4%, which is much smaller than the relative difference in between noise free and noisy data of approximately of 1.7%. We like to note that the relative difference The error linear to non-linear. Low rank bound as in [68]

## 5.3 Atmospheric Physics

Here we want to say how informative the data is and what we can about the ozone pressure and temperature profiles.

So all the samples as in Fig. 4.11 and Fig. 4.6, present valid solutions to the inverse problem. Hence, we can see that the variability of ozone in the upper atmosphere is large and that we do not capture the ozone peak around 80km. The posterior temperature profiles is similar to the prior profiles, as also seen in marginal posterior Fig. 4.20 to 4.24. We can already see that in the prior analysis, as the pressure temperature ratio does inherit the exponential structure of the pressure profile. So the posterior pressure profile is much more informative, see marginal for  $b$  in Fig. 4.24. So we can retrieve an informative pressure profile for the pressure but not for temperature.

Ideally we should do this iteratively update ozone and then temperature and pressure until proven convergence.

We could fix that by choosing a more restrictive prior for pressure but that would not be objective. Instead we should really work on a better model for ozone. We can get rid of the pressure skew at a SNR of  $\approx 10000$ .

### 5.3.1 Measurement Device

Then we can include more measurement specific details such as the pointing accuracy. Then we could sample measurement  $N_\Gamma$  geometries  $\Gamma^{(k)} \sim \pi(\Gamma)$  so that the posterior  $\pi(\mathbf{x}|\mathbf{y}) \approx 1/N_\Gamma \sum_\Gamma \pi(\mathbf{x}, \Gamma^{(k)}|\mathbf{y})$  and include other measurement device specific parameters.

### 5.3.2 Model

Since we have to truncate the full conditional at the end the model is not accurate enough to eliminate those values. This was to show that we can do a more comprehensive analysis compared to a regularised method. Ideally we like to use a more accurate model where we parametrise ozone, similar to the pressure and temperature profile. In doing so one would have to know much more about ozone in different altitudes. Then we possibly could employ a different graph Laplacian based on a different structure of ozone. And when we approximate the non-linear forward map with a affine map using a linear solver we could of course use other methods such as the machine learning methods.



## References

- [1] Bryan Duncan. *Aura at 20 Years*. <https://science.nasa.gov/science-research/earth-science/aura-at-20-years/>. [Online; accessed 31/08/25]. NASA's Goddard Space Flight Center (GSFC), 2024.
- [2] Susan L Ustin and Elizabeth McPhee Middleton. "Current and near-term Earth-observing environmental satellites, their missions, characteristics, instruments, and applications". In: *Sensors* 24.11 (2024), p. 3488.
- [3] Mallika Irene Suresh et al. "Multichannel upconversion of terahertz radiation in an optical disk resonator". In: *Opt. Express* 33.5 (2025), pp. 10302–10311.
- [4] Florian Sedlmeir et al. "Detecting THz in the telecom range: All resonant THz up-conversion in a whispering gallery mode resonator". In: *2014 39th International Conference on Infrared, Millimeter, and Terahertz waves (IRMMW-THz)*. 2014, pp. 1–2.
- [5] Clive D Rodgers. "Retrieval of atmospheric temperature and composition from remote measurements of thermal radiation". In: *Reviews of Geophysics* 14.4 (1976), pp. 609–624.
- [6] Nathaniel J. Livesey et al. *Earth Observing System (EOS) Microwave Limb Sounder (MLS) Version 5.0x Level 2 and 3 data quality and description document*. Version 5.0-1.1a. NASA Goddard Earth Sciences Data and Information Services Center, 2022.
- [7] Sze M Tan, Colin Fox, and Geoff K. Nicholls. *Course notes for ELEC 445 – Inverse Problems and Imaging*. <https://coursesupport.physics.otago.ac.nz/wiki/pmwiki.php/ELEC445/HomePage>. [Online; accessed 10/12/23]. Physics Department, University of Otago, 2016.
- [8] Nathaniel J Livesey et al. "Retrieval algorithms for the EOS Microwave limb sounder (MLS)". In: *IEEE Transactions on Geoscience and Remote Sensing* 44.5 (2006), pp. 1144–1155.
- [9] Schwartz M. et al. *MLS/Aura Level 2 Ozone (O<sub>3</sub>) Mixing Ratio V005*. [https://disc.gsfc.nasa.gov/datasets/ML203\\_005/summary?keywords=mls%20o3](https://disc.gsfc.nasa.gov/datasets/ML203_005/summary?keywords=mls%20o3). [Online; accessed 25/04/24]. NASA Goddard Earth Sciences Data and Information Services Center, 2020.
- [10] N. J. Livesey et al. "Validation of Aura Microwave Limb Sounder O<sub>3</sub> and CO observations in the upper troposphere and lower stratosphere". In: *Journal of Geophysical Research: Atmospheres* 113.D15 (2008).
- [11] L. Froidevaux et al. "Validation of Aura Microwave Limb Sounder stratospheric ozone measurements". In: *Journal of Geophysical Research: Atmospheres* 113.D15 (2008).
- [12] F. Werner et al. "Applying machine learning to improve the near-real-time products of the Aura Microwave Limb Sounder". In: *Atmospheric Measurement Techniques* 16.11 (2023), pp. 2733–2751.
- [13] Pasquale Sellitto et al. "Neural network algorithms for ozone profile retrieval from ESA-envisat sciamachy and NASA-AURA OMI satellite data". In: vol. 3. Aug. 2008, pp. III –170.
- [14] 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.

- [15] Gareth O. Roberts and Jeffrey S Rosenthal. “General state space Markov chains and MCMC algorithms”. In: *Probability Surveys* 1 (2004), pp. 20–71.
- [16] S. P. Meyn and R.L. Tweedie. *Markov Chains and Stochastic Stability. 2nd Edition*. New York: Cambridge University Press, 2009.
- [17] Havard Rue and Leonhard Held. *Gaussian Markov random fields: theory and applications*. London: CRC press, 2005.
- [18] Richard A Norton, J Andrés Christen, and Colin Fox. “Sampling hyperparameters in hierarchical models: improving on Gibbs for high-dimensional latent fields and large datasets”. In: *Communications in Statistics-Simulation and Computation* 47.9 (2018), pp. 2639–2655.
- [19] Charles W. Champ and Andrew V. Sills. “The Generalized Law of Total Covariance”. In: *preprint* (2022).
- [20] Charles J Geyer. “Practical markov chain monte carlo”. In: *Statistical science* (1992), pp. 473–483.
- [21] U. Wolff and B. Bunk. *Lecture Notes on Computational Physics II [in german]*. [www-com.physik.hu-berlin.de/comphys/comphys.htm](http://www-com.physik.hu-berlin.de/comphys/comphys.htm). [Online; accessed 29/08/25]. Humboldt University, Berlin, 2016.
- [22] A. Sokal. “Monte Carlo Methods in Statistical Mechanics: Foundations and New Algorithms”. In: *Functional Integration: Basics and Applications*. Ed. by Cecile DeWitt-Morette, Pierre Cartier, and Antoine Folacci. Boston, MA: Springer US, 1997, pp. 131–192.
- [23] Ulli Wolff. “Monte Carlo errors with less errors”. In: *Computer Physics Communications* 156.2 (2004), pp. 143–153.
- [24] Daniel Simpson, Finn Lindgren, and Håvard Rue. “Think continuous: Markovian Gaussian models in spatial statistics”. In: *Spatial Statistics* 1 (2012), pp. 16–29.
- [25] Gareth O. Roberts and Jeffrey S Rosenthal. “Harris recurrence of Metropolis-within-Gibbs and trans-dimensional Markov chains”. In: *The Annals of Applied Probability* 16.4 (2006), 2123–2139.
- [26] J. Andrés Christen and Colin Fox. “A general purpose sampling algorithm for continuous distributions (the t-walk)”. In: *Bayesian Analysis* 5.2 (2010), pp. 263 –281.
- [27] Johnathan M Bardsley. “MCMC-based image reconstruction with uncertainty quantification”. In: *SIAM Journal on Scientific Computing* 34.3 (2012), A1316–A1332.
- [28] Johnathan M Bardsley et al. “Randomize-then-optimize: A method for sampling from posterior distributions in nonlinear inverse problems”. In: *SIAM Journal on Scientific Computing* 36.4 (2014), A1895–A1910.
- [29] Dootika Vats et al. “Understanding Linchpin Variables in Markov Chain Monte Carlo”. In: *preprint* (2022). URL: <https://arxiv.org/pdf/2210.13574.pdf>.
- [30] Tiangang Cui and Sergey Dolgov. “Deep composition of tensor-trains using squared inverse rosenblatt transports”. In: *Foundations of Computational Mathematics* 22.6 (2022), pp. 1863–1922.
- [31] Philip J Davis and Philip Rabinowitz. *Methods of numerical integration*. San Diego, CA: Academic Press, Inc., 1984.
- [32] Sergey Dolgov et al. “Approximation and sampling of multivariate probability distributions in the tensor train decomposition”. In: *Statistics and Computing* 30 (2020), pp. 603–625.
- [33] Colin Fox et al. “Grid methods for Bayes-optimal continuous-discrete filtering and utilizing a functional tensor train representation”. In: *Inverse Problems in Science and Engineering* 29.8 (2021), pp. 1199–1217.

- [34] Ivan V Oseledets. “Tensor-train decomposition”. In: *SIAM Journal on Scientific Computing* 33.5 (2011), pp. 2295–2317.
- [35] Ivan Oseledets. “DMRG Approach to Fast Linear Algebra in the TT-Format”. In: *Computational Methods in Applied Mathematics* 11.3 (2011), pp. 382–393.
- [36] John Thickstun. *Kantorovich-rubinstein duality*. [https://courses.cs.washington.edu/courses/cse599i/20au/resources/L12\\_duality.pdf](https://courses.cs.washington.edu/courses/cse599i/20au/resources/L12_duality.pdf). [Online; accessed 31/08/25]. University of Washington, 2019.
- [37] Luigi Ambrosio, Elia Brué, and Daniele Semola. *Lectures on Optimal Transport*. Cham: Springer Nature Switzerland, 2024.
- [38] Jean Feydy. “Analyse de données géométriques, au delà des convolutions”. Thesis. Université Paris-Saclay, July 2020. URL: <https://theses.hal.science/tel-02945979>.
- [39] Marcel Berger. *Geometry I. 4th Edition*. Berlin Heidelberg: Springer-Verlag, 2009.
- [40] Katsumi Nomizu and Takeshi Sasaki. *Affine differential geometry*. Cambridge: Cambridge University Press, 1994.
- [41] Per Christian Hansen. “Regularization, GSVD and truncated GSVD”. In: *BIT numerical mathematics* 29.3 (1989), pp. 491–504.
- [42] Per Christian Hansen. “The L-Curve and its Use in the Numerical Treatment of Inverse Problems”. English. In: *Computational Inverse Problems in Electrocardiology*. Ed. by P. Johnston. WIT Press, 2001, pp. 119–142.
- [43] Yu-Xiang Wang et al. “Trend Filtering on Graphs”. In: *Journal of Machine Learning Research* 17.105 (2016), pp. 1–41.
- [44] Per Christian Hansen and Dianne Prost O’Leary. “The use of the L-curve in the regularization of discrete ill-posed problems”. In: *SIAM Journal on Scientific Computing* 14.6 (1993), pp. 1487–1503.
- [45] KC Santosh, Nibaran Das, and Swarnendu Ghosh. “Chapter 3 - Deep learning models”. In: *Deep Learning Models for Medical Imaging*. Ed. by KC Santosh, Nibaran Das, and Swarnendu Ghosh. Primers in Biomedical Imaging Devices and Systems. Academic Press, 2022, pp. 65–97.
- [46] C. Readings. *Envisat MIPAS An Instrument for Atmospheric Chemistry and Climate Research*. Noordwijk: ESA Publications Division, 2000.
- [47] Jae N. Lee and Dong L. Wu. “Solar Cycle Modulation of Nighttime Ozone Near the Mesopause as Observed by MLS”. In: *Earth and Space Science* 7.4 (2020).
- [48] Iouli E Gordon et al. “The HITRAN2020 molecular spectroscopic database”. In: *Journal of Quantitative Spectroscopy and Radiative Transfer* 277 (2022), p. 107949.
- [49] Marie Šimečková et al. “Einstein A-coefficients and statistical weights for molecular absorption transitions in the HITRAN database”. In: *Journal of Quantitative Spectroscopy and Radiative Transfer* 98.1 (2006), pp. 130–155.
- [50] George B. Rybicki and Alan P. Lightman. *Radiative processes in Astrophysics*. Weinheim: Wiley-VCH, 2004.
- [51] W.G. Read et al. “The clear-sky unpolarized forward model for the EOS aura microwave limb sounder (MLS)”. In: *IEEE Transactions on Geoscience and Remote Sensing* 44.5 (2006), pp. 1367–1379.
- [52] Colin Fox. *Blokkurs on computing MCMC for inverse problems*. unpublished. Physics Department, University of Otago, 2025.
- [53] Australian National Concurrent Design Facility. *CubeSat Microwave Radiometer Mission to Support Global Ozone Layer Monitoring. Concept Study - Summary Report*. unpublished, internal report. Canberra BC: UNSW Canberra Space, 2023.

- [54] Joe W Waters et al. "The earth observing system microwave limb sounder (EOS MLS) on the Aura satellite". In: *IEEE Transactions on Geoscience and Remote Sensing* 44.5 (2006), pp. 1075–1092.
- [55] *U.S. Standard Atmosphere, 1976*. Washington, D.C.: United States. National Oceanic and Atmospheric Administration, United States Committee on Extension to the Standard Atmosphere, 1976.
- [56] Gareth Roberts. *ST911 Fundamentals of Statistical Inference - Part III*. Department of Statistics, University of Warwick, 2015.
- [57] Drik Hesse. *py-uwerr; Python implementation of Monte Carlo error analysis a la Wolff*. <https://github.com/dhesse/py-uwerr>. [Online; accessed 09/09/25].
- [58] Ivan Oseledets and Eugene Tyrtyshnikov. "TT-cross approximation for multidimensional arrays". In: *Linear Algebra and its Applications* 432.1 (2010), pp. 70–88.
- [59] Sergey Dolgov and Robert Scheichl. "A Hybrid Alternating Least Squares–TT-Cross Algorithm for Parametric PDEs". In: *SIAM/ASA Journal on Uncertainty Quantification* 7.1 (2019), pp. 260–291.
- [60] I. V. Oseledets et al. *tpty - a Python implementation of the TT-toolbox*. <https://github.com/oseledets/tpty>. [accessed 23/06/25]. 2018.
- [61] Josef Dick, Frances Y. Kuo, and Ian H. Sloan. "High-dimensional integration: The quasi-Monte Carlo way". In: *Acta Numerica* 22 (2013), 133–288.
- [62] 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.
- [63] Per Christian Hansen. *Discrete Inverse Problems: Insight and Algorithms*. Philadelphia: SIAM, 2010.
- [64] Ville Satopää et al. "Finding a "Kneedle" in a Haystack: Detecting Knee Points in System Behavior". In: *2011 31st International Conference on Distributed Computing Systems Workshops*. IEEE. 2011, pp. 166–171.
- [65] Ulli Wolff. *UWerr.m Version6*. <https://www.physik.hu-berlin.de/de/com/ALPHAssoft>. [Online; accessed 5/11/23]. Humboldt-Universität to Berlin, 2004.
- [66] J. Andrés Christen and Colin Fox. *The t-walk software*. <https://www.cimat.mx/~jac/twalk/>. [Online; accessed 25/11/24]. CIMAT, Mexico, and University of Otago, New Zealand.
- [67] Jean Feydy. *GeomLoss – Geometric Loss functions between sampled measures, images and volumes*. <https://www.kernel-operations.io/geomloss/api/pytorch-api.html>. [Online; accessed 12/09/25].
- [68] Paul B. Rohrbach et al. "Rank Bounds for Approximating Gaussian Densities in the Tensor-Train Format". In: *SIAM/ASA Journal on Uncertainty Quantification* 10.3 (2022), pp. 1191–1224.
- [69] M. Capiński and P.E. Kopp. *Measure, Integral and Probability. Springer Undergraduate Mathematics Series*. London: Springer-Verlag London, 2004.
- [70] M. Simonnet. *Measures and Probabilities*. New York: Springer-Verlag, 1996.
- [71] Vesa Kaarnioja. *Inverse Problems. Eighth lecture*. <https://vesak90.userpage.fu-berlin.de/ip23/week8.pdf>. [Online; accessed 10/04/25]. Freie Universität Berlin, Department of Mathematics and Computer Science, 2023.
- [72] Greg Lawler. *Notes on probability*. <https://www.math.uchicago.edu/~lawler/probnotes.pdf>. [Online; accessed 10/04/25]. The University of Chicago, Department of Mathematics, 2016.

# Appendices



# A

## Theoretical and technical background

### A.1 Correlation Structure

In the book Gaussian Markov Random Fields [17], Rue and Held demonstrate that a strong correlation between the hyper-parameter  $\mu$  and the latent field  $\mathbf{x}$  can significantly slow down convergence when using samplers, in particular Gibbs samplers. They consider the hierarchical model

$$\mu \sim \mathcal{N}(0, 1) \quad (\text{A.1a})$$

$$\mathbf{x}|\mu \sim \mathcal{N}(\mu \mathbf{1}, \mathbf{Q}^{-1}), \quad (\text{A.1b})$$

and apply a Gibbs sampler based on the full conditional distributions

$$\mu^{(k)} | \mathbf{x}^{(k)} \sim \mathcal{N}\left(\frac{\mathbf{1}^T \mathbf{Q} \mathbf{x}^{(k-1)}}{1 + \mathbf{1}^T \mathbf{Q} \mathbf{1}}, \left(1 + \mathbf{1}^T \mathbf{Q} \mathbf{1}\right)^{-1}\right) \quad (\text{A.2})$$

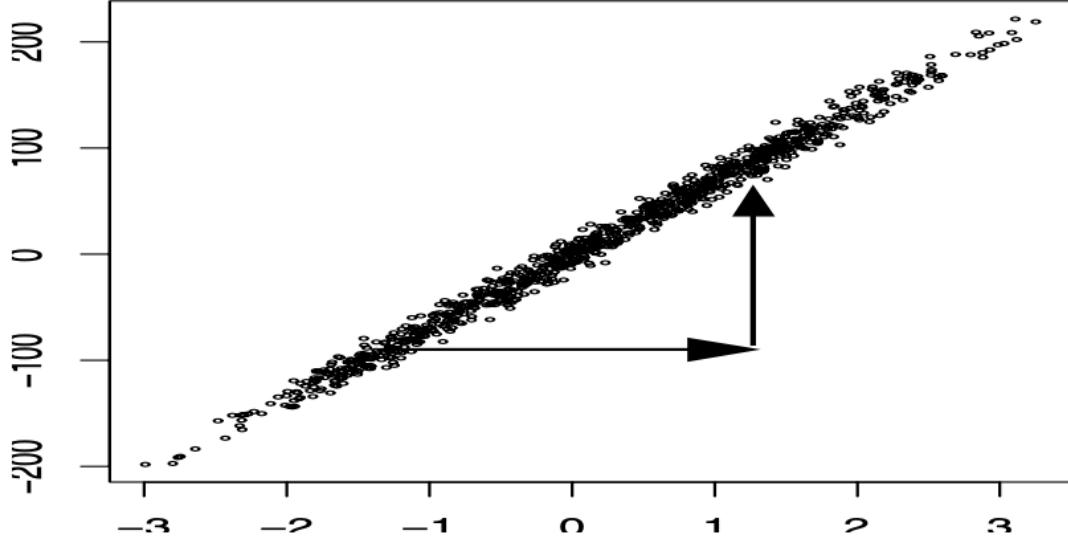
$$\mathbf{x}^{(k)} | \mu^{(k)} \sim \mathcal{N}(\mu^{(k)} \mathbf{1}, \mathbf{Q}^{-1}). \quad (\text{A.3})$$

As illustrated in Figure A.1, when the sampler is restricted to steps only in the  $\mu$ -direction (horizontal axis) or the  $\mathbf{x}$ -direction (vertical axis), it requires many iterations to adequately explore the parameter space. This inefficiency arises from the high correlation between  $\mu$  and  $\mathbf{x}$ , visible in Figure A.1 as a 'squeeze' of the distribution.

A solution to the slow mixing problem is to update  $(\mu, \mathbf{x})$  jointly. Since here  $\mu$  is one dimensional, effectively only marginal density of  $\mu$  is needed.

$$\mu^* \sim q(\mu^* | \mu^{(k-1)}) \quad (\text{A.4})$$

$$\mathbf{x}^{(k)} | \mu^* \sim \mathcal{N}(\mu^* \mathbf{1}, \mathbf{Q}^{-1}) \quad (\text{A.5})$$



**Figure A.1:** The figure taken from [17, Figure 4.1 (b)], shows samples from a marginal chain for  $\mu$  and  $\mathbf{1}^T \mathbf{Q} \mathbf{x}^{(k)}$  over 1000 iterations, based on the hierarchical model in Eq. A.1, with an autoregressive process encoded in  $\mathbf{Q}$ . The algorithm updates  $\mu$  and  $\mathbf{x}$  successively from their full conditional distributions. The plot displays  $(\mu^{(k)}, \mathbf{1}^T \mathbf{Q} \mathbf{x}^{(k)})$ , with  $\mu^{(k)}$  on the horizontal axis and  $\mathbf{1}^T \mathbf{Q} \mathbf{x}^{(k)}$  on the vertical axis. The slow mixing and convergence of  $\mu$  result from its strong dependence on  $\mathbf{1}^T \mathbf{Q} \mathbf{x}^{(k)}$ , while the sampler permits only axis-aligned (horizontal and vertical) and does not allow diagonal moves, as illustrated by the arrows.

With a simple MCMC algorithm targeting  $\mu$  one can explore the sample space efficiently and only draw a corresponding sample for  $\mathbf{x}$  from its full conditional once, for instance, the proposal  $\mu^*$  has been accepted.

## A.2 On the Monte-Carlo Error and Integrated Autocorrelation time

To assess the error  $\sigma^2$  of chain  $\mathcal{M}_i$ , we ignore systematic error due to initialisation bias (burn-in period), but we have to take into account that samples produced by any system or algorithm are correlated. To derive the integrated autocorrelation time (IATC), we follow the lecture notes [21]. In general, the error of a Monte-Carlo-based estimate from a sample set  $\mathcal{M}_i = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}, \dots, \mathbf{x}^{(s)}, \dots, \mathbf{x}^{(N)}\} \sim \pi(\mathbf{x}|\mathbf{y})$  is:

$$(\sigma^{(i)})^2 = \text{var}(\boldsymbol{\mu}_{\text{samp}}^{(i)}) = \text{var}(\mathbb{E}_{\mathbf{x}|\mathbf{y}}[h(\mathbf{x})]) = \left( \frac{1}{N} \sum_{k=1}^N h(\mathbf{x}^{(k)}) - \boldsymbol{\mu}^{(i)} \right)^2. \quad (\text{A.6})$$

Expanding this summation, we see that

$$(\sigma^{(i)})^2 = \frac{1}{N^2} \sum_{k,s=1}^N \Gamma(k-s) \quad (\text{A.7})$$

with the auto correlation coefficient  $\Gamma(k - s) = (h(\mathbf{x}^{(k)}) - \boldsymbol{\mu}^{(i)})(h(\mathbf{x}^{(s)}) - \boldsymbol{\mu}^{(i)})$ . Next we rewrite

$$\sum_{k,s=1}^N \Gamma(k - s) = \text{var}(\boldsymbol{\mu}_{\text{samp}}^{(i)}) \sum_{k,s=1}^N \frac{\Gamma(k - s)}{\Gamma(0)} = \text{var}(\boldsymbol{\mu}_{\text{samp}}^{(i)}) \sum_{k,s=1}^N \rho(k - s), \quad (\text{A.8})$$

with the normalised auto correlation coefficient  $\rho(k - s) = \Gamma(k - s)/\Gamma(0)$  at lag  $k - s$  and  $\Gamma(0) = \text{var}(\boldsymbol{\mu}_{\text{samp}}^{(i)})$  for  $k = s$ . Typically  $\Gamma(t)$  decays exponentially so that, for  $N \gg \tau$ ,  $\Gamma(t) \xrightarrow{t \rightarrow \infty} \exp\{-|t|/\tau\}$  and we can approximate

$$\sum_{k,s=1}^N \rho(k - s) = N \sum_{t=-(N-1)}^{N-1} \left(1 - \frac{|t|}{N}\right) \rho(t) \approx N \sum_{t=-\infty}^{\infty} \rho(t) := 2N\tau_{\text{int}}, \quad (\text{A.9})$$

see [22], and define the IATC as in [21, pp. 103-105]. If  $\tau \gg 1$  one can show that  $\tau_{\text{int}} \approx \tau$

$$\sum_{t=-\infty}^{\infty} \rho(t) = 1 + 2 \sum_{t=1}^{\infty} (e^{-1/\tau})^t = 1 + 2 \frac{e^{-1/\tau}}{1 - e^{-1/\tau}} \approx 1 + 2 \frac{1 - 1/\tau}{1/\tau} = 2\tau - 1 \approx 2\tau_{\text{int}} \quad (\text{A.10})$$

where we use the geometric power series  $\sum_{n=0}^{\infty} x^n = 1/(1+x)$  and the Taylor series  $e^x \approx 1 + x$  for small  $x$ . Consequently, the estimate for the Monte-Carlo error is:

$$(\sigma^{(i)})^2 \approx \frac{\text{var}(h(\mathbf{x}))}{N} \underbrace{\sum_{t=-\infty}^{\infty} \rho(t)}_{:= 2\tau_{\text{int}}} = \text{var}(h(\mathbf{x})) \frac{2\tau_{\text{int}}}{N}, \quad (\text{A.11})$$

where we define the IACT provides a good estimate of how many steps the sampling algorithm needs to take to produce one independent sample. More specifically, the effective sample size  $\frac{2\tau_{\text{int}}}{N}$  gives an estimate of how efficient a sampler is.

### A.3 Measure theory

Assume that the triple  $(\Omega, \mathcal{F}, \mathbb{P})$  defines a probability space, where  $\Omega$  denotes the complete sample space,  $\mathcal{F}$  is a  $\sigma$ -algebra consisting of a collection of countable subsets  $\{A_n\}_{n \in \mathbb{N}}$  with  $A_n \subseteq \Omega$ , and  $\mathbb{P}$  is a probability measure on  $\mathcal{F}$ . The formal conditions for  $\mathbb{P}$  to be a probability measure, and for  $\mathcal{F}$  to be a  $\sigma$ -algebra over  $\Omega$ , are given in Appendix A.3. We denote

$$\mathbb{P}(A) = \int_A d\mathbb{P} \quad (\text{A.12})$$

as the probability of an event  $A \in \mathcal{F}$ . By applying the Radon-Nikodym theorem [69], we can change variables

$$\mathbb{P}(A) = \int_A \frac{d\mathbb{P}}{d\mathbf{x}} d\mathbf{x} = \int_A \pi(\mathbf{x}) d\mathbf{x}, \quad (\text{A.13})$$

where  $d\mathbf{x}$  is a reference measure on the same probability space, commonly referred to as the Lebesgue measure. The Radon-Nikodym derivative  $\frac{d\mathbb{P}}{d\mathbf{x}}$  of  $\mathbb{P}$  with respect to  $\mathbf{x}$  is often interpreted as the probability density function (PDF)  $\pi(\mathbf{x})$ . Thus, we say that  $\mathbb{P}$  has a density  $\pi(\mathbf{x})$  with respect to  $\mathbf{x}$  [70, Chapter 10].

Now, let  $X : \Omega \rightarrow \mathbb{R}^d$  be a  $d$ -dimensional random variable mapping from the probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  to the measurable space  $(\mathbb{R}^d, \mathcal{X})$ , where  $\mathcal{X}$  is a collection of subsets in  $\mathbb{R}^d$  [71]. Then the associated PDF  $\pi(\mathbf{x})$  is a joint density of  $X$ , induced by the probability measure on  $\Omega$  [69, 71].

Recall the probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ , where  $\Omega$  denotes the sample space, and  $\mathcal{F}$  is a collection of countable subsets  $\{A_n\}_{n \in \mathbb{N}}$  of  $\Omega$ . Each  $A_n \subseteq \Omega$  is called an event, and a map  $\mathbb{P} : \mathcal{F} \rightarrow \mathbb{R}$  is referred to as a measure. In the following, we describe the conditions required for  $\mathcal{F}$  to be a  $\sigma$ -algebra, and for  $\mathbb{P}$  to qualify as a probability measure. We refer to [72] [69] for further reading.

### A.3.1 Probability Measure

For a probability measure, we require:

- $\mathbb{P}(\Omega) = 1$  and  $\mathbb{P}(\emptyset) = 0$
- $\mathbb{P}(A) \in [0, 1]$
- $\mathbb{P}(\bigcup_{j \in \mathbb{N}} A_j) = \sum_{j \in \mathbb{N}} \mathbb{P}(A_j)$  if we have pairwise disjoint sets or  $A_i \cap A_j = \emptyset$  for  $i \neq j$

In other words, the probability assigned to the entire sample space must be equal to one,  $\mathbb{P}(\Omega) = 1$ , and the probability of the empty set must be zero,  $\mathbb{P}(\emptyset) = 0$ . For any subset  $A \subseteq \Omega$ , the probability  $\mathbb{P}(A)$  must lie between zero and one, i.e.,  $\mathbb{P}(A) \in [0, 1]$ . If e.g. two subsets  $A$  and  $B$  are disjoint (i.e.,  $A \cap B = \emptyset$ ), then the probability of their union satisfies  $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B)$ . This property must also hold for a countable sequence of disjoint sets  $\{A_j\}_{j \in \mathbb{N}}$ , such that  $\mathbb{P}\left(\bigcup_{j \in \mathbb{N}} A_j\right) = \sum_{j \in \mathbb{N}} \mathbb{P}(A_j)$ .

### A.3.2 $\sigma$ -Algebra

A collections of subsets  $\mathcal{F}$  is called  $\sigma$ -algebra if:

- $\emptyset, \Omega \in \mathcal{F}$ ,
- if  $A \in \mathcal{F}$  then  $A^C := A/\Omega \in \mathcal{F}$
- if  $A_1, A_2, \dots \in \mathcal{F}$  then  $\bigcup_{j \in \mathbb{N}} A_j \in \mathcal{F}$

In other words, the empty set  $\emptyset$  and the entire sample space  $\Omega$  must always be elements of  $\mathcal{F}$ . If a set  $A \in \mathcal{F}$ , then its complement  $A^C = \Omega \setminus A$  must also be in  $\mathcal{F}$ . If, in terms of a probability measure, we are able to assign a probability  $\mathbb{P}(A)$  to an event  $A$ , we must also be able to assign a probability to the event “not  $A$ ”, i.e.,  $\mathbb{P}(A^C)$ . Finally, if a countable collection of sets  $A_1, A_2, \dots \in \mathcal{F}$ , then their union  $\bigcup_{j \in \mathbb{N}} A_j$  must also be in  $\mathcal{F}$ . These three properties define the requirements for  $\mathcal{F}$  to be a  $\sigma$ -algebra.

## A.4 Python Code

18.0pt plus 2.0pt minus 1.0pt

```

1 def MargBack(TTCore, univarGrid):
2     ''' Backward marginalisation, see Prop. 1 as in SIRT from Cui et al. [30] '''
3
4     dim = len(univarGrid)
5     B = dim * [None] # coeffTensor
6     B[-1] = TTCore[-1]
7     R = [None] * dim
8     C = [None] * dim
9
10    for k in range(dim - 1, 0, -1):
11        r_kmin1, n, r_k = np.shape(TTCore[k])
12        # Eq. , [30, Eq. 22] !! we set Lebesgue Measure to const = one
13        M = np.identity(n) * (univarGrid[k][-1] - univarGrid[k][0]) # Mass matrix
14        L = scy.linalg.cholesky(M)
15
16        # construct Tensor C Eq. (27)
17        C[k] = np.zeros((r_kmin1, n, r_k))
18        for alpha in range(0, r_kmin1):
19            for l in range(0, r_k):
20                C[k][alpha, :, l] = B[k][alpha, :, l] @ L[:, :]
21
22        # unfold along first coordinate and compute thin QR decomposition of C^T
23        # Eq. (28)
24        Q, R[k] = np.linalg.qr(C[k].reshape((r_kmin1, n * r_k), order='C').transpose(),
25                               mode='reduced')
26
27        # compute next coefficient tensor
28        # Eq. (29)
29        r_kmin2, n, r_kmin1 = np.shape(TTCore[k - 1])
30        B[k - 1] = np.zeros(np.shape(TTCore[k - 1]))
31        for alpha_2 in range(0, r_kmin2):
32            #for i in range(0, n):
33            for l_1 in range(0, r_kmin1):
34                B[k - 1][alpha_2, :, l_1] = TTCore[k - 1][alpha_2, :, :] @ R[k][l_1, :]
35
36    return B

```

**Listing A.1:** Pyhton code to calculate Backward marginals, as in Prop. 1 and [30].

```

1  def MargForw(TTCore, univarGrid):
2      ''' Forward marginalisation, see Prop.
3          2, similar to backward marginalisation as in Cui et al. [30] '''
4      # compute pre marginal coefficients sarting at dim = 1, k = 0
5      BPre = dim * [None] # coeffTensor
6      LebLam = 1 # !! Lebesgue Measure
7      BPre[0] = TTCore[0]
8      RPre = [None] * dim
9      CPre = [None] * dim
10
11     for k in range(0, dim-1):
12         r_kmin1, n, r_k = np.shape(TTCore[k])
13         # Eq. , [30, Eq. 22] !! we set Lebesgue Measure to const = one
14         M = np.identity(n) * (univarGrid[k][-1] - univarGrid[k][0]) # Mass matrix
15         L = scy.linalg.cholesky(M)
16
17         # construct Tensor C [30, Eq. (27)]
18         CPre[k] = np.zeros((r_kmin1, n, r_k))
19         for alpha in range(0, r_kmin1):
20             for l in range(0, r_k):
21                 CPre[k][alpha, :, l] = BPre[k][alpha, :, l] @ L[:, :]
22
23         # unfold along first coordinate and compute thin QR decomposition of C
24         # 3.1 [30, Eq. (28)]
25         Q, RPre[k] = np.linalg.qr(CPre[k].reshape((r_kmin1 * n, r_k)), order='C'), mode='reduced'
26
27         # compute next coefficient tensor [30, Eq. (29)]
28         r_k, n, r_kpls1 = np.shape(TTCore[k + 1])
29         BPre[k + 1] = np.zeros(np.shape(TTCore[k + 1]))
30         for alpha_1 in range(0, r_kpls1):
31             for l_1 in range(0, r_k):
32                 BPre[k + 1][l_1, :, alpha_1] = RPre[k][l_1, :] @ TTCore[k + 1][:, :, alpha_1]
33
34     return BPre

```

**Listing A.2:** Pyhton code to calculate forward marginals, as in Prop. 2.

```

1 def SIRT(seeds, SQTT, univarGrid, BackMarg, absError):
2     ''' do squared inverse rosenblatt transform (SIRT) as in Cui et al. [30] ''',
3     dim, numbSampl = seeds.shape
4     sampls = np.zeros(seeds.shape) # samples from approximated PDF
5     probVal = np.zeros(seeds.shape) # PDF values, for MH-correction step
6     Approx = np.zeros(seeds.shape[1]) # TT-Approx., to compare to true function
7
8     # lebesgue measure for quadrature Eq. 2.24
9     WholeLebLam = np.zeros(dim)
10    for k in range(0, dim):
11        WholeLebLam[k] = (univarGrid[k][-1] - univarGrid[k][0])
12    lamX = np.ones(dim)
13    for k in range(1, dim):
14        lamX[k - 1] = np.prod(WholeLebLam[k:])
15
16    # error as in Eq. 2.29 [30, Eq. (35)]
17    gamError = absError / np.prod(WholeLebLam)
18
19    # sample from first dimension [30, Eq. (30)]
20    firstMarg = gamError * lamX[0] + np.sum(BackMarg[0][0, :, :] ** 2, 1)
21    # cumulative distribution function, normalised numerically Eq. 2.43 [30, Eq. (17)]
22    firstCDF = np.cumsum(firstMarg / np.sum(firstMarg))
23    # draw samples as 'inverse transform'
24    sampls[0] = np.interp(seeds[0], firstCDF, univarGrid[0])
25    probVal[0] = np.interp(sampls[0], univarGrid[0], firstMarg / np.sum(firstMarg))
26
27    # sample from other dimensions
28    for n in range(0, numbSampl):
29        # interpolate linear on grid in first dimension Eq. 2.44 [32]
30        CurrApprCore = LinInterPolTT(SQTT[0], univarGrid[0], sampls[0][n])
31        for d in range(1, dim):
32            # marginal function condition on previous samples
33            rank_min, gridSize, rank_pls = BackMarg[d].shape
34            MargDep = np.zeros((BackMarg[d].shape))
35            for r in range(0, rank_min):
36                # condition on previous samples
37                MargDep[r, :, :] = CurrApprCore[0, r] * BackMarg[d][r, :, :]
38
39            # Eq. 2.45 [30, Eq. (31)]
40            currMarg = gamError * lamX[d] + np.sum(np.sum(np.copy(MargDep), axis=0)** 2,
41                axis=1)
42            # Eq. 2.43 [30, Eq. (17)]
43            currCDF = np.cumsum(currMarg / np.sum(currMarg))
44
45            # draw sample as 'inverse transform'
46            sampls[d][n] = np.interp(seeds[d][n], currCDF, univarGrid[d])
47            probVal[d][n] = np.interp(sampls[d][n], univarGrid[d],
48                currMarg / np.sum(currMarg))
49            # piecew. poly. interpol., Eq. 2.44 [32], cond. on sampl. for next PDF
50            CurrApprCore = np.copy(CurrApprCore) @ LinInterPolTT(SQTT[d], univarGrid[d],
51                sampls[d][n])
52
53            Approx[n] = gamError + CurrApprCore ** 2
54
55    return sampls, probVal, Approx

```

**Listing A.3:** Pyhton code to draw samples via SIRT, as in Alg. Box 3.

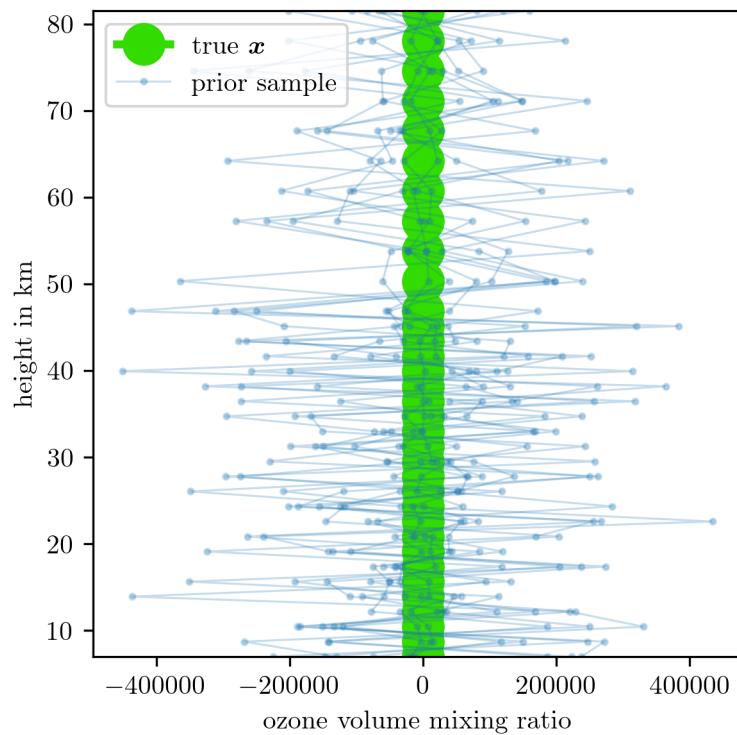
# B

## Additional Figures

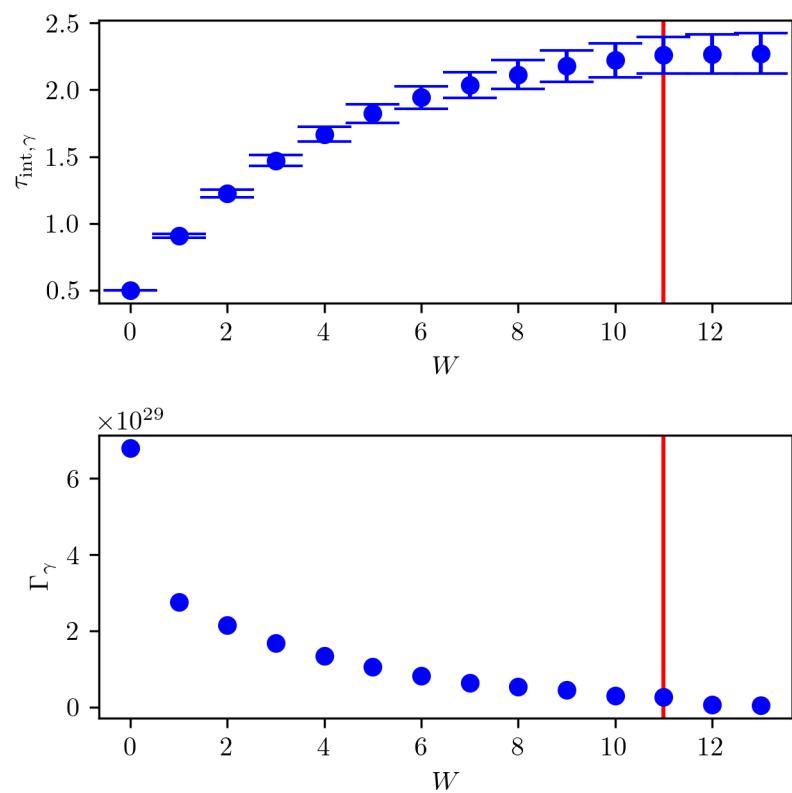
### B.1 Ozone

#### B.1.1 Ozone Prior

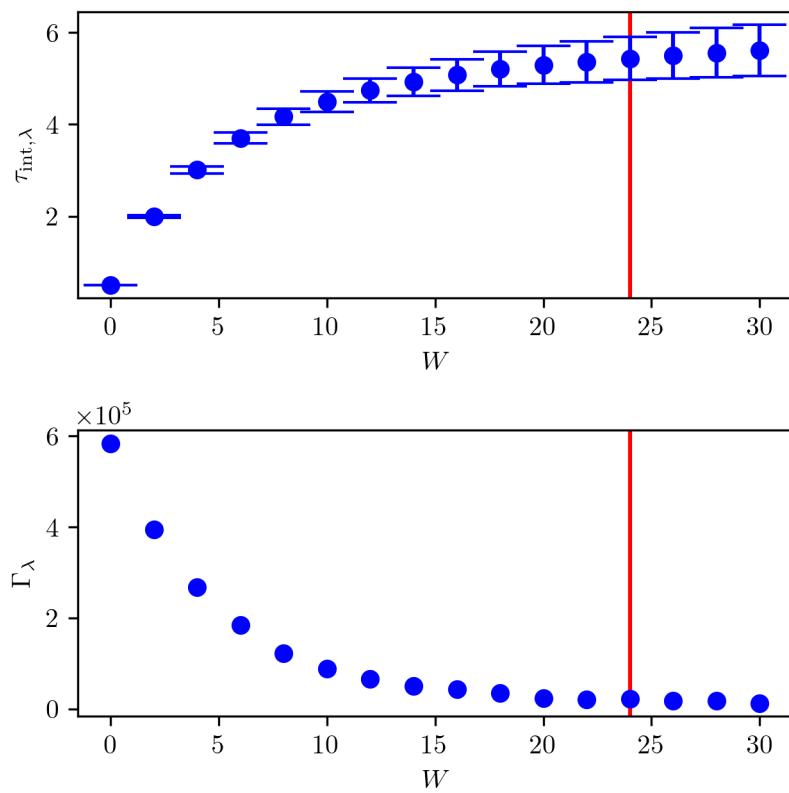
#### B.1.2 Integrated Autocorrelation plots



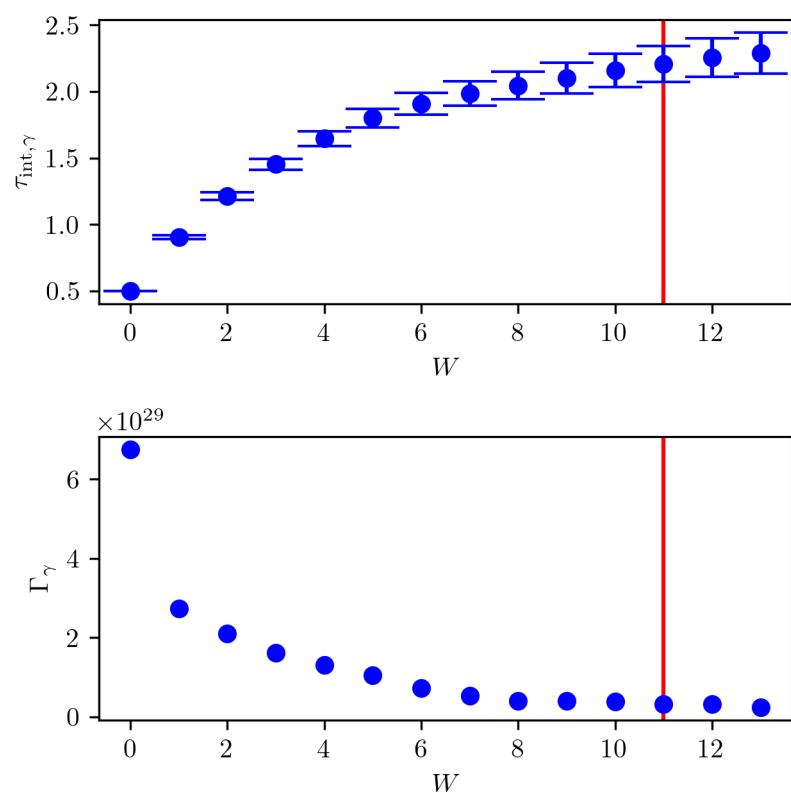
**Figure B.1:** We draw samples from ozone prior distribution  $\mathbf{x} \sim \mathcal{N}(0, \delta \mathbf{L})$  after generating a sample from the hyper-prior distribution  $\delta \sim \mathcal{T}(1, 10^{-10})$ . Note that since the variance of prior samples is very large compared to the ozone volume mixing ratios, the ozone profile appears to be constant, which is not the case, see e.g. Fig. 4.6.



**Figure B.2:** Here the autocorrelation function  $\Gamma_\gamma$  at different lags  $W$  is plotted as well as the IATC  $\tau_{\text{int},\gamma}$  for the samples from  $\pi(\gamma, \lambda | \mathbf{y})$  based on the linear forward model.



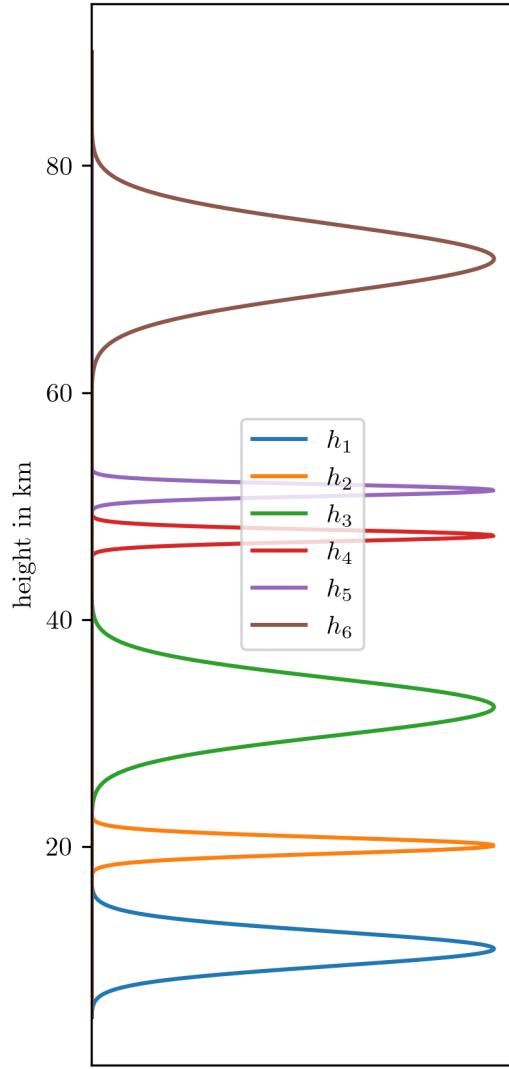
**Figure B.3:** IACT for samples  $\lambda \sim \pi(\cdot | \gamma, \mathbf{y})$  based on the approximated forward model.



**Figure B.4:** IACT for samples  $\gamma \sim \pi(\cdot | \lambda, \mathbf{y})$  based on the approximated forward model.

## B.2 Pressure and Temperature

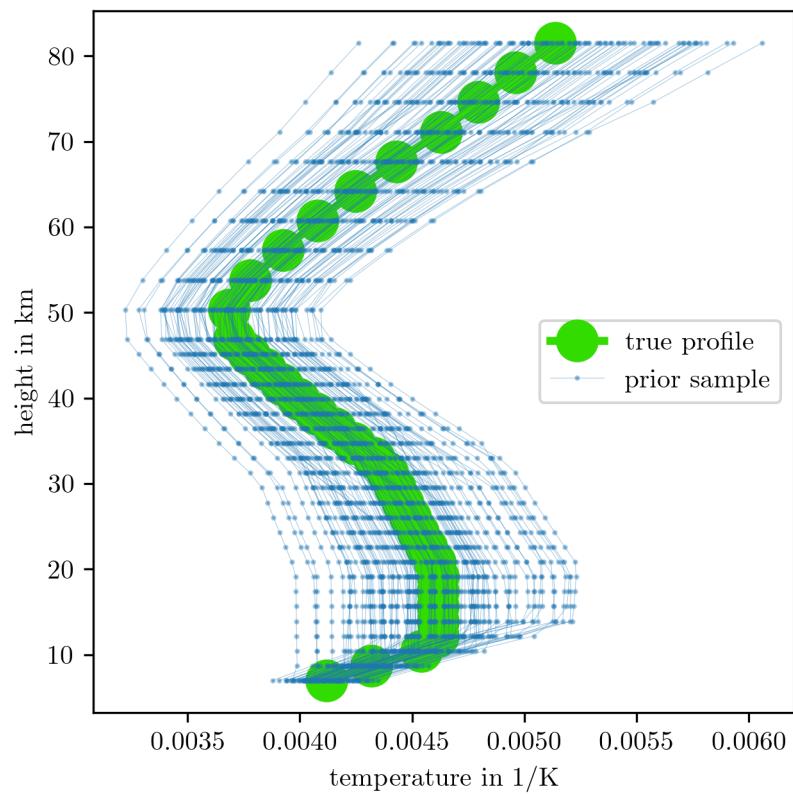
### B.2.1 Priors



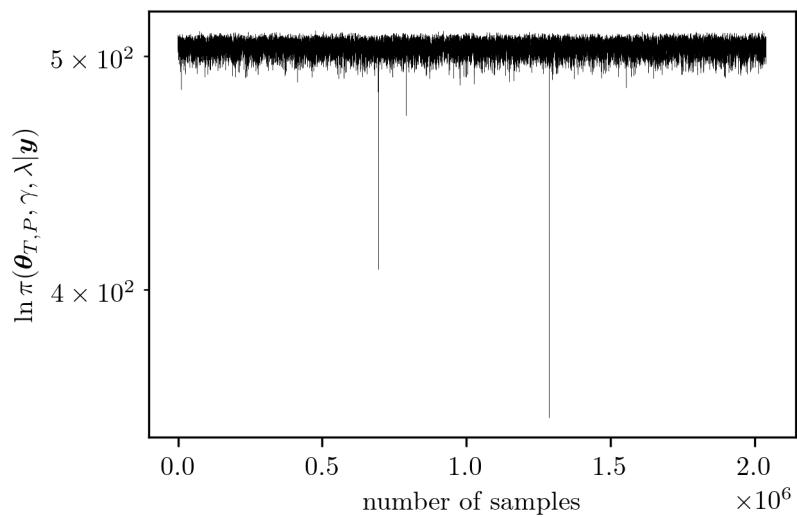
**Figure B.5:** Prior distributions  $\pi(\mathbf{h}_T)$ , which we choose so that they do not overlap and not conflict with the temperature function 4.3

### B.2.2 T-walk Trace

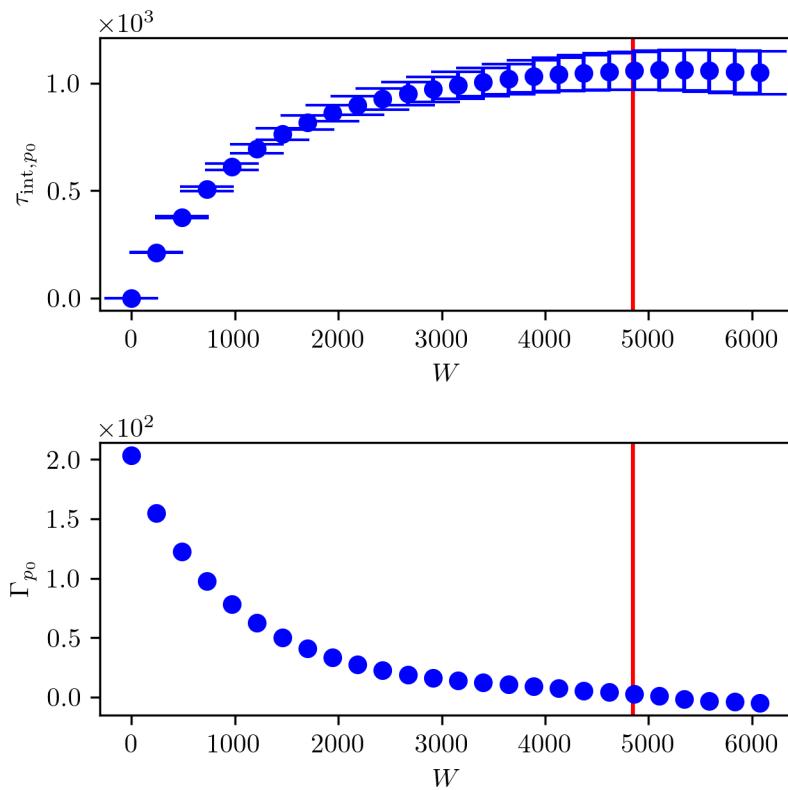
### B.2.3 Integrated Autocorrelation Plots



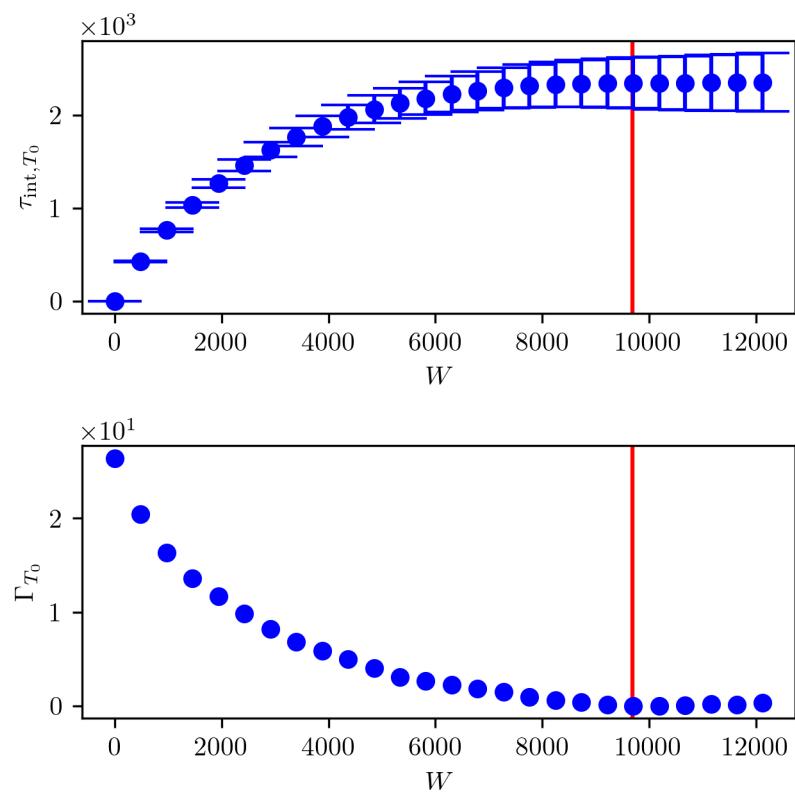
**Figure B.6:** Prior samples of the inverted temperature profile.



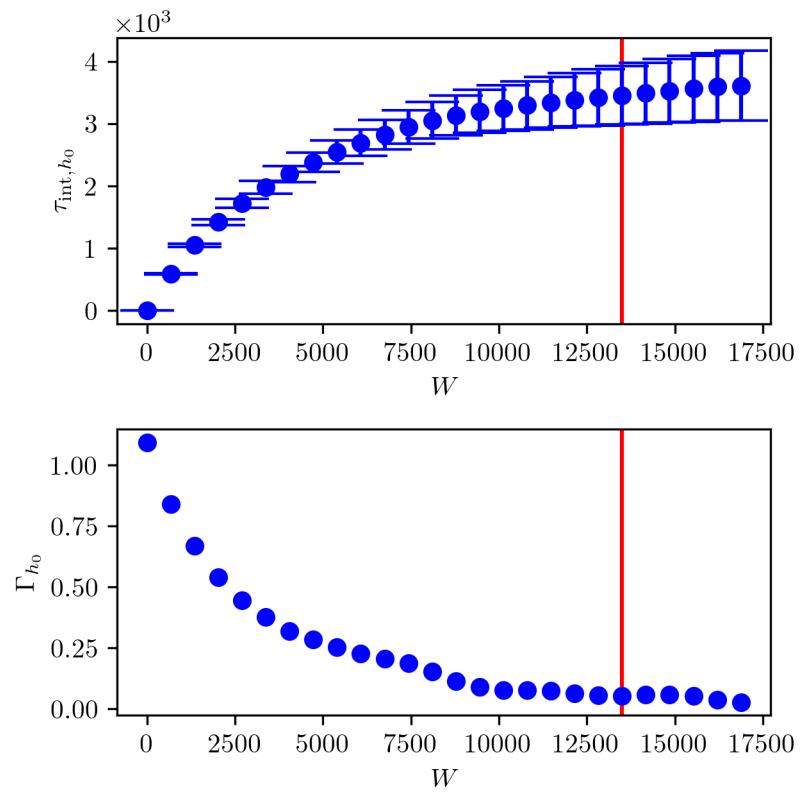
**Figure B.7:** Output trace of the t-walk on the posterior distribution  $\pi(p_0, b, \mathbf{h}_T, \mathbf{a} | \gamma, \mathbf{y})$ .



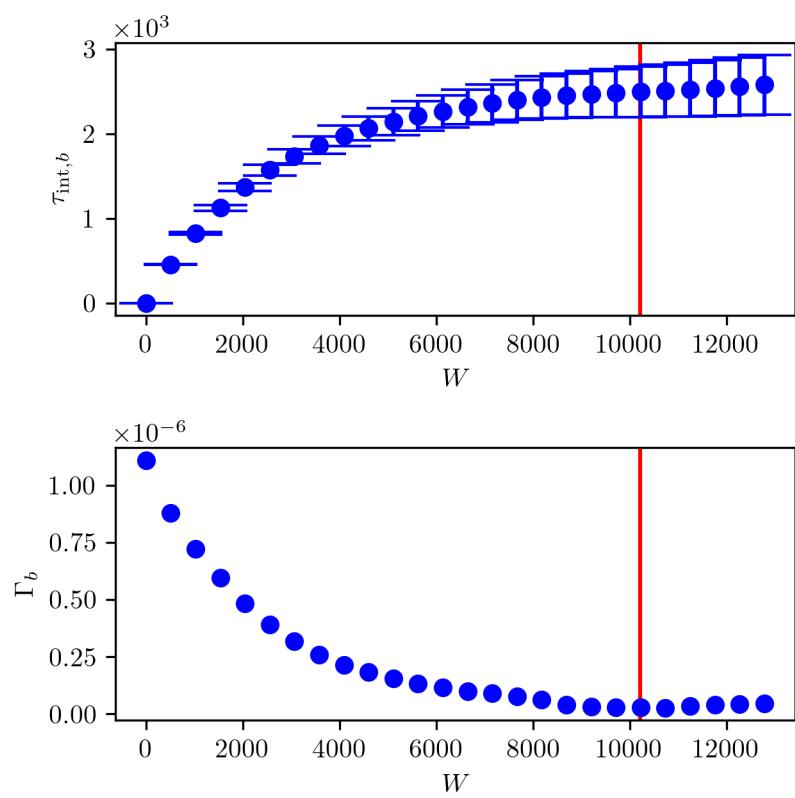
**Figure B.8:** IACT and autocorrelation function for samples  $h_1 \sim \pi(\cdot | h_{T,2}, h_{T,3}, h_{T,4}, h_{T,5}, h_{T,6}, a_0, a_1, a_2, a_3, a_4, a_5, a_6, T_0, b, p_0, \mathbf{y})$



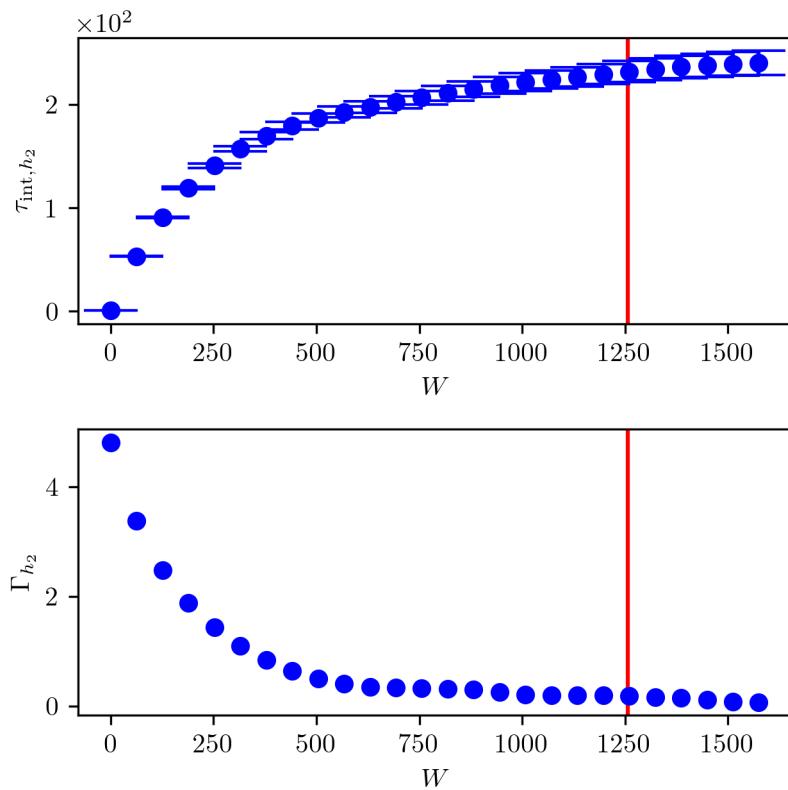
**Figure B.9:** IACT and autocorrelation function for samples  $h_2 \sim \pi(\cdot | h_1, h_3, h_4, h_5, h_6, a_0, a_1, a_2, a_3, a_4, a_5, a_6, T_0, b, p_0, \mathbf{y})$



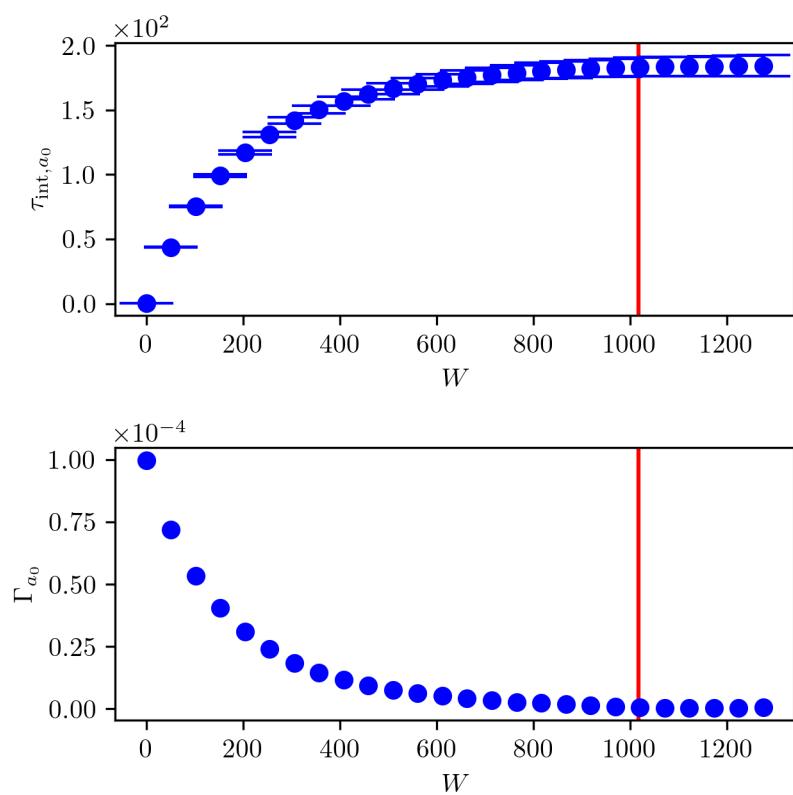
**Figure B.10:** IACT and autocorrelation function for samples  $h_3 \sim \pi(\cdot | h_1, h_2, h_4, h_5, h_6, a_0, a_1, a_2, a_3, a_4, a_5, a_6, T_0, b, p_0, y)$



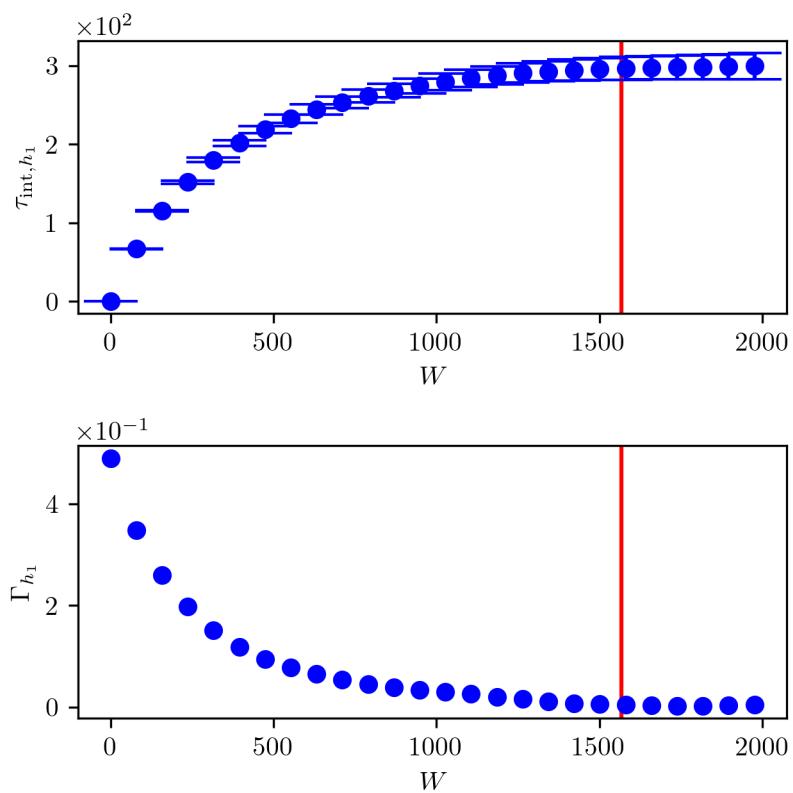
**Figure B.11:** IACT and autocorrelation function for samples  $h_4 \sim \pi(\cdot | h_1, h_2, h_3, h_5, h_6, a_0, a_1, a_2, a_3, a_4, a_5, a_6, T_0, b, p_0, \mathbf{y})$



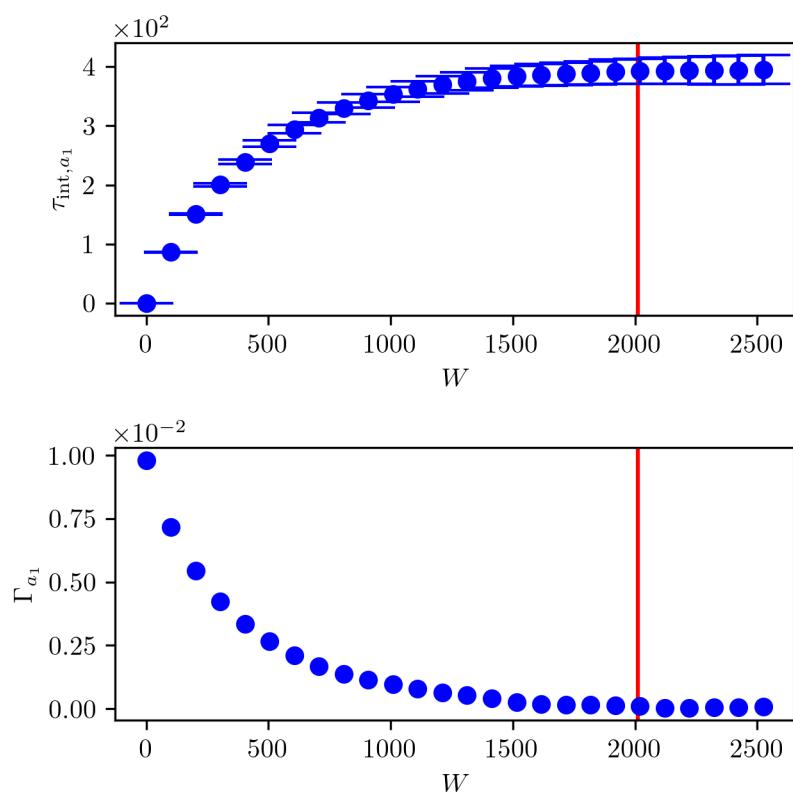
**Figure B.12:** IACT and autocorrelation function for samples  $h_5 \sim \pi(\cdot | h_1, h_2, h_3, h_4, h_6, a_0, a_1, a_2, a_3, a_4, a_5, a_6, T_0, b, p_0, y)$



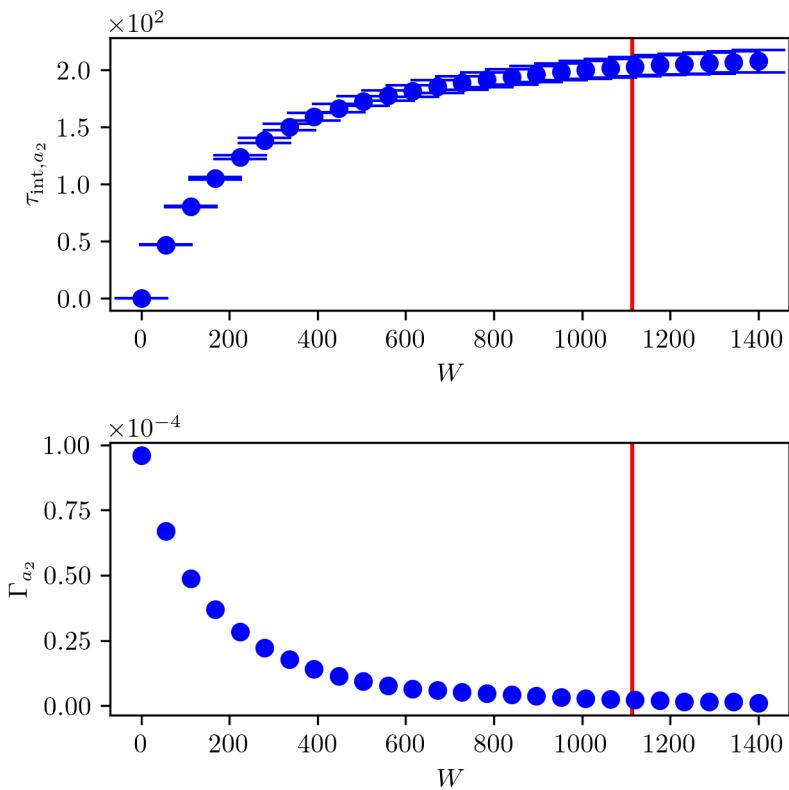
**Figure B.13:** IACT and autocorrelation function for samples  $h_6 \sim \pi(\cdot | h_1, h_2, h_3, h_4, h_5, a_0, a_1, a_2, a_3, a_4, a_5, a_6, T_0, b, p_0, \mathbf{y})$



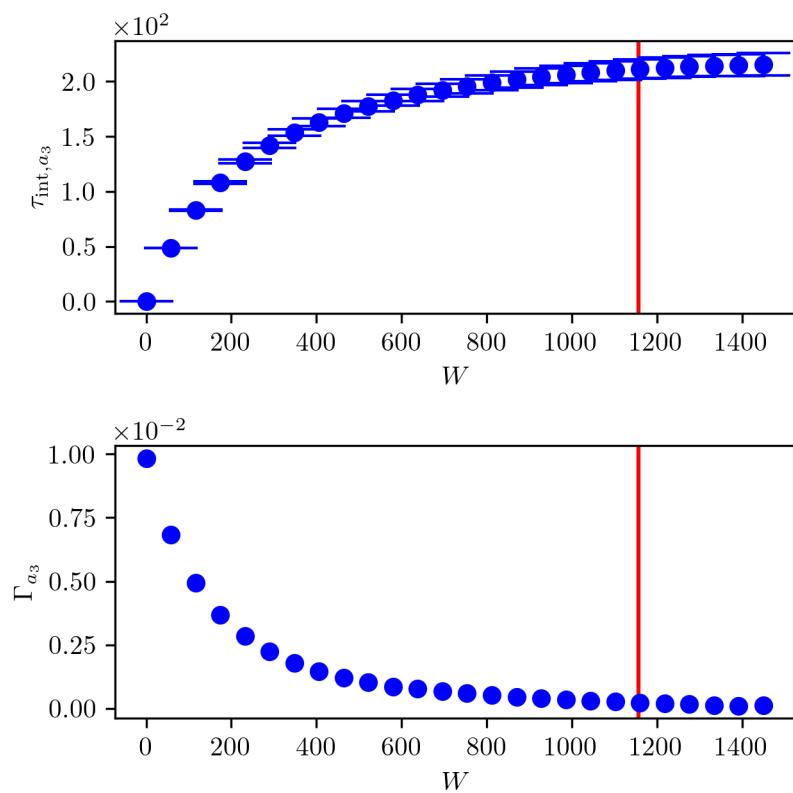
**Figure B.14:** IACT and autocorrelation function for samples  $a_0 \sim \pi(\cdot | h_1, h_2, h_3, h_4, h_5, h_6, a_1, a_2, a_3, a_4, a_5, a_6, T_0, b, p_0, \mathbf{y})$



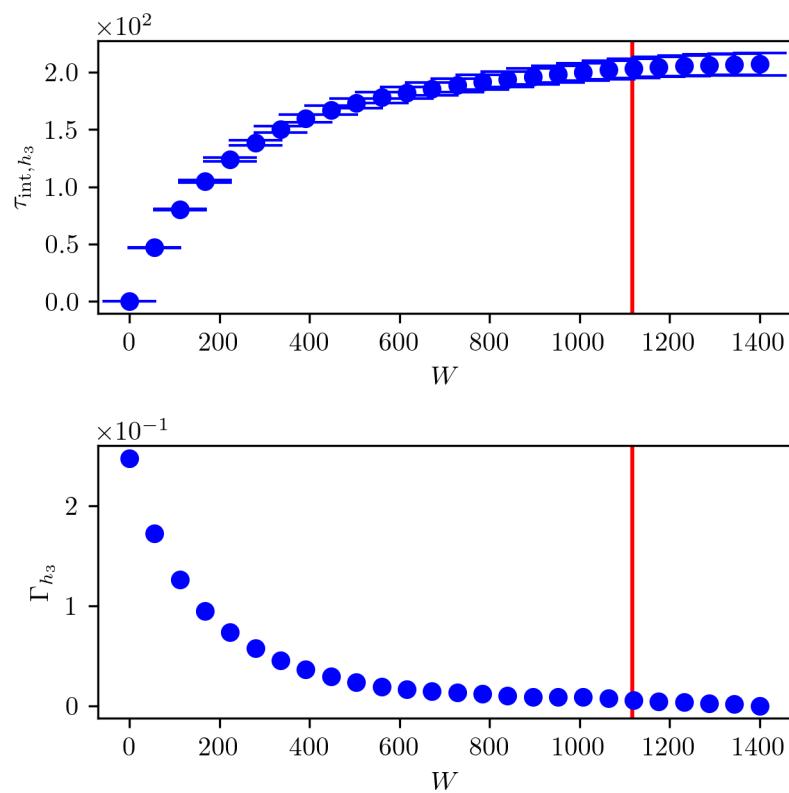
**Figure B.15:** IACT and autocorrelation function for samples  $a_1 \sim \pi(\cdot | h_1, h_2, h_3, h_4, h_5, h_6, a_0, a_2, a_3, a_4, a_5, a_6, T_0, b, p_0, \mathbf{y})$



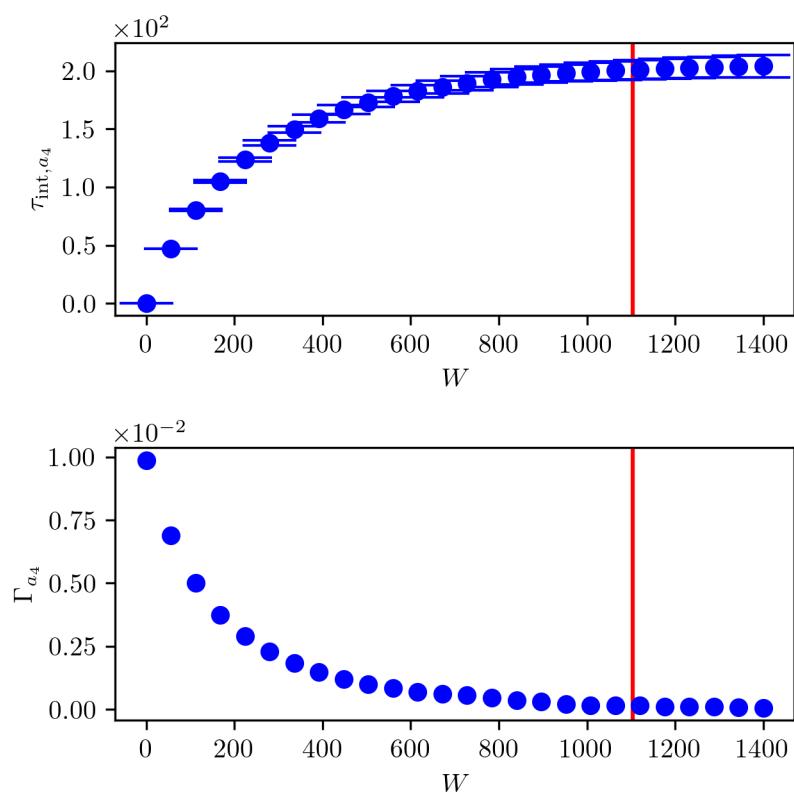
**Figure B.16:** IACT and autocorrelation function for samples  $a_2 \sim \pi(\cdot | h_1, h_2, h_3, h_4, h_5, h_6, a_0, a_1, a_3, a_4, a_5, a_6, T_0, b, p_0, \mathbf{y})$



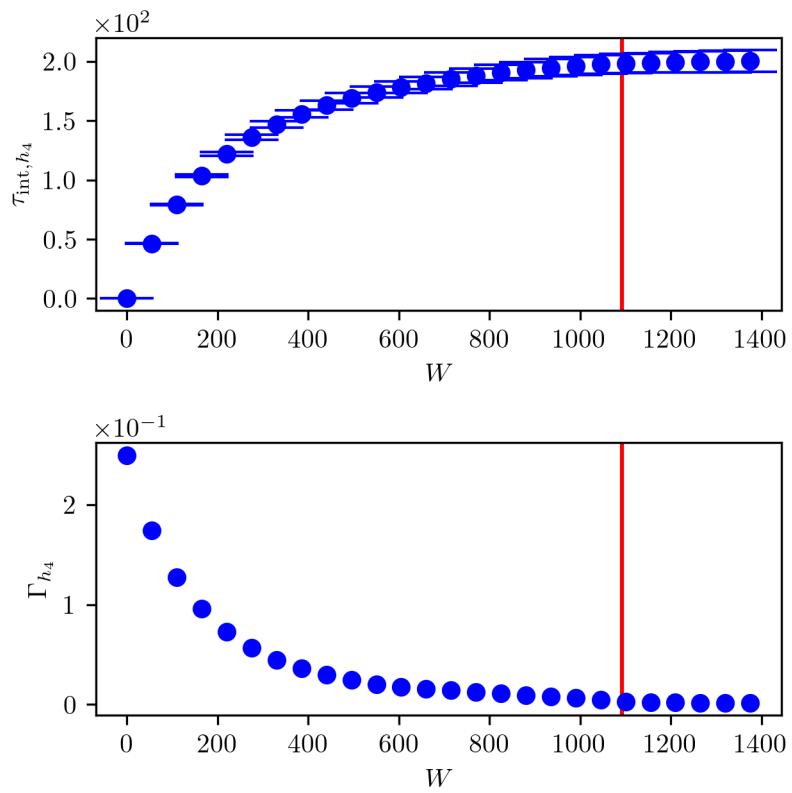
**Figure B.17:** IACT and autocorrelation function for samples  $a_3 \sim \pi(\cdot | h_1, h_2, h_3, h_4, h_5, h_6, a_0, a_1, a_2, a_4, a_5, a_6, T_0, b, p_0, \mathbf{y})$



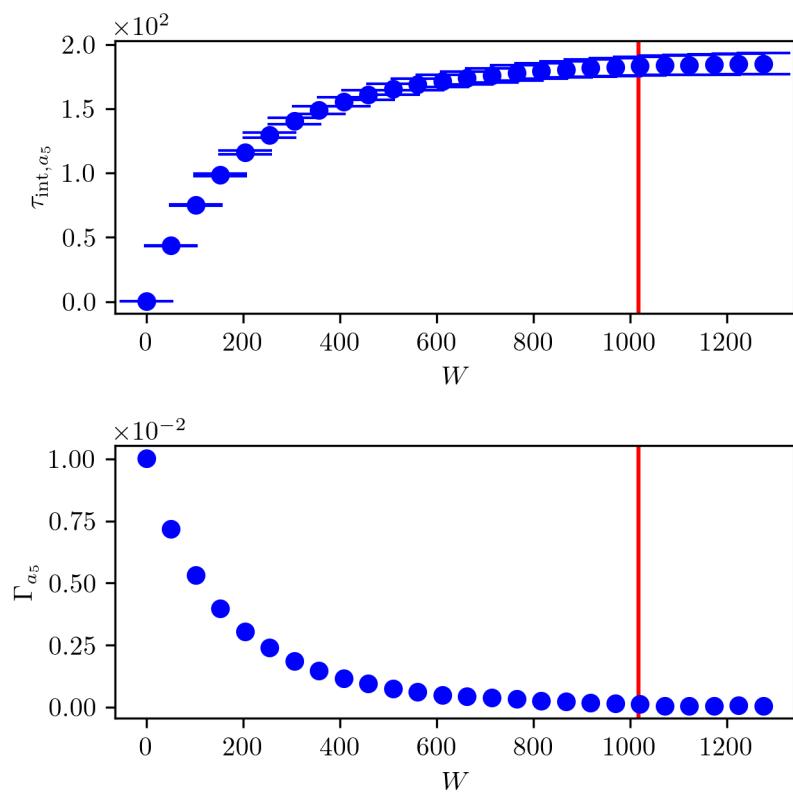
**Figure B.18:** IACT and autocorrelation function for samples  $a_4 \sim \pi(\cdot | h_1, h_2, h_3, h_4, h_5, h_6, a_0, a_1, a_2, a_4, a_5, a_6, T_0, b, p_0, \mathbf{y})$



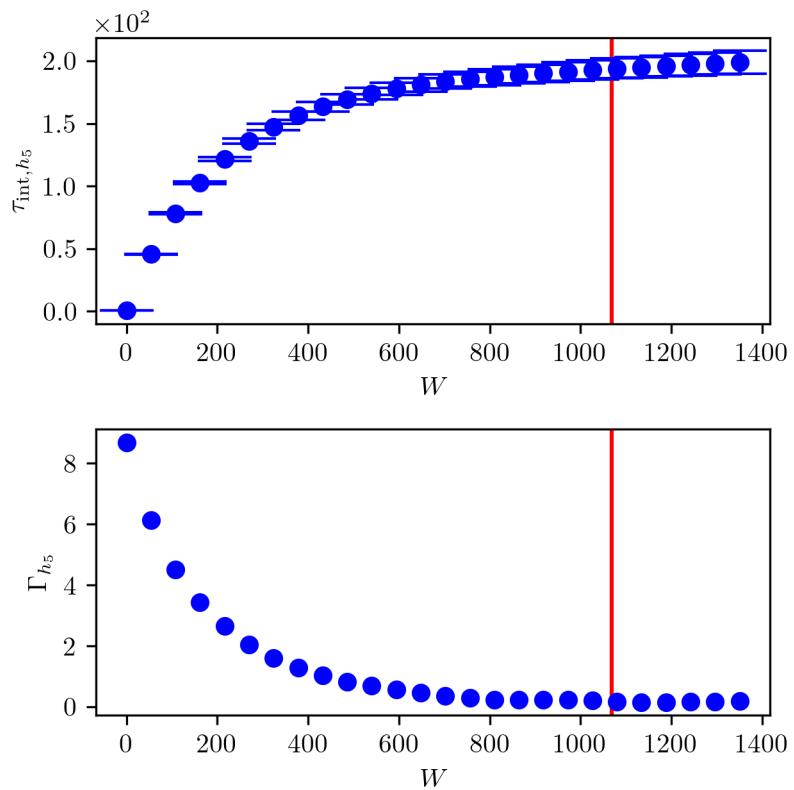
**Figure B.19:** IACT and autocorrelation function for samples  $a_5 \sim \pi(\cdot | h_1, h_2, h_3, h_4, h_5, h_6, a_0, a_1, a_2, a_3, a_4, a_5, T_0, b, p_0, \mathbf{y})$



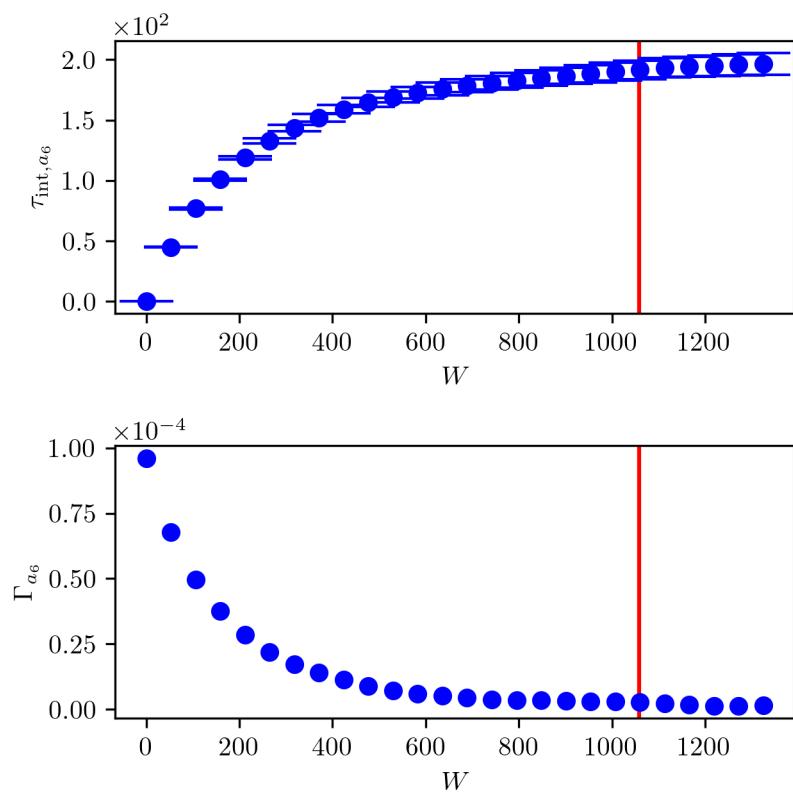
**Figure B.20:** IACT and autocorrelation function for samples  $a_6 \sim \pi(\cdot | h_1, h_2, h_3, h_4, h_5, h_6, a_0, a_1, a_2, a_3, a_4, a_5, T_0, b, p_0, \mathbf{y})$



**Figure B.21:** IACT and autocorrelation function for samples  $T_0 \sim \pi(\cdot | h_1, h_2, h_3, h_4, h_5, h_6, a_0, a_1, a_2, a_3, a_4, a_5, a_6, b, p_0, \mathbf{y})$



**Figure B.22:** IACT and autocorrelation function for samples  $b \sim \pi(\cdot|h_1, h_2, h_3, h_4, h_5, h_6, a_0, a_1, a_2, a_3, a_4, a_5, a_6, T_0, p_0, y)$



**Figure B.23:** IACT and autocorrelation function for samples  $p_0 \sim \pi(\cdot | h_1, h_2, h_4, h_5, h_6, a_0, a_1, a_2, a_3, a_4, a_5, a_6, T_0, b, \mathbf{y})$