

# State Space Models for Regional Epidemiological Indicators

Stefan Heyder

October 2023 – Draft v 0.1



# Abstract

A major challenge for scientific inquiry during an ongoing epidemic is the multitude of uncertainties one must consider. In this thesis, we detail the demands that epidemiological data impose and demonstrate that state space models (SSMs) offer a flexible class of statistical models capable of capturing these effects while delivering results that are straightforward to interpret. To facilitate inference and predictions in these models, we use importance sampling techniques, for which two popular choices are the Cross-Entropy method (CE-method) and Efficient Importance Sampling (EIS), which are based on Kullback-Leibler and least-squares loss, respectively. For these two methods we provide central limit theorems that shed light on the empirical observation that EIS provides better importance sampling approximations. Our theoretical results reveal that this is likely due to the lower asymptotic variance of EIS. To demonstrate the capabilities of such models, we fit models for three real-world applications using Germany's COVID-19 data. The first model illustrates how SSMs can be employed to account for the reporting process in detail, which may be used to handle reporting artifacts caused by factors such as holiday periods. Second, we provide a model that explicitly accounts for the exchange of cases between spatial regions. We use this model to perform one-week-ahead forecasts of reported cases for Germany and validate them against the ECDC's ForecastHub dataset, consisting of such forecasts made in real-time. Finally, we introduce a model for the delayed reporting of hospitalizations in Germany, which we use to nowcast the hospitalization incidence. We also compare the predictive performance of this model to real-time nowcasts provided by the German NowcastHub.

# Zusammenfassung

Kurze Zusammenfassung des Inhaltes in deutscher Sprache...



# Publications and Contributions

This thesis consists mostly of unpublished work. During my time as a PhD student, I have, however, been fortunate to collaborate with many scientists on problems in mathematical epidemiology with a focus on COVID-19, resulting in several preprints and publications. In this section I aim to clarify my contributions to these works and distinguish them from the contributions of the present thesis. Peer-reviewed publications are marked with a superscript<sup>¶</sup>. For the original contributions of this thesis, refer to the contribution statements located at the beginning of each chapter.

The first set of works follows the convention of sorting authors by their last name.

**Hotz et al., 2020** Thomas Hotz conceived the initial idea for this paper, derived most of the results, and wrote the initial manuscript. Matthias Glock and I managed the data cleaning, estimation of reproduction numbers and automation of the associated dashboard. Alexander Krämer, Anne Böhle, and Sebastian Semper provided consultation on epidemiological and practical relevance.

**Burgard et al., 2021<sup>¶</sup>** Thomas Hotz initially suggested applying small area estimation techniques to reproduction number estimates. I developed the model, estimator and simulation study, consulting with Thomas Hotz throughout the process. I wrote the first draft of the paper. Jan Pablo Burgard provided consultations on the small-area aspect of the paper, and Tyll Krüger on the epidemiological implications.

**Grundel et al., 2022<sup>¶</sup>** Sara Grundel and Karl Worthmann developed the idea of using optimal control techniques to balance testing with non-pharmaceutical interventions under societal constraints. Thomas Hotz and I designed the compartmental model, incorporating realistic parameters, and derived the epidemiological implications of the optimal strategies. Tobias Ritschel and Philipp Sauerteig established the mathematical and numerical results related to optimal control theory. The writing for the initial version of the paper was divided among Tobias Ritschel, Philipp Sauerteig and myself.

**Grundel et al., 2021<sup>¶</sup>** Sara Grundel and Karl Worthmann came up with the idea of applying the optimal control techniques from testing to vaccination strategies. Thomas Hotz and myself consulted on how to adapt the compartmental model to account for vaccination instead of testing and contributed to the epidemiological interpretation of the results. The initial version of the paper was written by Tobias Ritschel and Philipp Sauerteig.

**Heyder and Hotz, 2023** Thomas Hotz and I conceived with the idea of comparing different measures of COVID-19 spread with respect to ease of interpretation and communication. I then developed these ideas into an initial manuscript, except for the derivation of the renewal equation, which was contributed by Thomas Hotz.

The second set of works follows the convention of sorting authors by contribution.

**Bracher et al., 2021<sup>¶</sup>** and **Bracher et al., 2022<sup>¶</sup>** These two papers are a result of the joint efforts of the German and Polish ForecastHub<sup>1</sup>, organized by the Chair of Econometrics and Statistics at Karlsruhe Institute of Technology and the Computational Statistics Group at Heidelberg Institute for Theoretical Studies. Together with the authors of (Burgard et al., 2021), I contributed the ITWW\_country\_repro model, based on the same reference. This included automating the weekly submission of forecasts and actively participating in the weekly group discussions. Based on a pre-registered study protocol, these discussions, and extensive evaluations, the group from Karlsruhe, led by Johannes Bracher, wrote the initial manuscripts.

**Sherratt et al., 2022<sup>¶</sup>** This paper is based on the results of the European ForecastHub<sup>2</sup>, spearheaded by the European Centre for Disease Prevention and Control (ECDC). Again, Thomas Hotz and I contributed the ITWW model from the German and Polish ForecastHub, participated in weekly discussions and additionally contributed the ILM-EKF model, based on Thomas Hotz's initial idea.

**Brockhaus et al., 2023<sup>¶</sup>** This paper was conceived and written by Elisabeth Brockhaus and Johannes Bracher. My contribution includes the estimates of the Ilmenau model for reproduction numbers over time, and I participated in discussion and interpretation of results.

**Wolffram et al., 2023<sup>¶</sup>** The results of this work are based on the German NowcastHub<sup>3</sup>, which Thomas Hotz and I contributed daily nowcasts of the ILM-prop to. Again, the results of this paper are based on weekly group discussions, and the initial manuscript was prepared by Johannes Bracher and Daniel Wolffram.

## Own publications

- Bracher, J. et al. (Aug. 27, 2021). “A Pre-Registered Short-Term Forecasting Study of COVID-19 in Germany and Poland during the Second Wave.” In: *Nat Commun* 12.1 (1), p. 5173. ISSN: 2041-1723. DOI: [10.1038/s41467-021-25207-0](https://doi.org/10.1038/s41467-021-25207-0).
- Bracher, J. et al. (Oct. 31, 2022). “National and Subnational Short-Term Forecasting of COVID-19 in Germany and Poland during Early 2021.” In: *Commun Med* 2.1 (1), pp. 1–17. ISSN: 2730-664X. DOI: [10.1038/s43856-022-00191-8](https://doi.org/10.1038/s43856-022-00191-8).
- Brockhaus, E. K. et al. (Nov. 27, 2023). “Why Are Different Estimates of the Effective Reproductive Number so Different? A Case Study on COVID-19 in Germany.” In: *PLOS Computational Biology* 19.11, e1011653. ISSN: 1553-7358. DOI: [10.1371/journal.pcbi.1011653](https://doi.org/10.1371/journal.pcbi.1011653).
- Burgard, J. P. et al. (Aug. 31, 2021). “Regional Estimates of Reproduction Numbers with Application to COVID-19.” arXiv: [2108.13842 \[stat\]](https://arxiv.org/abs/2108.13842). URL: <http://arxiv.org/abs/2108.13842> (visited on 09/30/2021).
- Grundel, S. et al. (Jan. 1, 2021). “How to Coordinate Vaccination and Social Distancing to Mitigate SARS-CoV-2 Outbreaks.” In: *SIAM J. Appl. Dyn. Syst.* 20.2, pp. 1135–1157. DOI: [10.1137/20M1387687](https://doi.org/10.1137/20M1387687).
- (Apr. 1, 2022). “How Much Testing and Social Distancing Is Required to Control COVID-19? Some Insight Based on an Age-Differentiated Compartmental Model.” In: *SIAM J. Control Optim.* 60.2, S145–S169. ISSN: 0363-0129, 1095-7138. DOI: [10.1137/20M1377783](https://doi.org/10.1137/20M1377783).
- Heyder, S. et al. (Oct. 4, 2023). “Measures of COVID-19 Spread.” In: *Covid-19 pandisziplinär und international: Gesundheitswissenschaftliche, gesellschaftspolitische und philosophische Hintergründe*. Ed. by A. Kraemer et al. Medizin, Kultur, Gesellschaft. Wiesbaden: Springer Fachmedien, pp. 51–66. ISBN: 978-3-658-40525-0. DOI: [10.1007/978-3-658-40525-0\\_3](https://doi.org/10.1007/978-3-658-40525-0_3).

<sup>1</sup><https://github.com/KITmetricslab/covid19-forecast-hub-de>

<sup>2</sup><https://covid19forecasthub.eu/>

<sup>3</sup><https://covid19nowcasthub.de/>

- Hotz, T. et al. (Apr. 18, 2020). “Monitoring the Spread of COVID-19 by Estimating Reproduction Numbers over Time.” arXiv: [2004.08557](https://arxiv.org/abs/2004.08557) [q-bio, stat]. URL: <http://arxiv.org/abs/2004.08557> (visited on 07/20/2020).
- Sherratt, K. et al. (June 16, 2022). *Predictive Performance of Multi-Model Ensemble Forecasts of COVID-19 across European Nations*. DOI: [10.1101/2022.06.16.22276024](https://doi.org/10.1101/2022.06.16.22276024). Pre-published.
- Wolfram, D. et al. (Aug. 11, 2023). “Collaborative Nowcasting of COVID-19 Hospitalization Incidences in Germany.” In: *PLOS Computational Biology* 19.8, e1011394. ISSN: 1553-7358. DOI: [10.1371/journal.pcbi.1011394](https://doi.org/10.1371/journal.pcbi.1011394).

*Du musst bereit sein Dinge zu tun.*

— A meme on the internet, 2022.

# Acknowledgments

Put your acknowledgments here.



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Epidemiological considerations</b>	<b>3</b>
<b>3</b>	<b>Importance Sampling in State Space Models</b>	<b>5</b>
3.1	Gaussian Linear State Space Models . . . . .	9
3.2	Partially Gaussian state space models . . . . .	14
3.3	Importance Sampling . . . . .	18
3.3.1	Laplace approximation (LA) . . . . .	27
3.3.2	The Cross-Entropy method (CE-method) . . . . .	28
3.3.3	Efficient Importance Sampling (EIS) . . . . .	36
3.4	Interim discussion . . . . .	41
3.5	Gaussian importance sampling for state space models . . . . .	43
3.5.1	The Gaussian linear state space model (GLSSM)-approach . . . . .	43
3.5.2	The Markov-approach . . . . .	48
3.6	Inference in PGSSMs . . . . .	56
3.6.1	Maximum likelihood estimation . . . . .	56
3.6.2	Posterior inference . . . . .	62
3.7	Comparison of Importance Sampling method . . . . .	63
3.7.1	Breakdown of methods . . . . .	63
3.7.2	Computational complexity . . . . .	65
3.7.3	Asymptotic variance and relative efficiencies . . . . .	66
3.7.4	Performance of the optimal proposal . . . . .	72
3.8	Conclusion . . . . .	75
<b>4</b>	<b>Analysis of selected models</b>	<b>77</b>
4.1	Removing reporting delays and weekday effects . . . . .	78
4.1.1	Context . . . . .	78
4.1.2	Model . . . . .	78
4.1.3	Results . . . . .	81
4.1.4	Discussion . . . . .	84
4.2	Regional growth factor model . . . . .	86
4.2.1	Context . . . . .	86
4.2.2	Model . . . . .	87
4.2.3	Results . . . . .	88
4.2.4	Discussion . . . . .	88
4.3	Nowcasting hospitalizations . . . . .	88
4.3.1	Context . . . . .	88
4.3.2	Model . . . . .	90
4.3.3	Results . . . . .	93
4.3.4	Discussion . . . . .	96
<b>5</b>	<b>Discussion</b>	<b>97</b>

A Reproducibility and code	99
B Additional calculations	101
Bibliography	103
List of symbols	111
List of abbreviations	113

# Chapter 1

## Introduction

The Coronavirus disease 2019 (COVID-19) pandemic put the scientific community to the test: How infectious, morbid, and mortal was the disease? When and for how long did infected people become infectious? How effective are the countermeasures taken? How safe and effective are the vaccines that were developed at an unprecedented speed? Some of these questions, such as those about the epidemiology of COVID-19, are confined to well-established areas of research, while others, e.g. those about the efficacy of countermeasures, required collaboration across a wide range of disciplines — from infectious disease epidemiology, mathematical and statistical modeling, social and communication science, to non-scientific actors such as legislators, journalists and politicians.

Although there is still a significant amount of scientific and societal follow-up work to be done, given the magnitude of this challenge, it is astonishing how well the scientific community and society as a whole have dealt with the pandemic. A key factor in this accomplishment is the large-scale availability of data surrounding the pandemic. In many countries, including Germany, data on reported cases, deaths, vaccinations, and deaths were published daily by the respective national health authorities, i.e. the Robert Koch-Institut (RKI) (Robert Koch-Institut, 2021, 2022) in Germany. As the news reported daily on the number of newly reported cases and deaths, numerous dashboards with analyses of COVID-19 data were made available and an abundant number of scientific works were created, effectively communicating with the public, whose cooperation with countermeasures was critical, became increasingly important. To disseminate insights to the public, we need to understand and communicate the underlying dynamics of an epidemic.

An epidemic outbreak is inherently a random phenomenon (Diekmann, Heesterbeek, and Britton, 2013). Who becomes infected, how long they stay infectious, whom they meet while they are infectious and whom they ultimately infect are all aspects that depend to a certain degree on chance. If one is interested in large-scale phenomena, e.g. effects of immunization in a large population, one rely on a deterministic model (Britton et al., 2019), such as the classical S(E)IR model (Kermack and McKendrick, 1927) or its variants. However, as soon as one is interested in more detailed phenomena, as we are in this thesis, stochastic and statistical modeling becomes essential.

As statisticians, having access to vast amounts of data is both a blessing and a curse. While more, and ideally better, data enables us to formulate and address more relevant questions, the models we create to accommodate these data become increasingly complex. A major complexity of epidemic models are the dependencies across time, requiring the use of methods from time-series analysis. However, such models require more care in modeling, fitting, and interpretation, as there are more opportunities for error along the way. As we incorporate more detailed effects into our models, fitting the models to data becomes difficult to practically impossible using established techniques. While there are some remedies for this curse of dimensionality, e.g. exploiting as much available structure as possible, there is an ongoing need for new procedures enabling inference in these settings. As of the writing of this thesis, there is no comprehensive framework for creating such models that details how to incorporate the various phenomena practitioners are interested in. Additionally, we need both mathematical and practical insight into the performance of these procedures to make

informed decisions in applied settings: which methods should we prefer under which circumstances?

These considerations set the stage for this thesis. Driven by the need for good statistical models that allow us to answer urgent questions in infectious disease epidemiology, with COVID-19 as a prime example, we will start with an analysis of what is required of these sought-after models. We will define and discuss the role of several epidemiological indicators, namely quantities that have an interpretation related to the epidemic. It turns out that we will usually be interested in quantifying the speed at which the epidemic proliferates, and we discuss several popular indicators that measure this speed. A useful statistical analysis should provide interpretable insight into the problem at hand, so we focus on how straightforward this interpretation is, offering recommendations on when to use which indicator. To estimate these indicators from data, we must create statistical models that include them. Before we do so, we will compile a list of desiderata from the context of COVID-19 to identify a suitable class of such models.

Once we have a clear view of the epidemiological problems at hand, we demonstrate that many of the desiderata can be addressed by using SSMs, a flexible framework for modeling non-stationary time series. Unfortunately, we will require that these SSMs include integer-valued, non-Gaussian, observations, which makes fitting the models to data analytically impossible and numerically difficult, as one is faced with an intractable high-dimensional distribution; similar to challenges arising in Bayesian inference. Instead of analytical derivations, inference will have to rely on simulation methods, most notably importance sampling. To apply these methods, the practitioner has some flexibility in the so-called proposal distribution, a tractable approximation to sought-after distribution. Different disciplines have developed simulation-based techniques that allow the user to choose optimal proposals, where optimality is based on different performance criteria for different methods. In this thesis, we focus on two methods: the Cross-Entropy method (CE-method) and Efficient Importance Sampling (EIS), corresponding to Kullback Leibler divergence and least-squares loss respectively. In the literature, a comparison between these two methods is missing: there are neither mathematical nor empirical results comparing the two. We fill this gap by first proving central limit theorems for both methods, allowing for a theoretical comparison. Additionally, we provide extensive simulation studies comparing the methods on instructive univariate and SSMs examples. To this end, we also develop a new algorithm that allows the CE-method to be applied to state space models (SSMs).

Finally, we demonstrate how to solve a selection of infectious disease epidemiology problems using the mathematical insights we have gained. These examples focus on the COVID-19 epidemic in Germany and illustrate the modeling, computational, and applied aspects of this thesis. We focus on the following three challenges by providing models and insights for German data: How does one model and account for the complex and artifact-prone reporting process of incident cases? How does one incorporate cross-regional infections into these models? And, finally, how does one account for the long reporting delays present in hospitalizations?

In summary, this thesis contributes to both theoretical and applied problems in the context of statistical modeling of infectious diseases. We hope that its results will enhance our ability to respond to future epidemic outbreaks.

## Chapter 2

# Epidemiological considerations

Contributions of this chapter

??  
??  
??

The spread of infectious diseases, such as COVID-19, is a complex phenomenon. For COVID-19, this complexity arises from the interplay of many factors. Studying these factors will allow us to define the aims and challenges of epidemiological modeling in the context of this thesis. It will also guide us towards desirable and achievable outcomes of our efforts from an applied perspective.

First of all, there is considerable heterogeneity in the progression of the disease once an individual is infected (Salzberger et al., 2021). Some infectees may show few to no symptoms but are still highly infectious (Byambasuren et al., 2020), and disease progression is closely linked to age and preexisting comorbidities (Biswas et al., 2020). Additionally, different variants of SARS-CoV-2 differ in key epidemiological characteristics such as the reproduction number (Du et al., 2022) and mortality (Hughes et al., 2023).

Second, the spread is highly dependent on the contact behavior within the population, as the infector must be in close physical proximity to the infectee to infect them. These contact patterns are an essential component of any mathematical model for infectious diseases, as they define how the epidemic evolves. While there are some empirical studies (Mossong et al., 2008; Tomori et al., 2021), capturing the contact behavior at certain points in time, in the context of an ongoing epidemic, these patterns are subject to change, not only in intensity but also in shape (Tomori et al., 2021).

Finally, as the disease spreads in the population and vaccinations become available, the population develops partial immunity against the disease, if not against infection (Wu et al., 2023). This immunity affects the spread as well: if an infector has contact with a partially immune individual, the probability of transmission is reduced. Additionally, partial immunity may lead infectors to develop fewer or no symptoms so they may not be aware of being infectious, foregoing quarantine.

As statisticians, we face a challenging problem: Which of these factors should we include in our model and how should we incorporate them? The answer certainly depends on the specific epidemiological question under consideration and the availability and quality of the data.

Parts of this chapter, especially ??????, consist of the ideas published in (Heyder and Hotz, 2023), but have been rewritten to fit better into this thesis.

## Chapter 3

# Importance Sampling in State Space Models

### Contributions

The main contribution of this chapter consists of a rigorous comparison of two importance sampling frameworks: the Cross-Entropy method (CE-method) and Efficient Importance Sampling (EIS). Both methods determine optimal importance sampling proposals, but have, until now, been studied in separate communities: the CE-method is popular in rare-event estimation and engineering disciplines, while EIS is popular in the financial time series community.

The contributions of the individual sections are as follows:

??

**Gaussian Linear State Space Models** This section is a condensed introduction to Gaussian linear state space models (GLSSMs) and is loosely based on (Durbin and Koopman, 2012).

**Partially Gaussian state space models**

**Importance Sampling**

- We prove Lemma 3.5.
- Discussion surrounding (Chatterjee and Diaconis, 2018).
- We prove central limit theorems for both methods (Sections 3.3.2 and 3.3.3).
- Proof Proposition 3.5.

**Interim discussion**

**Gaussian importance sampling for state space models** derive an efficient algorithm to apply the CE-method to state space models (Section 3.5.2)

**Inference in PGSSMs**

**Comparison of Importance Sampling method** Extensively compare both methods on theoretical as well as practically relevant properties with instructive univariate and multivariate examples (Section 3.7).

In the last chapter, we have detailed the need for models capable of modeling complex epidemiological phenomena. Crucially, we must account for the time-series and discrete nature of the data at hand. State space models (SSMs) form a versatile class of statistical models that allow modeling of non-stationary time series data while providing a straightforward, mechanistic interpretation of the time series' dynamics. The main idea of these models is to introduce unobserved **latent states** whose joint distribution is governed by a Markov process, and to model the observed time series conditional on these states. By exploiting this structure, inference in SSMs becomes computationally efficient, as the complexity of algorithms is linear in the number  $n$  of time points considered. In this chapter, we provide a mathematical introduction to the theory of SSMs and the main tool we will use for inference: importance sampling. The foundations of SSMs presented in this chapter are, if not mentioned otherwise, based on (Chopin and Papaspiliopoulos, 2020; Durbin and Koopman, 2012).

Let us start from a very general definition of a SSM.

**Definition 3.1** (State Space Model). A **state space model** is a discrete time stochastic process  $(X_t, Y_t)_{t=0, \dots, n}$  taking values in the measurable space  $(\mathcal{X} \times \mathcal{Y}, \mathcal{B}_{\mathcal{X}} \otimes \mathcal{B}_{\mathcal{Y}})$  such that

- (i) The marginal distribution of the **states**  $(X_0, \dots, X_n)$  is a discrete time Markov process, i.e. for  $t = 1, \dots, n$

$$\mathbf{P}(X_t \in B | X_0, \dots, X_{t-1}) = \mathbf{P}(X_t \in B | X_{t-1}) \text{ a.s.} \quad (3.1)$$

for all measurable  $B \in \mathcal{B}_{\mathcal{X}}$ .

- (ii) Conditional on the state  $X_t$  and observation  $Y_{t-1}$ ,  $Y_t$  is independent of  $X_s$  and  $Y_{s-1}$ ,  $s < t$ , i.e.

$$\mathbf{P}(Y_t \in B | X_0, \dots, X_t, Y_0, \dots, Y_{t-1}) = \mathbf{P}(Y_t \in B | X_t, Y_{t-1})$$

for all measurable  $B \in \mathcal{B}_{\mathcal{Y}}$ .

For notational convenience, we will write  $X_{s:t} = (X_s, \dots, X_t)$  for the vector that contains all states from  $s$  to  $t$ ,  $s \leq t$ , dropping the first index if we consider the whole set of observations up to time  $t$ , so  $X_{:t} = X_{0:t}$ , and dropping the subscript if we consider all states at once,  $X = X_{:n}$ . Similarly we set  $Y_{s:t} = (Y_s, \dots, Y_t)$ ,  $Y_{:t} = Y_{0:t}$  and  $Y = Y_{:n}$ .

The models that we consider in this thesis will usually assume that densities for the state transitions with respect to a common dominating measure  $\mu_{\mathcal{X}}$  and similar for the observations with respect to some dominating measure  $\mu_{\mathcal{Y}}$ .

**Notation 3.1** (Densities, conditional densities). We will use the standard abuse of notation for densities that makes the type of density „obvious“ from the arguments used. This means that  $p(x)$  is the density for all states  $X$ , evaluated at  $x$ ,  $p(x_t | x_{t-1})$  the conditional density of  $X_t | X_{t-1}$ , evaluated at  $x_t$  and  $x_{t-1}$  and similarly for observations:  $p(y|x)$  is the density of the conditional distribution of all observations  $Y$  conditioned on all states  $X$ , evaluated at  $y$  and  $x$ .

Note that this notation also implicitly includes the time  $t$  and allows for changes in, e.g., the state transition over time.

When densities come from a parametric model parametrized by  $\theta \in \Theta \subseteq \mathbf{R}^l$  and the dependence of the model on  $\theta$  is of interest, i.e. because we try to estimate  $\theta$ , we indicate this by adding a subscript to the densities. If this dependence is not of interest, e.g. because  $\theta$  is fixed, we omit  $\theta$  for better readability.

In this notation, the joint density of a parametric SSM factorizes as

$$\begin{aligned} p_{\theta}(x, y) &= p_{\theta}(x_0, \dots, x_n, y_0, \dots, y_n) \\ &= p_{\theta}(x_0) \prod_{t=1}^n p_{\theta}(x_t | x_{t-1}) \prod_{t=0}^n p_{\theta}(y_t | x_t, y_{t-1}), \end{aligned} \quad (3.2)$$

where  $p_{\theta}(y_0 | x_0, y_{-1}) = p_{\theta}(y_0 | x_0)$ .



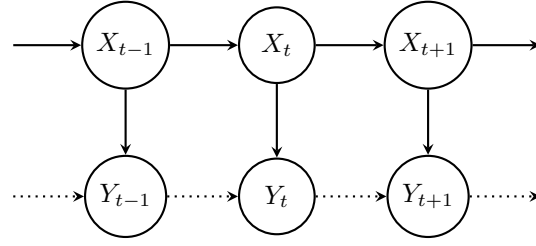


Figure 3.1: Dependency structure in a SSM as given by Definition 3.1. The dependencies between observations  $Y_{t-1}$  (indicated by dotted arrows) are usually not part of the standard definition of a SSM, but can be incorporated in a straightforward manner.

As inferences made in this thesis depend on the SSM only through the likelihood we identify almost sure versions of  $(X, Y)$  with themselves, i.e. all equations involving  $X$  or  $Y$  are understood almost surely.

**Remark 3.1** (dependence on  $Y_{t-1}$ , dimensions). Contrary to the standard definition of a SSM, as found in, e.g., (Chopin and Papaspiliopoulos, 2020, Chapter 2) or (Durbin and Koopman, 2012, Chapter 9), our Definition 3.1 allows  $Y_t$  to depend on  $Y_{t-1}$ . As the models considered in Chapter 4 will make extensive use of SSMs with this dependency structure we opt to use this non-standard definition here. This is of course not a limitation of the standard definition: given a SSM of the form described in Definition 3.1, we can transform it to the standard form by choosing states  $(X_t, Y_t) \in \mathcal{X} \times \mathcal{Y}$  and observations  $Y_t \in \mathcal{Y}$  such that the SSM becomes a stochastic process on  $(\mathcal{X} \times \mathcal{Y}) \times \mathcal{Y}$ .

Additionally, the goal of our inferences will always be functionals of the conditional distribution  $X|Y$  for a single, fixed, set of observations  $y$ . Assuming all densities exist, the conditional density  $p(x|y)$  is given, up to a constant not depending on  $x$ , by Equation (3.2):

$$p(x|y) \propto p(x, y) = p(x_0) \prod_{t=1}^n p(x_t|x_{t-1}) \prod_{t=0}^n p(y_t|x_t, y_{t-1}).$$

Thus, the dependence of  $Y_t$  on  $Y_{t-1}$  only affects our inferences through  $p(y_t|x_t, y_{t-1})$ , where, as  $Y_{t-1}$  is observed, the argument  $y_{t-1}$  is fixed. Consequently, all results we present in this chapter for SSMs where  $Y_t$  depends only on  $X_t$  that concern only the conditional distribution  $X|Y = y$  carry over to those given by Definition 3.1. We will reiterate this argument at appropriate points in this thesis.

In most SSMs we consider in this thesis we use  $\mathcal{X} = \mathbf{R}^m$  and  $\mathcal{Y} = \mathbf{R}^p$  or  $\mathcal{Y} = \mathbf{Z}^p$  so that  $\mathcal{X}$  is  $m$  dimensional and  $\mathcal{Y}$  is  $p$  dimensional and equip these spaces with the usual  $\sigma$ -Algebras. Unless noted otherwise, we use for  $\mu_{\mathcal{X}}$  the  $m$ -dimensional Lebesgue measure and for  $\mu_{\mathcal{Y}}$  either the  $p$ -dimensional Lebesgue measure ( $\mathcal{Y} = \mathbf{R}^p$ ) or the  $p$ -dimensional counting measure ( $\mathcal{Y} = \mathbf{Z}^p$ ).

Given data  $y = (y_t)_{t=0, \dots, n}$  that may be modeled with a SSM the practitioner is confronted with several tasks, which provide the structure of this chapter:

- (i) Choosing a suitable, usually parametric, class of SSMs that include the effects of interest.
- (ii) Fitting such a parametric model to the data at hand by either frequentist or Bayesian techniques.
- (iii) Infer the latent states  $X$  from the observations by determining, either analytically or through simulation, the smoothing distribution  $X|Y$ .

The first step, Item (i), requires that the practitioner specifies a joint probability distribution for the states and observations (see Chapter 4 for examples of this). Due to the assumed dependency structure, this boils down to specifying transition kernels for the states and observations. The setting given in Definition 3.1 is too abstract to perform inference in, so further assumptions on the types of distributions for the latent states and observations are needed. In this chapter, we

will discuss Gaussian linear state space model (GLSSM) (Section 3.1), where both the posterior distribution and the likelihood can be derived analytically. For the epidemiological application we have in mind, these are, however, insufficient due to the non-linear behavior of incidences and the low count per region (??). Such observations are better modeled with distributions on the natural numbers, i.e. with a Poisson or negative binomial distribution, both of which are exponential families of distributions. This will lead to the class of Partially Gaussian state space models (PGSSMs) (Section 3.2) which will become the main focus of our study.

Regarding the second step, Item (ii), a frequentist practitioner will want to perform maximum likelihood inference on  $\theta$ . While asymptotic confidence intervals for the maximum likelihood estimator (MLE)  $\hat{\theta}$  can be derived both theoretically and practically (Durbin and Koopman, 2012, Chapter 7), they are, in the context of this thesis, usually of little interest. For these asymptotic frequentist procedures to be meaningful, an appropriate central limit theorem must hold. However, as the time series we study are non-stationary and the dependence on parameters  $\theta$  is allowed to be arbitrary, it is in general not obvious that such a theorem holds for the model under consideration. Instead, we approach this fitting as an Empirical Bayes procedure and our main practical interest lies in analyzing the posterior distribution  $X|Y$  where we set  $\theta$  equal to  $\hat{\theta}$ .

To obtain the maximum likelihood estimates  $\hat{\theta}$  one needs access to the likelihood

$$p_{\theta}(y) = \int_{\mathcal{X}^n} p_{\theta}(x, y) dx = \int_{\mathcal{X}^n} p_{\theta}(y|x)p_{\theta}(x) dx \quad (3.3)$$

which is usually not analytically available. Direct numerical evaluation of Equation (3.3) is hopeless due to the high dimensionality of the state space  $\mathcal{X}^n$ . Instead, we will resort to simulation-based inference by importance sampling (see Section 3.3), a Monte-Carlo method that approximates  $p(y)$  by constructing a global tractable approximation to the integrand in Equation (3.3). Alternatively, sequential Monte Carlo (SMC) methods, i.e. particle filters, that perform importance sampling sequentially across the  $n + 1$  time steps can be used. We will not follow this approach for reasons described later, but refer the reader to the excellent reference (Chopin and Papaspiliopoulos, 2020) for an introduction to these methods.

check we do this below

The performance of these simulations depends crucially on our ability to construct distributions that are close to the posterior  $p(x|y)$  but are easy to sample from. To this end, we construct either Gaussian linear state space models (GLSSMs) (Section 3.5.1) in which sampling from the posterior is analytically possible, or Gaussian Markov processes (Section 3.5.2) which are directly amenable to simulation. These two approaches are motivated by what we term „optimal importance sampling“, where we use a proposal distribution that solves an optimization problem. Two popular approaches for choosing such a proposal are Efficient Importance Sampling and the Cross-Entropy method, which minimize an  $L^2$  or KL-divergence loss, respectively. Empirically, it has been shown that EIS outperforms the CE-method, to which we add theoretical insight in the form of two central limit theorems (Section 3.3): As both methods rely on importance sampling to determine an optimal proposal, the asymptotic variance of this procedure is of practical relevance, and we argue that this asymptotic variance is smaller for EIS. To this end, we provide extensive simulation studies investigating the asymptotic variance of the two methods in Section 3.7. To the best of the authors' knowledge, this is the first rigorous investigation comparing these two methods.

As an alternative to the MLE approach, a fully Bayesian approach would regard  $\theta$  as random and administer a prior distribution, say with density  $p(\theta)$ . In this setting, the main interest still lies in determining the posterior distribution of  $X|Y = y$ , but due to the prior put on  $\theta$ , its density, should it exist, is now given by

$$p(x|y) = \int p(x, \theta|y) d\theta,$$

where  $p(x, \theta|y)$  is the joint posterior of states and hyperparameters, conditional on observations  $y$ . To tackle this problem, one may again use importance sampling methods, see e.g. (Durbin and Koopman, 2012, Chapter 13.1), or use MCMC-methods tailored to SSMs, e.g. Particle-MCMC (Chopin and Papaspiliopoulos, 2020, Chapter 16).

To perform these tasks, the setting defined in Definition 3.1 is too general to have numerically tractable solutions, as such we restrict our studies to more structured SSMs. We begin with assuming a joint Gaussian distribution for states and observations. Subsequently, we will relax this assumption to allow the observations  $Y$  to have more general distributions.

### 3.1 Gaussian Linear State Space Models

Gaussian linear state space models (GLSSMs) are the working horses of most methods used in this thesis because many of the interesting quantities, e.g. the smoothing distribution, are analytically tractable and can be obtained computationally efficient. Indeed, for fixed dimension of states  $m$  and observations  $p$  the runtime of algorithms that we consider in this thesis is  $\mathcal{O}(n)$ , i.e. linear in the number of time points observed.

**Definition 3.2** (GLSSM). A Gaussian linear state space model (GLSSM) is a joint distribution over states and observations  $(X, Y)$  where states a.s. obey the transition equation

$$X_{t+1} = A_t X_t + u_t + \varepsilon_{t+1} \quad t = 0, \dots, n-1, \quad (3.4)$$

and observations a.s. obey the observation equation

$$Y_t = B_t X_t + v_t + \eta_t \quad t = 0, \dots, n. \quad (3.5)$$

Here  $A_t \in \mathbf{R}^{m \times m}$  and  $B_t \in \mathbf{R}^{p \times m}$  are matrices that specify the systems dynamics. The **innovations**  $(\varepsilon_{t+1})_{t=0, \dots, n-1}$  and **measurement noise**  $(\eta_t)_{t=0, \dots, n}$  and the starting value  $X_0 \sim \mathcal{N}(\mathbb{E}X_0, \Sigma_0)$  are jointly independent. Furthermore,  $\varepsilon_{t+1} \sim \mathcal{N}(0, \Sigma_t)$  and  $\eta_t \sim \mathcal{N}(0, \Omega_t)$  are centered Gaussian random variables and  $u_t \in \mathbf{R}^m, t = 0, \dots, n-1, v_t \in \mathbf{R}^p, t = 0, \dots, n$  are deterministic biases.

**Remark 3.2.** From Equation (3.4) it is easy to see that the states  $X = (X_0, \dots, X_n)$  form a Gaussian Markov process and that conditional on  $X_t, t \in \{0, \dots, n\}$ ,  $Y_t$  is independent of  $X_s$  and  $Y_s, s < t$ . Thus a GLSSM is indeed a SSM.

The defining feature of a GLSSM is that the joint distribution of  $(X, Y)$  is Gaussian, as  $(X, Y)$  may be written as an affine combination of the jointly Gaussian  $(X_0, \varepsilon_1, \dots, \varepsilon_n, \eta_0, \dots, \eta_n)$  and it is often useful to perform inferences in terms of innovations and measurement noise instead of states, see e.g. (Durbin and Koopman, 2012, Section 4.5).

As the joint distribution of  $(X, Y)$  is Gaussian, so are conditional distributions of states given any set of observations.

**Lemma 3.1** (Gaussian conditional distributions). *Let  $(X, Y)$  be jointly Gaussian with distribution  $\mathcal{N}(\mu, \Sigma)$  where*

$$\mu = (\mu_X, \mu_Y)$$

and

$$\Sigma = \begin{pmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{pmatrix},$$

where  $\mu$  and  $\Sigma$  are partitioned according to the dimensions of  $X$  and  $Y$ .

Then the following holds:

- (i) If  $\Sigma_{YY}$  is non-singular,  $X|Y = y$  follows a Gaussian distribution with conditional expectation

$$\mu_{X|Y=y} = \mathbb{E}(X|Y = y) = \mu_X + \Sigma_{XY}\Sigma_{YY}^{-1}(y - \mu_Y)$$

and conditional covariance matrix

$$\Sigma_{X|Y=y} = \text{Cov}(X|Y = y) = \Sigma_{XX} - \Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{YX}.$$

- (ii) In particular, let  $X \sim \mathcal{N}(\mu, \Sigma)$  and  $Y = BX + \varepsilon$  for a matrix  $B \in \mathbf{R}^{p \times m}$  and  $\mathbf{R}^p \ni \varepsilon \sim \mathcal{N}(0, \Omega)$  independent of  $X$  where  $\Omega \in \mathbf{R}^{p \times p}$ . Then, as  $\mathbb{E}Y = B\mu$ ,  $\text{Cov}(X, Y) = \text{Cov}(Y, X)^T = \Sigma B^T$  and  $\text{Cov}(Y) = B\Sigma B^T + \Omega$ , we have

$$\mathbb{E}(X|Y = y) = \mu + K(y - B\mu)$$

and

$$\text{Cov}(X|Y = y) = \Sigma - K\Sigma K^T = (I - KB)\Sigma,$$

as long as  $B\Sigma B^T + \Omega$  is non-singular. Here  $K = \Sigma B^T (B\Sigma B^T + \Omega)^{-1}$ .

(iii) If  $\Sigma_{XX}$  is non-singular, then  $Y - BX$  is independent of  $X$  for  $B = \Sigma_{YX}\Sigma_{XX}^{-1}$  and we may write

$$Y = BX + v + \eta$$

for an  $\eta \sim \mathcal{N}(0, \Omega)$  with covariance matrix  $\Omega = \Sigma_{YY} - \Sigma_{YX}\Sigma_{XX}^{-1}\Sigma_{XY}$  independent of  $X$ , and deterministic  $v = \mu_Y - B\mu_X$ .

(iv) Suppose that  $(X, Y, Z)$  is jointly Gaussian with mean  $\mu$  and covariance matrix  $\Sigma$ , partitioned similarly as before. If the conditional distribution of  $X$  given  $Y = y$  and  $Z = z$  is given by

$$X|Y = y, Z = z \sim \mathcal{N}(Ky + Gz + v, \Xi),$$

then the conditional distribution of  $X$  given only  $Y = y$  is

$$X|Y = y \sim \mathcal{N}(Ky + G\mu_{Z|Y=y} + v, \Xi + G\text{Cov}(Z|Y)G^T).$$

**Remark 3.3** (generalized inverse). If  $\Sigma_{YY}$  in Lemma 3.1 (i) is singular, the statement remains true if we choose as  $\Sigma_{YY}^{-1}$  a generalized inverse of  $\Sigma_{YY}$ , see (Rao, 2002, 8.a Note 3). A generalized inverse for a matrix  $A \in \mathbf{R}^{m \times p}$  is any matrix  $A^- \in \mathbf{R}^{m \times p}$  such that  $AA^-A = A$ . Given a singular value decomposition  $A = UDV^T$ , we may obtain the Moore-Penrose inverse  $A^\dagger = VD^-U^T$  of  $A$ , which is a generalized inverse of  $A$ , by inverting the non-zero diagonal elements of  $D$ .

*Proof.* For the first statement, we refer the reader to (Durbin and Koopman, 2012, Chapter 4, Lemma 1).

The second statement follows from substituting the value of  $K$ .

The third statement follows from noting that  $Y - BX = \begin{pmatrix} -B & I \end{pmatrix} \begin{pmatrix} X \\ Y \end{pmatrix}$  follows a Gaussian distribution. A quick calculation reveals that

$$\text{Cov}(Y - BX, X) = \Sigma_{YX} - B\Sigma_{XX} = \Sigma_{YX} - \Sigma_{YX} = 0,$$

showing the independence. Thus  $\eta = Y - BX - v$  follows a centered Gaussian distribution and equating covariance matrices, we see that  $\Omega$  has the desired form.

For the final statement, notice that  $\xi = X - KY - GZ - v$  fulfills

$$\xi|Y = y, Z = z \sim \mathcal{N}(0, \Xi)$$

which does not depend on  $y$  or  $z$ . Thus the unconditional distribution of  $\xi$  is  $\mathcal{N}(0, \Xi)$  as well, and  $\xi$  is independent of  $(Y, Z)$ . Rewriting  $X$  in terms of  $Y, Z$  and  $\xi$ , we obtain

$$X = KY + GZ + v + \xi,$$

and so

$$\mathbb{E}(X|Y = y) = Ky + G\mathbb{E}(Z|Y = y) + v,$$

as well as

$$\begin{aligned} \text{Cov}(X|Y = y) &= \text{Cov}(KY + GZ + v + \xi|Y = y) \\ &= \text{Cov}(GZ + \xi|Y = y) \\ &= \text{Cov}(GZ + \xi) - \text{Cov}(GZ + \xi, Y)\Sigma_{YY}^{-1}\text{Cov}(Y, GZ + \xi) \\ &= G\Sigma_{ZZ}G^T + \Xi - G\Sigma_{ZY}\Sigma_{YY}^{-1}\Sigma_{YZ}G^T \\ &= \Xi + G\text{Cov}(Z|Y)G^T. \end{aligned}$$

□

After having observed  $Y = y$ , our main interest lies in the conditional distribution of states  $X$  given  $Y = y$ , which we could obtain by applying Lemma 3.1, i.e. where  $B = \text{block-diag}(B_0, \dots, B_n)$  and  $\Omega = \text{block-diag}(\Omega_0, \dots, \Omega_n)$  are block-diagonal matrices. However, this would require inversion of the  $(n+1)p \times (n+1)p$  matrix  $(B\Sigma B + \Omega)$  which becomes numerical infeasible quickly. Instead, we can exploit the sequential structure of the GLSSM, which will allow us to perform conditioning on only a single observation at a time.

To this end, let us denote by  $\hat{X}_{t|s} = \mathbb{E}(X_t|Y_{:s} = y_{:s})$  the conditional expectation of  $X_t$  given a set of observations  $y_{:s}$  and by  $\Xi_{t|s} = \text{Cov}(X_t|Y_{:s} = y_{:s})$  the conditional covariance matrix of  $X_t$  given  $Y_{:s} = y_{:s}$ . Then

$$X_t|Y_{:s} = y_{:s} \sim \mathcal{N}(\hat{X}_{t|s}, \Xi_{t|s}).$$

For a given  $t$ , three values of  $s$  are of particular interest: If  $s = t - 1$  determining this conditional distribution is called a **prediction problem**, if  $s = t$  this is a **filtering problem** and if  $s = n$  a **smoothing problem**, and we call the distributions we seek the **predictive**, **filtering** or **smoothing distribution** respectively. Similarly we define  $\hat{Y}_{t|s} = \mathbb{E}(Y_t|Y_{:s} = y_{:s})$  to be the conditional expectation of  $Y_t$  given  $Y_{:s} = y_{:s}$ , note that  $\hat{Y}_{t|s} = Y_t$  if  $s \geq t$ . Finally, let  $\Psi_{t|s} = \text{Cov}(Y_t|Y_{:s} = y_{:s})$  be the conditional covariance matrix of  $Y_t$  given  $Y_{:s} = y_{:s}$ . Again  $\Psi_{t|s} = 0$  if  $s \geq t$ .

These distributions may be obtained efficiently using the celebrated Kalman filter (Algorithm 1) and smoother (Algorithm 2) algorithms, which we state here for completeness.

---

**Algorithm 1** Kalman filter, with runtime  $\mathcal{O}(n(m^2 + p^3))$

---

**Require:** GLSSM (Definition 3.2), observations  $y_0, \dots, y_n$ .

```

1:  $A_{-1} \leftarrow I \in \mathbf{R}^{m \times m}$  ▷ Identity Matrix
2:  $u_{-1} \leftarrow \mathbf{0} \in \mathbf{R}^m$ 
3:  $\hat{X}_{-1|-1} \leftarrow \mathbb{E}X_0$ 
4:  $\Xi_{0|-1} \leftarrow \mathbf{0}_{m \times m}$ 
5:  $\ell_{-1} \leftarrow 0$ 
6: for  $t \leftarrow 0, \dots, n$  do
7:    $\hat{X}_{t|t-1} \leftarrow A_{t-1}\hat{X}_{t-1|t-1} + u_{t-1}$  ▷ prediction
8:    $\Xi_{t|t-1} \leftarrow A_{t-1}\Xi_{t-1|t-1}A_{t-1}^T + \Sigma_t$ 
9:    $\hat{Y}_{t|t-1} \leftarrow B_t\hat{X}_{t|t-1} + v_t$ 
10:   $\Psi_{t|t-1} \leftarrow B_t\Xi_{t|t-1}B_t^T + \Omega_t$ 
11:   $K_t \leftarrow \Xi_{t|t-1}B_t^T\Psi_{t|t-1}^{-1}$  ▷ filtering
12:   $\hat{X}_{t|t} \leftarrow \hat{X}_{t|t-1} + K_t(y_t - \hat{Y}_{t|t-1})$ 
13:   $\Xi_{t|t} \leftarrow \Xi_{t|t-1} - K_t\Psi_{t|t-1}K_t^T$ 
14:   $\ell_t \leftarrow \ell_{t-1} + \frac{p}{2} \log(2\pi) + \frac{1}{2} \log \det \Psi_{t|t-1} + \frac{1}{2} (y_t - \hat{Y}_{t|t-1})^T \Psi_{t|t-1}^{-1} (y_t - \hat{Y}_{t|t-1})$  ▷ NLL
15: end for
```

---

In Algorithm 1 every time point  $t = 0, \dots, n$  is processed in the same way, with a two-step procedure: first we predict the new observation  $Y_t$  based on  $Y_{:t-1}$ . Using the linearity of the system as well as the assumed conditional independence, this is achieved by applying the system dynamics to the current conditional expectation and covariance matrices. After  $Y_t$  has been observed, we can update the conditional distribution of the states by appealing to Lemma 3.1. For a rigorous derivation of the Kalman filter, we refer the reader to (Durbin and Koopman, 2012, Chapter 4) or the excellent monograph of (Schneider, 1986).

The Kalman filter is very efficient: each loop iteration requires inversion of the  $p \times p$  matrix  $\Psi_{t|t-1}$ . Assuming this operation dominates the time complexity, e.g. because  $m \approx p$ , the time complexity of the Kalman filter is  $\mathcal{O}(np^3)$ , a drastic improvement over the naïve  $\mathcal{O}(n^3m^3)$ , obtained by applying Lemma 3.1 to the joint distribution of  $(X, Y)$ . Similarly, the space complexity of Algorithm 1 is  $\mathcal{O}(n(m^2 + p^2))$ , and grows only linearly in  $n$ .

Notice that the Kalman filter iteratively calculates the negative log-likelihood  $\ell_t$

$$\ell_t = -\log p(y_{:t}) = -\log \sum_{s=0}^t \log p(y_s | y_{:(s-1)})$$

while filtering. This is possible because of the dependency structure of the GLSSM, which makes the increments in  $\ell_t$  tractable, as

$$Y_s | Y_{:(s-1)} \sim \mathcal{N}(\hat{Y}_{s|s-1}, \Psi_{s|s-1}),$$

for  $s = 0, \dots, n$ , which is shown in the derivation of the Kalman filter. Thus, the Kalman filter enables us to perform MLE by giving us access to  $\ell_n$ .

From this discussion we can see how we may alter the Kalman filter to accommodate a similar dependency structure as proposed in Definition 3.1 (depicted in Figure 3.1): If we allow to have

$$Y_t = B_t X_t + C_{t-1} Y_{t-1} + \eta_t \quad t = 0, \dots, n$$

we would still be able to perform the filtering step of Algorithm 1 by determining the conditional distribution of  $Y_t$  given  $Y_{:(t-1)}$  using Lemma 3.1 — .

Depending on the situation at hand, one of the many variants of the basic algorithm presented in Algorithm 1 may be used. If the inversion of  $\Psi_{t|t-1}$  is numerically unstable, the filtered covariance matrices  $\Xi_{t|t}$  may become numerically non-positive definite. In this case, the square root filter and smoother (Morf and Kailath, 1975) may be used. It is based on Cholesky roots of the involved covariance matrices, ensuring them to be PSD.

When the dimension of observations is much larger than that of the states,  $p \gg m$ , the information filter (Fraser and Potter, 1969) can be used. Instead of performing operations on the covariance matrices, i.e.  $\Xi_{t|t-1}$  and  $\Psi_{t|t-1}$ , the information filter operates on their inverses, the precision matrices  $\Xi_{t|t-1}^{-1}$  and  $\Psi_{t|t-1}^{-1}$  as well as rescaled states  $\Xi_{t|t-1}^{-1} \hat{X}_{t|t-1}$  and observation  $\Psi_{t|t-1}^{-1} \hat{Y}_{t|t-1}$  estimates. This makes the filtering step more efficient, as the most intensive step is the calculation of  $\Psi_{t|t-1}^{-1}$ . However, the price one pays is that the prediction step now requires inversion of a  $m \times m$  matrix, and as such the computational gains only manifest when  $p$  is sufficiently large compared to  $m$  (Assimakis, Adam, and Douladiris, 2012).

If the dimensions of the model are so large that calculating the  $m \times m$  and  $p \times p$  covariance matrices becomes an issue, the simulation based Ensemble Kalman filter (EnKF) (Evensen, 1994) can be used. Instead of calculating the covariance matrices analytically, the EnKF stores a particle approximation to the Gaussian filtering distribution and iteratively performs a prediction and update step with a particle approximation, similar to the analytical update the Kalman filter performs. Despite being based on linear Gaussian dynamics, the EnKF is successfully employed in many high-dimensional non-linear non Gaussian problems (Katzfuss, Stroud, and Wikle, 2016).

For non-linear problems of moderate dimension, i.e. those where we replace the right-hand side of both state (Equation (3.4)) and observation (Equation (3.5)) equations by non-linear functions, other variants such as the Extended Kalman filter (EKF) (Jazwinski, 1970) and the unscented Kalman filter (UKF) (Julier and Uhlmann, 1997) may be used. The EKF applies the Kalman filter to a linearization of the non-linear system around the current conditional means  $\hat{X}_{t|t-1}$  and  $\hat{X}_{t|t}$ . If the systems dynamics are highly non-linear, this approximation can fail. Alternatively, the UKF, which is based on the unscented transform, directly approximates the predicted means and covariance matrix, by constructing a set of deterministic points that are propagated through the systems dynamics.

The Kalman smoother (Algorithm 2) computes the marginal distributions  $X_t | Y$  for  $t = 0, \dots, n$ . Upon closer inspection, the mean and covariance updates resemble that of the Kalman filter (Algorithm 1). This is no coincidence: By the assumed dependence structure (Figure 3.1, except for the dashed arrows), we obtain the following lemma, which will allow us to prove the recursions.

**Lemma 3.2** (conditional independence from future observations). *Let  $t \in \{0, \dots, n-1\}$  and  $s > t$ . In a GLSSM, conditional on  $X_{t+1}$ ,  $X_t$  is independent of  $Y_s$ ,  $s > t$ .*

*Proof.* As  $s > t$ , we have

$$p(x_t, y_s | x_{t+1}) = p(y_s | x_{t+1}, x_t) p(x_t | x_{t+1}) = p(y_s | x_{t+1}) p(x_t | x_{t+1})$$

where the second equality follows from the dependency structure of the model.  $\square$

algorithmen konsistent mit gets und =, return value

algorithmen konsistent t, t-1

---

**Algorithm 2** Kalman smoother. Note that the Kalman filter already outputs the smoothed last state  $\hat{X}_{n|n}$  and covariance  $\Xi_{n|n}$ .

---

**Require:** GLSSM (Definition 3.2), outputs from Kalman filter (Algorithm 1)

```

1: for  $t \leftarrow n-1, \dots, 0$  do
2:    $G_t = \Xi_{t|t} A_t \Xi_{t+1|t}^{-1}$ 
3:    $\hat{X}_{t|n} = \hat{X}_{t|t} + G_t (\hat{X}_{t+1|n} - \hat{X}_{t+1|t})$ 
4:    $\Xi_{t|n} = \Xi_{t|t} - G_t (\Xi_{t+1|t} - \Xi_{t+1|n}) G_t^T$ 
5: end for
```

---

We can now sketch the proof for the Kalman smoother recursions, based on the arguments in (Chopin and Papaspiliopoulos, 2020, Chapter 7.3). By the preceding lemma, the conditional distribution of  $X_t$  given  $Y_{:n}$  and  $X_{t+1}$  is the same as that given  $Y_{:t}$  and  $X_{t+1}$ . We may now regard  $X_{t+1} = A_t X_t + u_t + \varepsilon_{t+1}$  as an additional observation at time  $t$ , and use the Kalman filter update to determine this conditional distribution:

$$X_t | Y_{:n} = y_{:n}, X_{t+1} = x_{t+1} \sim \mathcal{N}(\hat{X}_{t|t} + G_t(x_{t+1} - \hat{X}_{t+1|t}), \Xi_{t|t} - G_t \Xi_{t+1|t} G_t^T).$$

As  $\hat{X}_{t|t}$  and  $\hat{X}_{t+1|t}$  are linear functions of  $Y_{:n}$  (actually  $Y_{:t}$ ), we may apply the last statement of Lemma 3.1, to see that, conditional on  $Y_{:n} = y_{:n}$ , the distribution of  $X_t$  is Gaussian with mean

$$\hat{X}_{t|t} + G_t (\hat{X}_{t+1|n} - \hat{X}_{t+1|t})$$

and covariance matrix

$$\Xi_{t|t} - G_t \Xi_{t+1|t} G_t^T + G_t \Xi_{t+1|n} G_t^T = \Xi_{t|t} - G_t (\Xi_{t+1|t} - \Xi_{t+1|n}) G_t^T.$$

These quantities are calculated by the Kalman smoother (Algorithm 2).

Going back to the proof of the last statement in Lemma 3.1, we see that we can represent  $X_t$  as

$$X_t = \hat{X}_{t|t} + G_t(X_{t+1} - \hat{X}_{t+1|t}) + \xi_t, \quad (3.6)$$

for a  $\xi_t \sim \mathcal{N}(0, \Xi_{t|t} - G_t \Xi_{t+1|t} G_t^T)$  which is independent of  $Y_{:n}$  and  $X_{t+1}$ . This recurrence may be used to generate samples from the joint smoothing distribution, which is useful if one is interested in non-linear functionals of the smoothing distribution that involve multiple states at once, such as a moving median or maximum. It is based on the following decomposition of the smoothing density

$$p(x|y) = p(x_n|y) \prod_{t=n-1}^0 p(x_t | x_{t+1}, y_{:t}).$$

The resulting algorithm is called the Forwards Filter, Backwards Sampling (FFBS) (Algorithm 3) and was first described in (Frühwirth-Schnatter, 1994) in the context of a data augmentation algorithm for Bayesian analysis of GLSSM.



---

**Algorithm 3** Forwards filter, backwards smoother (Frühwirth-Schnatter, 1994, Proposition 1)

---

**Require:** GLSSM (Definition 3.2), outputs from Kalman filter (Algorithm 1)

- 1: Simulate  $\hat{X}_{n|n} \sim \mathcal{N}(\hat{X}_{n|n}, \Xi_{n|n})$
  - 2: **for**  $t \leftarrow n - 1, \dots, 0$  **do**
  - 3:      $G_t = \Xi_{t|t} A_t \Xi_{t+1|t}^{-1}$
  - 4:     Simulate  $\xi_t \sim \mathcal{N}(0, \Xi_{t|t} - G_t \Xi_{t+1|t} G_t^T)$
  - 5:     Set  $\tilde{X}_{t|n} = \hat{X}_{t|t} + G_t (\hat{X}_{t+1} - \hat{X}_{t+1|t}) + \xi_t$
  - 6: **end for**
- 

**Remark 3.4** (regularity of  $\Sigma_t$  and  $\Omega_t$ ). Throughout this section, we have assumed, either explicitly or implicitly, that the innovation and observation covariance matrices  $\Sigma_t$  and  $\Omega_t$  are non-singular, i.e. SPD.

For the Kalman filter we require that for every  $t$ ,  $\Psi_{t|t-1}$  is non-singular, i.e. that we can apply Lemma 3.1 (i). This is fulfilled as soon as  $\Omega_t$  is non-singular, which is a reasonable assumption in most models. Following the remark after Lemma 3.1, we could also replace  $\Psi_{t|t-1}^{-1}$  in Algorithm 1 by its Moore-Penrose inverse.

A similar argument can be made for singular  $\Xi_{t+1|t}$ , where we replace  $\Xi_{t+1|t}^{-1}$  by its Moore-Penrose inverse in the Kalman smoother (Algorithm 2) and the FFBS (Algorithm 3).

In the context of COVID-19, variants of the Kalman filter have been employed to analyse the time-varying behavior of epidemiological parameters. Usually the models start from some theoretical, e.g. compartmental, model of how the epidemic spreads. After time-discretization and possibly linearization, one ends up with a GLSSM, to which the Kalman filter or one of its variants may be applied. In (Arroyo-Marioli et al., 2021) the authors construct a simple GLSSM to reconstruct the time-varying reproduction number from observed growth factors, exploiting the linear relationship between the two quantities in the SIR compartmental model and using the Kalman filter and smoother to perform inference. (Song et al., 2021; Zhu et al., 2021) directly apply the EKF to time-discretized compartmental models, fitting them either to simulated (Zhu et al., 2021) or real (Song et al., 2021) data. Similarly, (Engbert et al., 2020) use the EnKF to fit a stochastic compartmental model to German regional data, where the EnKF allows to deal with the non-linear and non-Gaussian properties on these small spatial scales.

The attractive feature of GLSSMs is that a large part of inference is analytically feasible: we can calculate the likelihood, smoothing distribution and sample from it. However, the modeling capacity of GLSSMs is limited: most interesting phenomena in the context of this thesis follow neither linear dynamics nor are well modeled by a Gaussian distribution.

Nevertheless, linearization of non-linear dynamics suggests that GLSSMs can have some use as approximations to these more complicated phenomena, provided they are sufficiently close to Gaussian models, e.g. unimodal and without heavy tails. We start to move away from linear Gaussian models by allowing observations that are non-Gaussian.

## 3.2 Partially Gaussian state space models

For the applications considered in this thesis the distribution of observations is never Gaussian — see ?? — and all we can hope for is that the data-generating mechanism is close enough to a Gaussian distribution that inferences made in a GLSSM may carry over. For epidemiological models, Gaussian distributions may be appropriate if incidences are high, e.g. during large outbreaks in a whole country. When case numbers are small, the discrete nature of incidences is better captured by a distribution on  $\mathbf{N}_0$ , and standard distributions used are the Poisson and negative binomial distributions, see e.g. (Lloyd-Smith et al., 2005). We thus want SSMs where observations are allowed to follow these non-Gaussian distributions.

Concerning the distribution of states, we keep the linear Gaussian assumption, i.e. Equation (3.4).



As demonstrated in Chapter 4, using Gaussian states and transitions allows for flexible modeling of many epidemiological desiderata. Furthermore, keeping the states Gaussian will enable us to use Efficient Importance Sampling (EIS) effectively, by constructing approximations via GLSSM which possess the same state dynamics. Alternatively,  $t$ -distributed innovations or more general transition kernels could be employed and we refer the interested reader to (Durbin and Koopman, 2012, Part II) for a selection of these models. The following definition is that of (Koopman, Lit, and Nguyen, 2019), which itself is an extension of earlier work of (Shephard, 1994). (Shephard, 1994) considered only SSMs where, conditional on another Markov process  $Z = (Z_t)_{t=0, \dots, n}$ , model is a full GLSSM, which allows for efficient inference if the conditional distribution  $Z|(X, Y)$  is tractable. As their definition involves a conditional GLSSM, the observations still take values in  $\mathbf{R}^p$ , not  $\mathbf{N}^p$  as is necessary for our endeavors. Thus we opt for the definition presented in (Koopman, Lit, and Nguyen, 2019), where we replace the Gaussian observations (Equation (3.5)) with arbitrary distributions.

**Definition 3.3** (Partially Gaussian state space model (PGSSM)). A Partially Gaussian state space model (PGSSM) is a joint distribution for  $(X, Y)$  where states  $X$  follow Equation (3.4), i.e.

$$X_{t+1} = A_t X_t + u_t + \varepsilon_{t+1} \quad t = 0, \dots, n-1,$$

with  $X_0 \sim \mathcal{N}(0, \Sigma_0)$ ,  $\varepsilon_t \sim \mathcal{N}(0, \Sigma_t)$ ,  $u_t \in \mathbf{R}^m$  for  $t = 1, \dots, n$  and  $X_0, (\varepsilon_t)_{t=1, \dots, n}$  jointly independent.

Furthermore, the observations  $Y$  form a conditional Markov process, conditional on states  $X$ , where the conditional densities of observations, given states admit the following form

$$p(y|x) = \prod_{t=0}^n p(y_t|x_t, y_{t-1}),$$

with respect to the dominating measure  $\bigotimes_{t=0}^n \mu_Y$ . Here  $p(y_t|x_t, y_{t-1})$  are allowed to take any arbitrary density<sup>1</sup>.

It is straightforward to check that a PGSSM is indeed a SSM.

**Remark 3.5.** Recalling Remark 3.1, if our main interest lies in the conditional distribution  $X|Y = y$  for a fixed set of observations  $y$ , it will suffice to consider models where

$$p(y|x) = \prod_{t=0}^n p(y_t|x_t)$$

holds, and we will do so in the following to enhance readability. At points where this distinction matters, e.g.

add example

, we will give appropriate remarks.

Both the Poisson and negative binomial distribution belong to the class of exponential family distributions. As such, their densities have a convenient structure, allowing only for a linear interaction between the natural parameter and the densities' argument. We refer to (Brown, 1986) for a comprehensive treatment of exponential families and use their definitions throughout this section.

**Definition 3.4** (exponential family). Let  $\mu$  be a  $\sigma$ -finite measure on  $\mathbf{R}^p$  and denote by

$$\Psi = \left\{ \psi \in \mathbf{R}^p : \int \exp(\psi^T y) \, d\mu(y) < \infty \right\}$$

the set of parameters  $\psi$  such that the moment-generating function of  $\mu$  is finite. For every  $\psi \in \Psi$

$$p_\psi(y) = Z(\psi)^{-1} \exp(\psi^T y)$$

<sup>1</sup>Recall that we have not specified  $\mu_Y$ , so it is always possible to use  $p = \mathbf{1}_Y$ , the constant function.

defines a probability density with respect to the measure  $\mu$ , where

$$Z(\psi) = \int \exp(\psi^T x) \, d\mu(y)$$

is the normalizing constant. We call both the densities  $p_\psi$  and induced probability measures

$$\mathbf{P}_\psi(A) = \int_A p_\psi(y) \, d\mu(y),$$

for measurable  $A \subset \mathbf{R}^p$ , a **standard exponential family**.

Conversely, let  $\mathbf{P}_\psi, \psi \in \Psi$  be a given parametric family of probability measures on some space  $\mathcal{Y}$  that is absolutely continuous with respect to a common dominating measure  $\mu$ . Suppose there exist a reparametrization  $\eta : \Psi \rightarrow \mathbf{R}^p$ , a statistic  $T : \mathcal{Y} \rightarrow \mathbf{R}^p$  and functions  $Z : \Psi \rightarrow \mathbf{R}$ ,  $h : \mathcal{Y} \rightarrow \mathbf{R}$ , such that

$$p_\psi(y) = \frac{d\mathbf{P}_\psi}{d\mu} = Z(\psi)h(y) \exp(\eta(\psi)^T T(y)),$$

then we call  $(\mathbf{P}_\psi)_{\psi \in \Psi}$  and  $(p_\psi)_{\psi \in \Psi}$  a  **$p$ -dimensional exponential family**. If  $\eta(\psi) = \psi$  is the identity, we call  $\psi$  the natural parameter. If  $T(y) = y$ , we call  $y$  the natural observation. If  $\psi$  is the natural parameter and  $y$  the natural observation, we call  $(\mathbf{P}_\psi)_{\psi \in \Psi}$  a **natural exponential family**. By reparametrization (in  $\psi$ ) and sufficiency (in  $y$ ) every  $p$ -dimensional exponential family can be written as an equivalent standard exponential family, see the elaborations in (Brown, 1986, Chapter 1).

Exponential families have the attractive property that they are log-concave in their parameters. As such the Fisher-information is always positive semidefinite, which will be crucial in defining surrogate Gaussian models in Section 3.5.

**Lemma 3.3** (log-concavity of exponential family distributions). *Let  $(p_\psi)_{\psi \in \Psi}$  be a natural  $p$ -dimensional exponential family and  $\Psi$  convex and open in  $\mathbf{R}^p$ . In this case  $\psi \mapsto \log p_\psi(y)$  is concave for every  $y \in \mathbf{R}^p$ .*

*Proof.* As  $\log p_\psi(y) = -\log Z(\psi) + \psi^T y$  it suffices to show that  $\psi \mapsto \log Z(\psi)$  is convex. However,

$$\psi \mapsto \log Z(\psi) = \log \int \exp(\psi^T y) \, d\mu(y)$$

is the cumulant generating function of the base measure  $\mu$ , which is convex (Billingsley, 1995, p. 144f).  $\square$

Additionally, the moment generating function  $\psi \mapsto Z(\psi)$  is smooth on the interior of  $\Psi$  and allows to switch the order of integration and differentiation.

**Theorem 3.1** ((Brown, 1986, Theorem 2.2, Corollary 2.3)). *Let  $\psi \in \text{int } \Psi$  be an interior point. Then the moment generating function  $Z : \Psi \rightarrow \mathbf{R}$  is infinitely often differentiable with derivatives*

$$\frac{\partial^{|\alpha|}}{\partial^\alpha \psi} Z(\psi) = \int y^\alpha \exp(\psi^T y) \, d\mu(y)$$

for any multi-index  $\alpha \in \mathbf{N}^k$ .

Additionally, the gradient of  $\log Z$ ,  $\nabla_\psi \log Z(\psi)$  is given by

$$\nabla_\psi \log Z(\psi) = \mathbb{E}T(X),$$

and the Hessian of  $\log Z$ ,  $H_\psi \log Z(\psi)$  by

$$H_\psi \log Z(\psi) = \text{Cov}(T(X)),$$

where  $X \sim \mathbf{P}_\psi$ .

**Example 3.1** (Poisson & negative binomial distribution). Both the family of Poisson distributions, parameterized by rate  $\lambda$  and the negative binomial distribution, parameterized by success probability  $q$  with fixed overdispersion  $r$  form an exponential family.

The log-density of the Poisson distribution with rate  $\lambda$ ,  $\text{Pois}(\lambda)$  w.r.t. the counting measure on  $\mathbf{N}_0$  is

$$\log p_\lambda(x) = -\lambda + x \log \lambda - \log x!.$$

Thus the Poisson distribution forms an exponential family with natural parameter  $\log \lambda$ , natural statistic  $\text{id}$  (the identity), base measure  $h(x) = \frac{1}{x!}$  and moment-generating function  $Z(\lambda) = \exp(-\lambda)$ .

The log-density of the negative binomial distribution with overdispersion parameter  $r$  and success probability  $q$   $\text{NegBinom}(q, r)$  is

$$\log p_{q,r}(x) = \log \binom{x+r-1}{x} + x \log(1-q) + r \log q.$$

For fixed  $r$  these distributions form an exponential family with natural parameter  $\log(1-q)$ , natural statistic  $T = \text{id}$ , base measure  $h(x) = \binom{x+r-1}{x}$  and moment-generating function  $Z(q) = r \log q$ .

In this parametrization the mean of the  $\text{NegBinom}(q, r)$  distribution is  $\mu = r \frac{1-q}{q}$  and its variance is  $r \frac{1-q}{q^2}$ . An alternative parametrization that will become useful Chapter 4 is that by the log mean  $\xi = \log \mu$  and overdispersion  $r$ . As  $q = \frac{r}{r+\mu}$ , this parametrization has log-density

$$\log p_{r,\xi}(x) = \log \binom{x+r-1}{x} + x\xi - (r+x) \log(\exp \xi + r) - r \log r,$$

which does not form a natural exponential family. However, it retains the log-concavity of Lemma 3.3, as a quick calculation reveals that

$$\partial_{\xi^2}^2 \log p_{r,\xi}(x) = -(r+x) \frac{r \exp(-\xi)}{(r \exp(-\xi) + 1)^2} < 0$$

for all  $x \in \mathbf{N}_0$ .

The models we study in Chapter 4 belong, for the most part, to the following subclass of PGSSM models.

**Definition 3.5** (Exponential Family Partially Gaussian state space model (EGSSM)). An Exponential Family Partially Gaussian state space model (EGSSM) is a PGSSM where the conditional distribution of  $Y_t$  given  $X_t$  comes from an exponential family with respect to a base measure  $\mu_t$ , i.e.

$$p(y_t|x_t) = h_t(y_t) Z_t(x_t) \exp(\eta_t(x_t)^T T_t(y_t))$$

for suitable functions  $h_t, Z_t, \eta_t, T_t$ . If  $Y_t$  in the PGSSM is allowed to depend on the previous  $Y_{t-1}$ , the functions  $h_t, Z_t, \eta_t$  and  $T_t$  may depend on  $y_{t-1}$ .

If, additionally, matrices  $B_t \in \mathbf{R}^{p \times m}$  exist, such that for the signal  $S_t = B_t X_t \in \mathbf{R}^p$ ,  $Y_t$  only depends on  $X_t$  through  $S_t$ , i.e. it holds

$$p(y_t|x_t) = \prod_{i=1}^p h_t^i(y_t^i) Z_t^i(s_t) \exp(\eta_t^i(s_t^i) T(y_t^i)),$$

for functions  $h_t^i : \mathbf{R} \rightarrow \mathbf{R}, Z_t^i : \mathbf{R} \rightarrow \mathbf{R}, \eta_t^i : \mathbf{R} \rightarrow \mathbf{R}, T : \mathbf{R} \rightarrow \mathbf{R}, i = 1, \dots, p$ , we say the EGSSM has a **linear signal**, similar to the treatment in (Durbin and Koopman, 2012, Part II).

**Remark 3.6.** To simplify notation we will usually assume that the functions  $h, Z$  and  $T$  are the same for all  $t$  (and  $i$ , if the EGSSM has a linear signal) and drop in our notation the dependence of  $h, Z$ , and  $T$  on  $t$  (and  $i$ ). Similarly, we assume that the base measure  $\mu_t$  is the same for all  $t$ .

From Lemma 3.3, we immediately obtain the following results (Durbin and Koopman, 2012, Section 10.6.4)

**Lemma 3.4** (log-concavity of the smoothing distribution). *Consider an EGSSM, where  $\eta_t = \text{id}$  for all  $t$ . Then  $x \mapsto \log p(x|y)$  is concave for  $\mu_Y$ -a.e.  $y$ .*

*Proof.* We may write

$$\log p(x|y) = \log p(y|x) + \log p(x) - \log p(y),$$

where the last term does not depend on  $x$ .  $\log p(x)$  is concave in  $x$ , as  $p(x)$  is the joint density of a multivariate Gaussian distribution. Furthermore

$$\log p(y|x) = \sum_{t=0}^n \log p(y_t|x_t, y_{t-1}),$$

which, by Lemma 3.3 is concave in  $x$ . □

Notice that the dependence of  $Y_t$  on  $Y_{t-1}$  does not influence the statement of this lemma, as we are interested in properties of  $x \mapsto p(x|y)$ .

As in the previous chapter, after having observed  $Y$ , one is interested in the conditional distribution of states  $X$ , given  $Y$ . If the observations are not Gaussian, this is a difficult task as the distribution is not analytically tractable. Instead, approximations, e.g. the Laplace approximation (LA), which will exploit the log-concavity developed here or simulation-based inference, e.g. importance sampling (Sections 3.3 and 3.5), sequential Monte Carlo (Chopin and Papaspiliopoulos, 2020) or MCMC-methods (Brooks et al., 2011) are used. Similarly, fitting hyperparameters  $\psi$  by maximum likelihood inference becomes more difficult as evaluating  $\ell(\psi) = p(y) = \int p(x, y) dx$  is not analytically available, thus requiring numerical or simulation methods for evaluation and gradient descent or EM-techniques for optimization, see Section 3.6.

In this thesis, we will focus on importance sampling methods, which are the focus of the next section.

### 3.3 Importance Sampling

Importance sampling is a simulation technique that allows us to approximate integrals w.r.t a measure of interest, the target, by sampling from a tractable approximation, the proposal, instead, thus performing Monte-Carlo integration. To account for the fact that we did not sample from the correct probability measure, we weight samples according to their importance. As the user has freedom in the choice of approximation (except for some technical conditions), importance sampling also acts as a variance reduction technique with better approximations resulting in smaller Monte-Carlo variance. Thus the role that importance sampling plays is twofold: first, it enables Monte-Carlo integration even if sampling from the target is not possible, and second it allows us to do so in an efficient way by choosing, to be defined precisely below, the approximation in an optimal way.

Alternative approaches to importance sampling for performing inference in SSMs include Markov chain Monte Carlo (MCMC) and SMC. Recall from the introduction to this chapter that this inference concerns three objectives: maximum likelihood estimation, i.e. evaluation and optimization of the likelihood, access to the posterior distribution  $X_{:n}|Y_{:n}$  and prediction of future states and observations. Let us give a concise comparison of these alternative approaches, weighing their advantages and disadvantages over importance sampling, in particular for the SSMs that this thesis deals with.

MCMC (Brooks et al., 2011) is a simulation technique that allows to simulation of correlated samples from a target distribution by constructing an ergodic Markov chain that has as its invariant distribution the desired distribution. If one is able to simulate from such a Markov chain, one can generate samples whose marginal distributions are close to the target distribution. Thus, these samples can be used in MC-integration to estimate expectations of interest, though one has to be mindful of autocorrelation of these samples. For standard variants of MCMC, such as Metropolis-Hastings MCMC or Hamiltonian Monte Carlo, one needs access to the density of the

sought after distribution up to a constant to simulate a step in the Markov chain. While these methods are very general, in high dimensions, these are affected by the curse of dimensionality.

Let us argue for our choice of using IS over MCMC for estimating conditional expectations of the form  $\mathbb{E}(f(X)|Y)$  for PGSSMs. For the models we consider in this thesis, the dimension of  $X$  ( $(n+1)m$ ) can become quite large, so MCMC suffers from the aforementioned curse of dimensionality. IS can also suffer from this curse, especially if the proposal is far from the target. If, however, the proposal is close to the target, IS can perform surprisingly well, see e.g. (Chopin and Ridgway, 2017) where it is used as the gold standard method against which other methods are benchmarked.

As IS is based on independent samples, it can be parallelized easily, whereas parallelizing MCMC is more involved, using e.g. (Neiswanger, Wang, and Xing, 2014). Additionally, analysis of convergence is much simpler than that of MCMC, which requires consideration of burn-in samples, autocorrelation of samples and investigating trace plots for the chain getting stuck.

SMC (Chopin and Papaspiliopoulos, 2020) or particle filters, use sequential importance sampling to provide a particle approximation to the filtering distributions  $X_t|Y_{:t}$ , essentially decomposing the problem into a  $n$  importance sampling steps. To avoid particle collapse, SMC is usually equipped with a resampling step once the effective sample size of the current set of particles drops below a specified level. Once the final filtering distribution  $X_n|Y_{:n}$  is approximated, the smoothing distribution may be obtained in several ways, e.g., backwards sampling or a two-filter approach, see (Chopin and Papaspiliopoulos, 2020, Chapter 12).

Conveniently, SMC allows us to approximate the likelihood  $\ell(\theta)$  for a single parameter by a single pass of the particle filter. However, the discrete nature of resampling makes the approximated likelihood non-continuous, complicating maximum likelihood inference. (Chopin and Papaspiliopoulos, 2020, Chapter 14) discusses several strategies: the first amounts to importance sampling of the order as discussed in this thesis, where one fixes a reference parameter  $\theta_0$  to perform importance sampling with  $p_{\theta_0}(x|y)$  against  $p_{\theta}(x|y)$ . The second strategy only works in the univariate case and consists of approximating the non-continuous inverse CDFs appearing in the resampling step by continuous ones. Finally, if the dependence on the hyperparameters  $\theta$  allows for application of the EM-algorithm, it may be used to perform the optimization. Contrary to SMC, the global importance sampling approach we discuss in Sections 3.5 and 3.6 allows us to perform importance sampling in an optimal way, and allows for use of numerical differentiation as the dependence of  $\log p_y(\theta)$  on  $\theta$  is smooth, as there is no resampling involved.

This chapter proceeds with a general treatment of importance sampling, loosely based on (Chopin and Papaspiliopoulos, 2020, Chapter 8) and (Durbin and Koopman, 2012, Chapter 11). Subsequently, we will focus our attention on methods to obtain good importance sampling proposals.

Suppose we have a function  $h : \mathcal{X} \rightarrow \mathbf{R}$  whose integral w.r.t. to some measure  $\mu$ ,

$$\zeta = \int_{\mathcal{X}} h(x) d\mu(x),$$

exists and whose value we want to compute. Furthermore, suppose that we can write

$$\int_{\mathcal{X}} h(x) d\mu(x) = \int_{\mathcal{X}} f(x) d\mathbf{P}(x) = \mathbf{P}[f],$$

for a probability measure  $\mathbf{P}$  and function  $f : \mathcal{X} \rightarrow \mathbf{R}$ , e.g. because  $\mathbf{P} = p\mu$  and  $h(x) = f(x)p(x)$   $\mu$ -a.s.. Here, and in the remainder of this chapter, we use the operator shorthand notation  $\mathbf{P}[f] = \int f d\mathbf{P}$  for a measure  $\mathbf{P}$  and a  $\mathbf{P}$ -integrable function  $f$ . Let  $\mathbf{G}$  be a another probability measure on  $\mathcal{X}$  such that  $f\mathbf{P}$  is absolutely continuous with respect to  $\mathbf{G}$ ,  $f\mathbf{P} \ll \mathbf{G}$ , and let  $v = \frac{d f\mathbf{P}}{d \mathbf{G}}$  be the corresponding Radon-Nikodym derivative. Then

$$\zeta = \mathbf{P}[f] = \int_{\mathcal{X}} f(x) d\mathbf{P}(x) = \int_{\mathcal{X}} \left( \frac{d f\mathbf{P}}{d \mathbf{G}} \right) (x) d\mathbf{G}(x) = \mathbf{G}[v]$$

which suggests to estimate  $\zeta$  by Monte-Carlo integration:

$$\hat{\zeta} = \frac{1}{N} \sum_{i=1}^N v(X^i),$$

the importance sampling estimate of  $\zeta$ . The importance samples  $X^i, i = 1, \dots, N$  have distribution  $\mathbf{G}$ , and will usually be i.i.d. For this procedure to work, we want  $\hat{\zeta}$  to fulfill a law of large numbers and a central limit theorem, so we will want  $v \in L^2(\mathbf{G})$ , where  $L^p(\nu)$  is the space of  $p$ -times  $\nu$ -integrable functions for a measure  $\nu$ . The i.i.d. assumption could also be dropped, e.g. when we employ antithetic variables, see (Ripley, 2009, Section 5.3) and Section 3.6. Here we call  $\hat{\zeta}$  the importance sampling estimate of  $\zeta$ .

If  $v \in L^2(\mathbf{G})$  and under i.i.d. sampling the Monte-Carlo variance of  $\hat{\zeta}$  is  $\frac{\text{Var}(v(X^i))}{N}$ , and so naturally we want  $\text{Var}(v(X^i))$  to be small to ensure fast convergence of  $\hat{\zeta}$ . As  $v$  depends on the proposal  $\mathbf{G}$ , and we have flexibility in choosing  $\mathbf{G}$ , importance sampling acts as a variance reduction technique: the better  $\mathbf{G}$  approximates  $f\mathbf{P}$ , in the sense that the variance of  $v$  w.r.t.  $\mathbf{G}$  is small, the faster importance sampling will converge.

A classical result is that the minimum MSE proposal  $\mathbf{G}^*$  has a closed form. Indeed it is given by the total variation measure of  $f\mathbf{P}$ , renormalized to be a probability measure, which can be shown by a simple application of Jensen's inequality.

**Proposition 3.1** (Chopin and Papaspiliopoulos, 2020, Proposition 8.2). *[minimum MSE proposal] The proposal  $\mathbf{G}^*$  that minimizes the MSE of importance sampling is given by*

$$\mathbf{G}^* = \frac{|f|}{\mathbf{P}[|f|]} \mathbf{P}.$$

Unfortunately, this optimality result has no practical use, indeed if  $f$  is positive we would need to obtain  $\mathbf{P}[f]$  first, the overall target of our endeavor. Additionally, sampling from  $\mathbf{G}^*$  is not guaranteed to be practically feasible.

If the Radon-Nikodym derivative  $w = \frac{d\mathbf{P}}{d\mathbf{G}}$  exists, then  $v = fw$ , which, for the problems we will study, is the case. Then

$$\hat{\zeta} = \frac{1}{N} \sum_{i=1}^N f(X^i)w(X^i),$$

where  $w(X^i)$  is called the importance weight, or just weight, of the  $i$ -th sample. If the samples are clear from the context we sometimes write  $w^i = w(X^i)$ . This motivates us to regard

$$\hat{\mathbf{P}}_N = \frac{1}{N} \sum_{i=1}^N w(X_i) \delta_{X_i}, \quad (3.7)$$

as a particle approximation of  $\mathbf{P}$ , in the sense that for sufficiently well behaved test functions  $f$ , as  $N \rightarrow \infty$

$$\hat{\mathbf{P}}_N[f] = \frac{1}{N} \sum_{i=1}^N f(X^i)w(X^i) \rightarrow \mathbf{P}[f].$$

We will return to the question of which functions  $f$  to consider further below and assume in the following discussion  $fw \in L^2(\mathbf{G})$ .

To perform importance sampling one must be able to evaluate  $w$ . In the context of PGSSMs this is usually not possible: if  $\mathbf{P}$  is the intractable conditional distribution of  $X|Y$ , then the integration constant of its density  $p(y)$  is not analytically available. Still, we can usually evaluate the weights up to a constant, i.e.

$$\tilde{w}(x) \propto \frac{d\mathbf{P}}{d\mathbf{G}}(x)$$

is available. The missing constant is then  $\mathbf{G}\tilde{w}$ , which is itself amenable to importance sampling: we may estimate it by  $\frac{1}{N} \sum_{i=1}^N \tilde{w}(X^i)$ . This leads to the so-called self-normalized importance sampling weights

$$W_i = \frac{w(X^i)}{\sum_{i=1}^N w(X^i)},$$

Monte Carlo estimates

$$\hat{\zeta} = \sum_{i=1}^N W_i f(X^i),$$

and particle approximation

$$\hat{\mathbf{P}}_N = \sum_{i=1}^N W_i \delta_{X^i}.$$

Unless  $\tilde{w}$  is degenerate, i.e. constant,

$$\hat{\zeta} = \frac{\sum_{i=1}^N \tilde{w}(X^i) f(X^i)}{\sum_{i=1}^N \tilde{w}(X^i)}$$

is a ratio of two non-constant, unbiased estimators and so is itself biased. Nevertheless, noticing that the rescaled denominator  $\frac{1}{N} \sum_{i=1}^N \tilde{w}(X^i)$  consistently estimates the integration constant  $\mathbf{G}\tilde{w}$ , allows us to apply Slutsky's lemma and obtain a central limit theorem for  $\hat{\zeta}$  (recall that we assumed  $f w \in L^2(\mathbf{G})$ ).

The class for test functions  $f$  for which this holds depends on  $\mathbf{P}$  and  $\mathbf{G}$ . (Agapiou et al., 2017) study the behavior of uniformly bounded test functions  $\|f\| \leq 1$ . For these functions it suffices that  $w \in L^2(\mathbf{G})$  to ensure asymptotic normality of  $\zeta$ . Thus an important quantity is

$$\rho = \frac{1}{(\mathbf{G}\tilde{w})^2} \mathbf{G}[\tilde{w}^2] = \mathbf{G}[w^2] = \mathbf{P}[w],$$

the second moment of the importance sampling weights. (Agapiou et al., 2017) show that the bias

$$\left| \mathbb{E}(\hat{\mathbf{P}}_N - \mathbf{P})[f] \right|$$

and mean-squared error (MSE)

$$\mathbb{E} \left( (\hat{\mathbf{P}}_N - \mathbf{P})[f] \right)^2$$

of importance sampling are both, for bounded  $f$ , of order  $\mathcal{O}(\frac{\rho}{N})$ . Here the expectation  $\mathbb{E}$  is with respect to the random particles  $X^1, \dots, X^N$ . Consequently, for bounded functions, keeping  $\frac{\rho}{N}$  small produces importance sampling estimates with small bias and MSE. This can be achieved in two ways: either we choose  $\mathbf{G}$  „close enough“ to  $\mathbf{P}$  to ensure small  $\rho$ , or we choose  $N$  large enough to compensate for a large  $\rho$ .

Applying Jensen's inequality, we see that

$$\mathcal{D}_{\text{KL}}(\mathbf{P}||\mathbf{G}) = \mathbf{P}[\log w] \leq \log \mathbf{P}[w] = \log \rho,$$

so small  $\rho$  implies a small KL-divergence between  $\mathbf{P}$  and  $\mathbf{G}$  as well. Conversely, the following theorem of Chatterjee and Diaconis implies that a small KL-divergence is both sufficient and necessary for importance sampling to perform well.

**Theorem 3.2** (Chatterjee and Diaconis, 2018, Theorem 1.1). *Let  $\mathbf{P}$  and  $\mathbf{G}$  be probability measures on a measurable space  $(\mathcal{X}, \mathcal{B})$  such that  $\mathbf{P} \ll \mathbf{G}$  and let  $f \in \mathbf{L}^2(\mathbf{P})$  be a function with  $\|f\|_{L^2(\mathbf{P})} = (\mathbf{P}f^2)^{1/2} < \infty$ . Let  $Y$  be an  $\mathcal{X}$  valued random variable with law  $\mathbf{P}$ .*

*Let  $L = \mathcal{D}_{\text{KL}}(\mathbf{P}||\mathbf{G}) = \mathbb{E} \log w(Y)$  be the KL-divergence between  $\mathbf{P}$  and  $\mathbf{G}$ , and let*

$$\hat{\mathbf{P}}_N = \sum_{i=1}^N w(X^i) \delta_{X^i}$$

*be the particle approximations of  $\mathbf{P}$  based on samples  $X^1, \dots, X^N \stackrel{i.i.d}{\sim} \mathbf{G}$ ,  $N \in \mathbb{N}$ .*

If the sample size  $N$  is given by  $N = \exp(L + t)$  for a  $t \geq 0$ ,

$$\mathbb{E} \left| \hat{\mathbf{P}}_N[f] - \mathbf{P}[f] \right| \leq \|f\|_{L^2(\mathbf{P})} \left( \exp(-t/4) + 2\sqrt{\mathbb{P}(\log w(Z) > L + t/2)} \right). \quad (3.8)$$

Conversely, if  $N = \exp(L - s)$  for  $s \geq 0$ , then for any  $\delta \in (0, 1)$

$$\mathbb{P}(\hat{\mathbf{P}}_N[\mathbf{1}] \geq 1 - \delta) \leq \exp\left(-\frac{s}{2}\right) + \frac{\mathbb{P}(\log w(Z) \leq L - \frac{s}{2})}{1 - \delta}, \quad (3.9)$$

where  $\mathbf{1}$  is the constant function  $x \mapsto 1$ .

Notice the boldface  $\mathbb{P}$  and  $\mathbb{E}$  to differentiate the measures  $\mathbf{P}$  and  $\mathbf{G}$  from expectations and probabilities with respect to the abstract probability space  $(\Omega, \mathcal{A}, \mathbb{P})$  where the random variables  $X_1, \dots, X_N$  and  $Y$  live.

The proof of this theorem is based on splitting  $\mathcal{X}$  into  $\{\log w \leq L + \frac{t}{2}\}$  and its complement and straightforward, it may be found in the Appendix of (Chatterjee and Diaconis, 2018). Theorem 1.2 in the same paper provides a qualitatively similar result for autonormalised importance sampling.

Let us consider the implications of Theorem 3.2, starting with Equation (3.8), by devising heuristics to decide when  $\mathbf{G}$  is a good proposal for fixed sample size  $N$ , and assume for simplicity that  $\|f\|_{L^2(\mathbf{P})} = 1$ . First of all, as  $t = \log N - L$ , we have  $\exp(-t/4) = \exp(L/4)N^{-\frac{1}{4}}$ , so for large  $N$  this term becomes negligible, and the interesting term in inequality (3.8) is the second one. As  $\mathbb{E} \log w(Z) = L$ , this term is a tail probability and we can use standard mass-concentration inequalities to analyze its behavior as  $t$  (and so  $N$ ) grows. Markov's inequality tells us that

$$\mathbb{P}\left(\log w(Z) > L + \frac{t}{2}\right) \leq \frac{L}{L + t/2} = \frac{2}{1 + \frac{\log N}{L}}.$$

Second, if, additionally,  $\log w(Z)$  has finite variance, Chebyshev's inequality yields

$$\mathbb{P}\left(\log w(Z) > L + \frac{t}{2}\right) \leq \frac{4 \operatorname{Var}(\log w(Z))}{t^2} = \frac{4 \operatorname{Var}(\log w(Z))}{(\log N - L)^2}.$$

In both upper bounds provided by the concentration inequalities, all else being equal, a smaller KL-divergence will yield a tighter bound. However, in Chebyshev's inequality, the variance of log weights also plays a role, and will surely be different for different proposals. Assuming  $\mathbf{G} \ll \mathbf{P}$ , we have  $\frac{d\mathbf{G}}{d\mathbf{P}} = \frac{1}{w}$  and so

$$\mathbb{E} \exp(-\log w(Z)) = \mathbb{E} \frac{1}{w(Z)} = \mathbf{P} \left[ \frac{d\mathbf{G}}{d\mathbf{P}} \right] = 1,$$

If the log-weights are bounded from above and below, the following lemma shows that as the variance of  $U = -\log w(Z)$  goes to 0, their mean,

$$\mathbb{E}U = \mathbb{E} -\log w(Z) = -\mathcal{D}_{\text{KL}}(\mathbf{P}||\mathbf{G})$$

goes to 0 as well.

**Lemma 3.5.** *For  $a, b \in \mathbf{R}$ , let  $U \in [a, b]$  be a bounded random variable with variance  $\sigma^2$  and  $\mathbb{E} \exp U = 1$ . Let  $\mu = \mathbb{E}U$  be the mean of  $U$ . Then there exists a  $\delta \in [\exp(a), \exp(b)]$ , such that*

$$0 \geq \mu = \log \left( 1 - \delta \frac{\sigma^2}{2} \right).$$

If, additionally,  $\sigma^2 < \frac{2}{\exp(b)}$  then

$$\mu \geq \log \left( 1 - \exp(b) \frac{\sigma^2}{2} \right).$$



*Proof.* As  $U$  is bounded, all involved expectations exist and are finite. That  $\mu \leq 0$  follows from Jensen's inequality. We perform a first-order Taylor expansion of  $\exp(U - \mu)$ , where the random variable  $\xi$  is between  $U - \mu$  and 0:

$$1 = \exp(\mu) \mathbb{E} \exp(U - \mu) = \exp(\mu) \left( 1 + \mathbb{E}(U - \mu) + \mathbb{E} \left( \frac{(U - \mu)^2}{2} \exp(\xi) \right) \right).$$

Then  $\xi' = \xi + \mu$  is in  $[a, b]$ , and note that, unless  $U = 1$  a.s.,  $\mathbb{E} \exp U = 1$  forces  $a < 0 < b$ . Thus

$$1 = \exp(\mu) + \mathbb{E} \left( \frac{(U - \mu)^2}{2} \exp(\xi') \right),$$

and as  $\xi' \in [a, b]$ , the expectation is in  $\left[ \exp(a) \frac{\sigma^2}{2}, \exp(b) \frac{\sigma^2}{2} \right]$ , i.e.  $\mathbb{E} \left( \frac{(U - \mu)^2}{2} \exp(\xi') \right) = \delta \frac{\sigma^2}{2}$  for some  $\delta \in [\exp(a), \exp(b)]$ . Solving for  $\mu$ , we get

$$\mu = \log \left( 1 - \delta \frac{\sigma^2}{2} \right),$$

as promised.

The second statement follows from  $\delta \leq \exp(b)$  and the monotonicity of log, where the condition ensures that the argument is positive.  $\square$

**Corollary 3.1.** *Let  $\mathbf{P}$  and  $\mathbf{G}$  be equivalent probability measures with bounded Radon-Nikodym derivative  $w = \frac{d\mathbf{P}}{d\mathbf{G}} \in [a, b]$ ,  $a, b \in \mathbf{R}$  and KL-divergence  $\mathcal{D}_{KL}(\mathbf{P}||\mathbf{G}) = \mathbf{P}[\log w]$ .*

*If  $\log w \in L^2(\mathbf{P})$  with variance  $\sigma^2 = \mathbf{P}[(\log w - L)^2]$ , and  $\sigma^2 < \frac{2}{\exp(b)}$ , then*

$$\mathcal{D}_{KL}(\mathbf{P}||\mathbf{G}) \leq -\log \left( 1 - \exp(b) \frac{\sigma^2}{2} \right).$$

Under the assumptions of this corollary, we see that a small variance of the log-weights implies a small KL-divergence, which in turn implies good importance sampling performance.

Let us now discuss the implications of Equation (3.9). We see that for large  $s$ , i.e.  $N \ll \exp(L)$ , the right-hand side is small, and so the probability that importance sampling fails for the constant function is practically relevant. Observe that here

$$\hat{\mathbf{P}}_N[\mathbf{1}] = \frac{1}{N} \sum_{i=1}^N w_i$$

is the mean of weights, which does not have to sum to 1. As a result, Chatterjee and Diaconis recommend to choose  $N = \mathcal{O}(\exp(\mathcal{D}_{KL}(\mathbf{P}||\mathbf{G})))$ .

Based on this discussion, we see that choosing  $\mathbf{G}$  such that either the KL-divergence or the variance of the log-weights is small is sensible. Making the variance small has the additional advantage that it, at least for bounded log-weights, also implies an upper bound for the KL-divergence. We will return to this train of thought when we discuss optimal ways of performing importance sampling, such as the CE-method (minimizing the KL-divergence) and EIS (minimizing the variance of log-weights) in the following sub-chapters.

In practice, we will want to judge whether for an actual sample  $X^1, \dots, X^N \stackrel{\text{i.i.d.}}{\sim} \mathbf{G}$  importance sampling has converged, and there are several criteria available in the literature. The classic effective sample size (ESS)(Kong, Liu, and Wong, 1994)

$$\text{ESS} = \frac{1}{\sum_{i=1}^N W_i^2} \in [1, N]$$

arises from an analysis of the asymptotic efficiency of importance sampling estimates: Consider additional  $Y^1, \dots, Y^N \stackrel{\text{i.i.d.}}{\sim} \mathbf{P}$ , a test function  $f \in L^2(\mathbf{P})$  and assume that  $\rho < \infty$ . We may then estimate  $\zeta = \mathbf{P}f$  in two ways: either by using the importance sampling estimate

$$\hat{\zeta}_{\text{IS}} = \hat{\mathbf{P}}_N(f) = \sum_{i=1}^N W_i f(X^i) = \frac{1}{N} \sum_{i=1}^N (NW_i) f(X^i),$$

or by standard Monte-Carlo integration

$$\hat{\zeta}_{\text{MC}} = \frac{1}{N} \sum_{i=1}^N f(Y^i).$$

(Kong, 1992) applies the delta method to  $\text{Var}(\hat{\zeta}_{\text{IS}})$ , obtaining

$$\text{Var}(\hat{\zeta}_{\text{IS}}) \approx \text{Var}(\hat{\zeta}_{\text{MC}}) (1 + \text{Var}(NW_1)).$$

Note that this approximation does not depend on the specific  $f$  considered, and it is not guaranteed that for large  $N$  the remainder goes to 0, as (Kong, 1992) mentions. In particular, the approximation has to fail whenever  $\text{Var}(\hat{\zeta}_{\text{IS}}) < \text{Var}(\hat{\zeta}_{\text{MC}})$ , i.e. when importance sampling actually performs variance reduction. Nevertheless, whenever the approximation is valid, we may interpret

$$\frac{N}{1 + \text{Var}(NW_1)}$$

as an effective sample size, in the sense that  $N$  times the relative efficiency of  $\hat{\zeta}_{\text{MC}}$  relative to  $\hat{\zeta}_{\text{IS}}$  is approximately given by this expression. As the self-normalized weights  $W_1, \dots, W_N$  are exchangeable and sum to 1, their expected value is  $\mathbb{E}W_1 = \frac{1}{N}$ . Estimating  $\text{Var}(W_1)$  by the unadjusted sample covariance  $\frac{1}{N} \sum_{i=1}^N W_i^2 - \frac{1}{N^2}$  then results in the promised

$$\text{ESS} = \frac{N}{1 + N^2 \left( \frac{1}{N} \sum_{i=1}^N W_i^2 - \frac{1}{N^2} \right)} = \frac{1}{\sum_{i=1}^N W_i^2}.$$

Notice that as the self-normalized weights sum to 1, the ESS is at least 1, as  $0 \leq W_i \leq 1$  and at most  $N$  by the Cauchy-Schwarz inequality.

If we write the ESS in terms of the unnormalized weights  $\tilde{w}$  we see that the efficiency factor (EF)  $\text{EF} = \frac{\text{ESS}}{N}$  fulfills, as  $N \rightarrow \infty$ ,

$$\text{EF} = \frac{\text{ESS}}{N} = \frac{\left( \frac{1}{N} \sum_{i=1}^N \tilde{w}_i \right)^2}{\frac{1}{N} \sum_{i=1}^N \tilde{w}_i^2} \xrightarrow{a.s.} \frac{(\mathbf{G}[\tilde{w}])^2}{\mathbf{G}[\tilde{w}^2]} = \rho^{-1},$$

if  $\tilde{w} \in L^2(\mathbf{G})$  (Agapiou et al., 2017, Section 2.3.2). Thus, asymptotically, a large ESS leads to small bias and MSE for bounded functions  $f$ . Additionally, the above derivations allow us to interpret the second moment

$$\rho = \mathbf{G}[(NW_1)^2] = (\mathbf{G}[NW_1])^2 + \text{Var}(NW_1) = 1 + \text{Var}(NW_1) \approx \frac{\text{Var}(\hat{\zeta}_{\text{IS}})}{\text{Var}(\hat{\zeta}_{\text{MC}})}$$

as the asymptotic relative efficiency of the two estimators, as long as this approximation is valid. In practice, a small ESS can be an indicator that importance sampling with  $\mathbf{G}$  may be inadequate. Note that relying solely on the empirical ESS may lead to problems, see the following example. To prepare, we prove a lemma regarding  $\rho$  for Gaussian targets and proposals.

**Lemma 3.6.** *Let  $\mathbf{P} = \mathcal{N}(\mu, \Sigma)$  and  $\mathbf{G} = \mathcal{N}(\nu, \Omega)$  be two  $p$ -dimensional Gaussian distributions with means  $\mu, \nu \in \mathbf{R}^p$  and SPD covariance matrices  $\Sigma, \Omega \in \mathbf{R}^{p \times p}$ . Then  $\rho$  is finite if, and only if,  $\Omega \succ \frac{1}{2}\Sigma$ .*

*Proof.* For the weights  $w = \frac{p}{g}$  we have

$$\begin{aligned}\rho = \mathbf{G}[w^2] &= \int \frac{p^2(x)}{g^2(x)} g(x) dx = \int \frac{p^2(x)}{g(x)} dx \\ &= \int \frac{\sqrt{\det \bar{\Omega}}}{\sqrt{(2\pi)^p \det \Sigma}} \exp \left( -(x - \mu)^T \Sigma^{-1} (x - \mu) + \frac{1}{2} (x - \nu)^T \Omega^{-1} (x - \nu) \right) dx.\end{aligned}$$

The exponent is a quadratic form in  $x$ , and so the integral is finite if, and only if, the matrix of coefficients,  $-\Sigma^{-1} + \frac{1}{2}\Omega^{-1}$  is negative definite. Rearranging terms, we see that this is equivalent to  $\Omega \succ \frac{1}{2}\Sigma$ .  $\square$

**Example 3.2** (failure of the ESS). Consider the Gaussian scale mixture

$$\mathbf{P} = \frac{1}{2} (\mathcal{N}(0, 1) + \mathcal{N}(0, \varepsilon^{-2}))$$

and proposal  $\mathbf{G} = \mathcal{N}(0, 1)$ . The weights are then given by

$$w(x) = \frac{1}{2} \left( 1 + \frac{\varepsilon}{\sqrt{2\pi}} \exp \left( -\frac{x^2}{2} (\varepsilon^2 - 1) \right) \right)$$

and their second moment w.r.t.  $\mathbf{G}$

$$\rho = \int w^2(x) \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{x^2}{2} \right) dx$$

is finite if, and only if,  $\varepsilon^2 > \frac{1}{2}$ , by the preceding lemma. Thus, for  $\varepsilon^2 \leq \frac{1}{2}$  interpreting the ESS or EF is not sensible. Nevertheless, given samples  $X^1, \dots, X^N \stackrel{\text{i.i.d.}}{\sim} \mathbf{G}$ , we may calculate the ESS in the usual way. If  $N$  is only moderately large, there is a high probability that most samples do not lie in a region where weights are small, i.e. in the tails of the second component. Thus, unless  $N$  is large, the empirical ESS will be large, deceiving us to think that importance sampling with  $\mathbf{G}$  is feasible.

We illustrate this by a simulation study, where we calculate the EF  $M = 100$  times for different values of  $N$  and  $\varepsilon$ . We used  $N = 100, 1000, 10000$  and  $\varepsilon^2 = 0.01, 0.1, 0.5$ ; the results may be found in Figure 3.2. Notice that for all values of  $\varepsilon$  considered, we have  $\rho = \infty$ . We see that even for  $N = 1000$  and  $\varepsilon = \frac{1}{2}$  the upper quartile of EFs is 71%, which seems reasonable to declare importance sampling to perform well.

Let us note that having access to the normalized weights  $w$  here allows us to spot this deficiency of the ESS by recognizing that while ESS is high, the weights  $w$  are not close to 1, but rather  $\frac{1}{2}$ .

As an alternative, we may want to assess whether importance sampling has converged through the empirical variance of  $\hat{\zeta}_N$ ,<sup>2</sup> i.e.,

$$\widehat{\text{Var}}(\hat{\zeta}_N) = \frac{1}{N} \left( \frac{1}{N} \sum_{i=1}^N w_i^2 f(X^i)^2 - \hat{\zeta}_N^2 \right)$$

is, while seemingly natural, flawed (Chatterjee and Diaconis, 2018). Indeed, the authors show that for any given threshold  $\epsilon$  we may find an  $N$  which only depends on  $\epsilon$ , such that the probability that the empirical variance exceeds  $\epsilon$  for this  $N$  is small. This is summarized in the following theorem.

**Theorem 3.3** (Chatterjee and Diaconis, 2018, Theorem 2.1). *Given any  $\epsilon > 0$ , there exists*

*lower bound on  $N$ ?*

$N \leq \epsilon^{-2} 2^{1+\epsilon^{-3}}$  such that the following is true. Take any  $\mathbf{G}$  and  $\mathbf{P}$  as in Theorem 3.2, and any  $f : \mathcal{X} \rightarrow \mathbf{R}$  such that  $\|f\|_{L^2(\mathbf{P})} \leq 1$ . Then

$$\mathbb{P} \left( \widehat{\text{Var}}(\hat{\zeta}_N) < \epsilon \right) \geq 1 - 4\epsilon.$$

<sup>2</sup>As the following arguments depend on the sample size  $N$ , we mark this dependency by adding  $N$  to the subscript of the estimator.

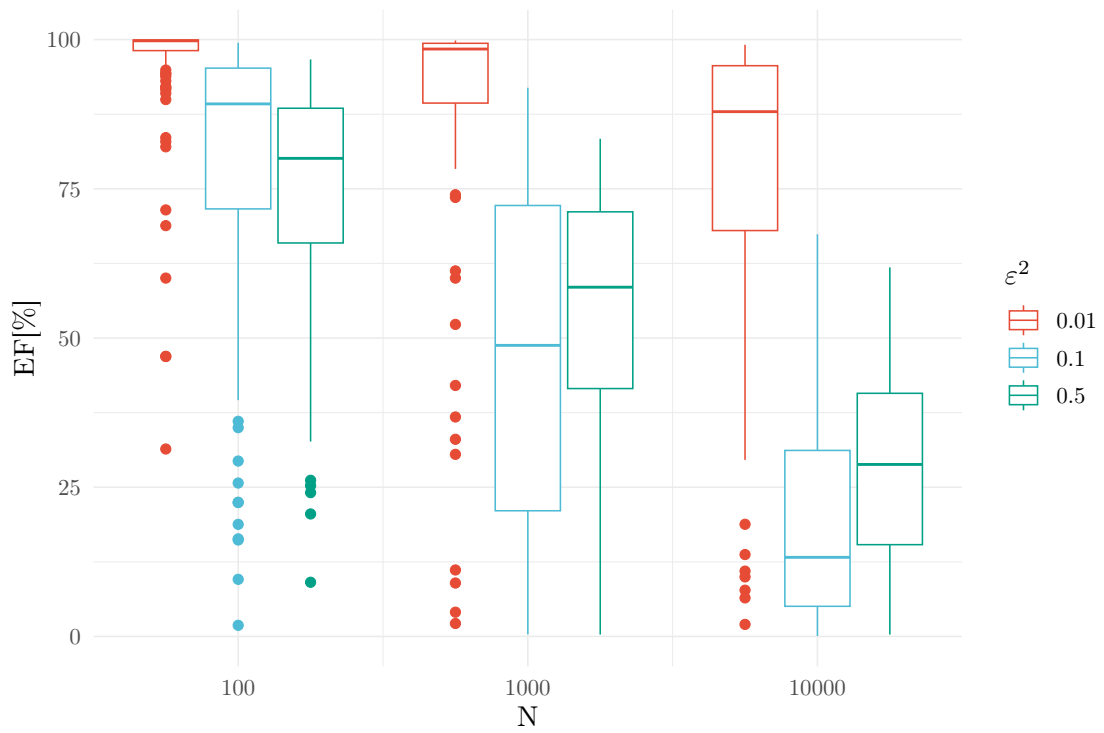


Figure 3.2: Empirical EF for the setup of Example 3.2 for varying sample sizes  $N$  and  $\varepsilon^2$  and  $M = 100$  replications. Here  $\mathbf{G} = \mathcal{N}(0, 1)$  and  $\mathbf{P} = \frac{1}{2} (\mathcal{N}(0, 1) + \mathcal{N}(0, \varepsilon^{-2}))$ . In all scenarios the second moment  $\rho$  is infinite, thus high EFs are misleading us to believe that importance sampling performs well when it does not.

The problem here is that  $N$  does not depend on  $\mathbf{G}$  and  $\mathbf{P}$ , so we may choose  $\mathbf{G}$  almost singular to  $\mathbf{P}$ . As an example, take  $\mathbf{P} = \mathcal{N}(0, 1)$  and  $\mathbf{G} = \mathcal{N}(0, \sigma^2)$  for  $\sigma^2 > \frac{1}{2}$ . The weights are then given by

$$w(x) = \sigma \exp\left(-\frac{x^2}{2}\left(1 - \frac{1}{\sigma^2}\right)\right),$$

and for  $X \sim \mathbf{G}$  the variance of  $w(X)X$  is

$$\tau^2 = \text{Var}(w(X)X) = \frac{\sigma^4}{(2\sigma^2 - 1)^{\frac{3}{2}}} \quad (3.10)$$

which goes to  $\infty$  as  $\sigma^2$  does, see the appendix for details. Thus, for a pre-specified  $\epsilon > 0$ , let  $N$  be as in Theorem 3.3 and choose  $\sigma^2$  such that  $\text{Var}(\hat{\zeta}_N) = \frac{\tau^2}{N}$  is larger than, say,  $10\epsilon$ . By the preceding theorem, we would, with large probability, observe a small empirical variance and thus declare  $\hat{\zeta}_N$  to have converged, whereas, in reality, we would need a sample size that is 100 times as large.

Thus using the empirical variance as a threshold for convergence should be avoided, at least for importance sampling where the weights can be evaluated exactly. For self-normalized importance sampling, the authors do not provide such a theorem. As a remedy (Chatterjee and Diaconis, 2018) suggest the heuristic  $q_N = \mathbb{E}Q_N$  where

$$Q_N = \max_{1 \leq i \leq N} W_i \in [0, 1].$$

This judges whether importance sampling has collapsed to just a few particles and is itself amenable to Monte-Carlo integration, by repeatedly sampling  $N$  samples from  $\mathbf{G}$  and calculating the weights. As this requires multiple runs of importance sampling, it may, however, be prohibitively expensive in practice.

In the following sections, we will predominantly take the position that we are interested in finding a good particle approximation  $\hat{\mathbf{P}}_N$  of the form Equation (3.7) over finding the optimal proposal  $\mathbf{G}^*$  from Proposition 3.1 and assume that the importance sampling weights can only be evaluated up to a constant. This has several reasons: First of all, for most problems considered in this thesis  $\mathbf{P}$  is usually a conditional distribution, e.g.  $\mathbf{P} = \mathbb{P}^{X|Y=y}$  for states  $X$  and observations  $Y$  in the SSM context. Should the appropriate densities exist, evaluating the weights amounts to calculating

$$\frac{d\mathbb{P}^{X|Y=y}}{d\mathbf{G}}(x) = \frac{p(x|y)}{g(x)} = \frac{p(y|x)p(x)}{g(x)p(y)} \propto \frac{p(y|x)p(x)}{g(x)}.$$

In these situations  $p(y) = \int p(x, y) dx$  is usually intractable. For  $\mathbf{G}^*$  we are in the same situation, where the evaluation of the integration constant  $\mathbf{P}|f|$  is infeasible, but the density  $|f(x)|p(x)$  is available. Second, focusing on the particle approximation allows us to consider multiple test functions  $f$ , e.g. focus on different marginals of  $\mathbf{P}$ , which is usually what practitioners are interested in. Finally, this allows us to simplify the notation used in this thesis.  $\mathbf{P}$  will always be the probability measure of interest and  $\mathbf{G}$  the proposal. In later parts of this thesis, we will predominantly perform Gaussian importance sampling, i.e.  $\mathbf{G} = \mathcal{N}(\mu, \Sigma)$ , hence a handy mnemonic is to think of  $\mathbf{G}$  as a Gaussian proposal.

Let us now turn towards the problem of finding a good proposal  $\mathbf{G}$  for a given  $\mathbf{P}$ .

### 3.3.1 Laplace approximation (LA)

The Laplace approximation (LA) goes back to Laplace (Laplace, 1986) who invented the technique to approximate moments of otherwise intractable distributions. Since (Tierney and Kadane, 1986; Tierney, Kass, and Kadane, 1989) rediscovered its use to approximate posterior means and variances, it has been a staple method for approximate inference. The method is based on a second-order Taylor series expansion of the log target density  $\log p(x)$  around its mode  $\hat{x}$ , i.e. matching mode and curvature. Assuming the density is sufficiently smooth, we have

$$\log p(x) \approx \log p(\hat{x}) + \underbrace{\nabla_x \log p(\hat{x})}_{=0} (x - \hat{x}) + \frac{1}{2} (x - \hat{x})^T H(x - \hat{x}) \quad (3.11)$$

where  $H$  is the Hessian of  $\log p$  evaluated at  $\hat{x}$ . As  $\log p(\hat{x})$  does not depend on  $x$ , the right-hand side can be seen (up to additive constants) as the density of a Gaussian distribution with mean  $\hat{x}$  and covariance matrix  $\Sigma = -H^{-1}$ . Thus using  $\mathbf{G} = \mathcal{N}(\hat{x}, -H^{-1})$  as a proposal in importance sampling seems promising. If  $\hat{x}$  is the unique global mode of  $p$  and  $H$  is negative definite, the LA yields an actual Gaussian distribution. To obtain the LA in practice, a Newton-Raphson scheme may be used, which conveniently tracks  $H$  as well. Furthermore, if  $\mathbf{P}$  includes more structure, e.g. it is the smoothing density in the SSM context, we may be able to exploit this structure to design efficient Newton-Raphson schemes, see Section 3.5.1.

The main advantage of the LA is that it is usually fast to obtain and, for sufficiently well-behaved distributions on a moderate dimensional space, provides reasonably high ESS. Additionally, the Newton-Raphson iterations to find the mode and Hessian are robust and require no simulation, unlike the other methods discussed further below. For the SSMs we consider in this thesis, the numerical methods can be implemented using the Kalman filter and smoother (Durbin and Koopman, 1997; Shephard and Pitt, 1997), even in the degenerate case where  $H$  is indefinite (Jungbacker and Koopman, 2007), see also Section 3.5.1.

However, as the LA is a local approximation, it may be an inappropriate description of the global behavior of the target, see Example 3.3 for a breakdown of LA, and the simulation studies presented in Section 3.7. Additionally, even if the LA works in principle, its ESS will usually degenerate quickly once the dimension increases whereas the Cross-Entropy method (CE-method) and Efficient Importance Sampling (EIS) do so at a slower pace.

### 3.3.2 The Cross-Entropy method (CE-method)

Recall from our discussion surrounding Theorem 3.2 that a small KL-divergence between the target  $\mathbf{P}$  and the proposal  $\mathbf{G}$  implies good performance for importance sampling. As the KL-divergence depends on global properties of  $\mathbf{P}$ , i.e. the Radon-Nikodym derivative  $\frac{d\mathbf{P}}{d\mathbf{G}}$ , minimizing it leads to a global approximation of  $\mathbf{P}$ , improving on the local-approximation provided by the LA.

The Cross-Entropy method (CE-method) (Rubinstein, 1999; Rubinstein and Kroese, 2004) implements this idea and selects from a parametric family  $(\mathbf{G}_\psi)_{\psi \in \Psi}$  of proposals the one that minimizes the Kullback Leibler divergence (KL-divergence) to the target. Here  $\Psi$  is usually a subset of  $\mathbf{R}^k$ , which may be open, closed or neither. Thus, the CE-method aims at solving the following optimization problem

$$\min_{\psi \in \Psi} \mathcal{D}_{\text{KL}}(\mathbf{P} \parallel \mathbf{G}_\psi),$$

for the optimal  $\psi_{\text{CE}}$ , should the minimum exist. The existence and uniqueness of  $\psi_{\text{CE}}$  will depend heavily on the choice of parametric family  $(\mathbf{G}_\psi)_{\psi \in \Psi}$  and  $\mathbf{P}$ .

We will assume the existence of a common dominating measure  $\mu$  for both  $\mathbf{P}$  and all  $\mathbf{G}_\psi$ ,  $\psi \in \Psi$  with corresponding densities  $p$  and  $g_\psi$ ,  $\psi \in \Psi$ . The importance sampling weights are then given by

$$w_\psi(x) = \frac{p(x)}{g_\psi(x)},$$

$x \in \mathcal{X}$ , or, if at least one of  $p$  and  $g_\psi$  is only available up to a constant, by

$$\tilde{w}_\psi(x) \propto \frac{p(x)}{g_\psi(x)}.$$

If the dependence on  $\psi$  is not of interest or the particular  $\psi$  is obvious from the context, we may drop the subscript.

The KL-divergence is given by

$$\mathcal{D}_{\text{KL}}(\mathbf{P} \parallel \mathbf{G}_\psi) = \mathbf{P}[\log w_\psi],$$

and can be infinite, e.g. if  $\mathbf{P}$  does not possess second moments and  $\mathbf{G}_\psi$  are Gaussian distributions. If the KL-divergence is infinite for all  $\psi \in \Psi$ , the CE-method becomes uninteresting. As such we will require that the KL-divergence is finite for at least one  $\psi \in \Psi$ , and restrict  $\Psi$ , without loss of generality, to those  $\psi$  where the KL-divergence is finite.

As the densities w.r.t. a common dominating measure exist, we may reformulate the optimization problem to maximize the cross-entropy between  $p$  and  $g_\psi$  instead:

$$\begin{aligned} \operatorname{argmin}_{\psi \in \Psi} \mathcal{D}_{\text{KL}}(\mathbf{P} \parallel \mathbf{G}_\psi) &= \operatorname{argmin}_{\psi \in \Psi} \mathbf{P}[\log p] - \mathbf{P}[\log g_\psi] \\ &= \operatorname{argmax}_{\psi \in \Psi} \mathbf{P}[\log g_\psi]. \end{aligned} \quad (3.12)$$

As the KL-divergence is non-negative by the information inequality, the cross-entropy  $\mathbf{P}[\log g_\psi]$  is bounded from above by the differential entropy of  $\mathbf{P}$ ,  $\mathbf{P}[\log p]$ . For centered distributions with covariance matrix  $\Sigma$  the differential entropy is bounded above by the maximum entropy distribution in this setting, the Gaussian  $\mathcal{N}(0, \Sigma)$  (Cover and Thomas, 2006, Example 12.2.8). Thus, if second moments of  $\mathbf{P}$  exist, the cross-entropy is bounded from above, and so a maximizer exists if the supremum over  $\Psi$  is attained. This would be the case if  $\Psi$  is compact and  $\psi \mapsto \mathbf{P}[\log g_\psi]$  is continuous, however compact  $\Psi$  is too restrictive for our purposes. Instead, we are going to focus on more realistic assumptions.

Suppose now that  $\psi \mapsto \log g_\psi(x)$  is (strictly) concave for  $\mathbf{P}$ -almost every  $x \in \mathcal{X}$  and  $\Psi$  is a convex subset of  $\mathbf{R}^k$ . Then  $\psi \mapsto \mathbf{P}[\log g_\psi]$  is (strictly) concave as well. As a consequence, we may apply the usual results from convex optimization, i.e. every local maximum is a global one and if  $\psi \mapsto \log g_\psi(x)$  is strictly concave for  $\mathbf{P}$ -almost every  $x$ , there is at most one maximizer (Bazaraa, Sherali, and Shetty, 2006, Theorem 3.4.2).

As we have seen in Lemma 3.3, the densities of exponential families are log-concave in the natural parameter, and as such they will be the primary candidates for our investigations of the CE-method. If we use proposals from an exponential family, we may get rid of the base measure term  $h(x)$  in the densities, as the following lemma shows.

**Lemma 3.7.** *Let  $\mathbf{P}$  be a probability measure on  $\mathcal{X} = \mathbf{R}^p$  and let  $(\mathbf{G}_\psi)_{\psi \in \Psi}$  be a natural exponential family on  $\mathcal{X}$  such that  $\mathbf{P} \ll \mathbf{G}_\psi$  for all  $\psi \in \Psi$ . Let  $\mu$  be the dominating measure of the exponential family, such that*

$$\frac{d\mathbf{G}_\psi}{d\mu}(x) = \frac{h(x)}{Z(\psi)} \exp(\psi^T T(x)),$$

with  $h \geq 0$   $\mu$ -a.s.

Then  $h\mu$  is a dominating measure for both  $\mathbf{P}$  and  $\mathbf{G}_\psi$  for every  $\psi$  in  $\Psi$ .

*Proof.* Let  $A \subseteq \mathbf{R}^p$  be measurable. As  $h$  is a.s. non-negative,  $(h\mu)(A) = 0$  implies that  $h\mathbf{1}_A = 0$   $\mu$ -a.s. Thus  $\mathbf{G}_\psi(A) = \int \mathbf{1}_A(x) \frac{h(x)}{Z(\psi)} \exp(\psi^T T(x)) d\mu = 0$  for all  $\psi$  as well. As  $\mathbf{G}_\psi \gg \mathbf{P}$  and  $\gg$  is transitive,  $h\mu$  dominates  $\mathbf{P}$  as well.  $\square$

As a consequence, when performing importance sampling with target  $\mathbf{P}$  and proposal  $\mathbf{G}_\psi$  from an exponential family, we will assume in the following that  $h \equiv 1$ , achieved by taking  $h\mu$  as the joint dominating measure.

An additional attractive property of the CE-method for exponential families with natural parameter  $\psi \in \mathbf{R}^k$  is that the optimal  $\psi_{\text{CE}}$  only depends on the expected value  $\mathbf{P}[T]$ . We first show, that if the covariance of the sufficient statistic is positive definite, the expected value of  $T$  under  $\mathbf{G}_\psi$  uniquely determines  $\psi \in \Psi$ , see also (Brown, 1986, Corollary 2.5) for a similar result in minimal exponential families.

**Lemma 3.8.** *Let  $(\mathbf{G}_\psi)_{\psi \in \Psi}$  form a  $k$ -dimensional natural exponential family with log-densities*

$$\log g_\psi(x) = \psi^T T(x) - \log Z(\psi),$$

and convex parameter space  $\Psi \subseteq \mathbf{R}^k$ . Let  $\psi, \psi' \in \text{int } \Psi$  with  $\mathbf{G}_\psi[T] = \mathbf{G}_{\psi'}[T]$ . If  $\text{Cov}_{\mathbf{G}_\psi} T$  is positive definite, then  $\psi$  and  $\psi'$  coincide.

*Proof.* Consider the function  $b : \Psi \rightarrow [-\infty, \infty)$

$$\xi \mapsto b(\xi) = \mathbf{G}_{\psi'}[\log g_\xi] = \xi^T \mathbf{G}_{\psi'}[T] - \log Z(\xi).$$

By Theorem 3.1,  $\mathbf{G}_{\psi'}[T]$  is finite and  $b$  possesses derivatives of every order. Then  $\psi$  is a critical point of this map, as the gradient at  $\psi$  is

$$\mathbf{G}_{\psi'}[T] - \nabla_{\psi} \log Z(\psi) = \mathbf{G}_{\psi'}[T] - \mathbf{G}_{\psi}[T] = 0.$$

The Hessian of this function at  $\xi$  is, see Theorem 3.1,

$$-H_{\xi} \log Z(\xi) = -\text{Cov}_{\mathbf{G}_{\xi}}[T],$$

which is negative semi-definite, so  $b$  is concave. At  $\xi = \psi$  it is negative definite, so the critical point  $\psi$  is a strict local maximum. By concavity, it is the unique global maximum, and thus the unique critical point, so  $\psi = \psi'$ .  $\square$

**Proposition 3.2** (The CE-method for exponential families). *Let  $(\mathbf{G}_{\psi})_{\psi \in \Psi}$  form a  $k$ -dimensional natural exponential family with log-densities*

$$\log g_{\psi}(x) = \psi^T T(x) - \log Z(\psi),$$

*and convex parameter space  $\Psi \subseteq \mathbf{R}^k$ . Suppose  $T \in L^1(\mathbf{P})$ .*

*If there is a  $\psi_{CE} \in \Psi$  such that*

$$\mathbf{P}[T] = \mathbf{G}_{\psi_{CE}}[T],$$

*then  $\psi_{CE}$  is a maximizer of Equation (3.12). Furthermore, if  $\text{Cov}_{\mathbf{G}_{\psi_{CE}}} T$  is positive definite the maximizer is unique.*

*Proof.* The target may be rewritten as

$$\psi \mapsto f(\psi) = \mathbf{P}[\log g_{\psi}(x)] = -\log Z(\psi) + \psi^T \mathbf{P}[T].$$

As  $\log Z(\psi)$  is the cumulant-generating function of  $\mathbf{G}_{\psi}$  it is twice differentiable, and so is  $f$ . The gradient of  $\log Z(\psi)$  is

$$\nabla_{\psi} \log Z(\psi) = \mathbf{G}_{\psi}[T]$$

and its Hessian is

$$H_{\psi} \log Z(\psi) = \text{Cov}_{\mathbf{G}_{\psi}}(T)$$

the covariance of  $T$  under  $\mathbf{G}_{\psi}$ . Thus the Hessian of  $f$  is

$$H_{\psi} f = -\text{Cov}_{\mathbf{G}_{\psi}}(T),$$

which is negative-semi-definite. Therefore  $f$  is concave, and any local maximizer  $\psi$  is a global maximizer. The gradient of  $f$  is

$$\nabla_{\psi} f(\psi) = \mathbf{P}[T] - \mathbf{G}_{\psi}[T],$$

which is equal to 0 if, and only if,  $\psi$  solves

$$\mathbf{P}[T] = \mathbf{G}_{\psi}[T].$$

Uniqueness follows from the preceding Lemma 3.8.  $\square$

As a consequence, the CE-method for natural exponential families reduces to matching the moments of the sufficient statistic of the target and proposal. In many cases, this system of equations can be solved analytically or by gradient descent algorithms. Let us discuss the assumptions and applicability of this proposition. Assuming that  $T \in L^1(\mathbf{P})$  is necessary for the target to be finite, it cannot be dropped. As  $T$  typically consists of polynomial, rational or exponential functions, this is not too restrictive, provided the target does not exhibit heavy tails. The proof of uniqueness relies on  $\text{Cov}_{\mathbf{G}_{\psi}} T$  being positive definite, to ensure that  $\psi \mapsto \log Z(\psi)$  is strictly convex. This could also be achieved by requiring the exponential family to be minimal, see (Brown, 1986, Theorem 1.13 (iv)). The existence of a  $\psi$  such that  $\mathbf{P}[T] = \mathbf{G}_{\psi}[T]$  is not restrictive for most commonly used



distributions: for the (multivariate) normal, Poisson, negative binomial and binomial distribution there is always a unique solution, as the sufficient statistics consist of means and covariances.

While  $\mathbf{P}[T]$  is usually not available, it is itself amenable to importance sampling. Given a proposal  $\mathbf{G}$  we may estimate  $\mathbf{P}[T]$  by  $\hat{\mathbf{P}}_N T = \sum_{i=1}^N W^i T(X^i)$  for  $X^1, \dots, X^N \stackrel{\text{i.i.d.}}{\sim} \mathbf{G}$  and auto-normalized importance sampling weights  $W^i$  and in turn, applying Proposition 3.2, estimate  $\psi_{\text{CE}}$  by  $\hat{\psi}_{\text{CE}}$  solving

$$\hat{\mathbf{P}}_N[T] = \mathbf{G}_{\hat{\psi}_{\text{CE}}}[T]. \quad (3.13)$$

As  $T \in L^1(\hat{\mathbf{P}}_N)$ , the only conditions we have to check to apply the above proposition are that this equation has a unique solution  $\mathbf{G}$ -almost surely in the interior of  $\Psi$  and that  $\Psi$  is convex.

To apply the CE-method in practice, one usually iterates the sampling and estimation steps, using the previously found  $\hat{\psi}_{\text{CE}}$  to sample in the current iteration and starting the iteration with a proposal from the same exponential family  $\mathbf{G} = \mathbf{G}_{\psi^0}$ . To ensure numerical convergence, a popular device is that of common random numbers (CRNs), i.e. using the same random number seed in all iterations. A basic version of the CE-method is presented in Algorithm 4.

---

**Algorithm 4** The basic CE-method algorithm for exponential families

---

**Require:** exponential family  $(\mathbf{G}_{\psi})_{\psi \in \Psi}$ , initial  $\psi^0$ , sample size  $N$ , unnormalized weights  $\tilde{w}$

- 1: set  $l = 0$
- 2: store random number seed
- 3: **repeat**
- 4:   restore random number seed
- 5:   sample  $X^1, \dots, X^N \sim \mathbf{G}_{\psi^l}$
- 6:   calculate self-normalized weights  $W^i$  for  $i = 1, \dots, N$
- 7:   estimate  $\hat{\psi}_{\text{CE}}$  ▷ Equation (3.13)
- 8:   set  $\psi^{l+1} = \hat{\psi}_{\text{CE}}$
- 9:   set  $l = l + 1$
- 10: **until**  $\hat{\psi}^l$  converged
- 11: **return**  $\hat{\psi}_{\text{CE}} = \hat{\psi}^l$

---

literature review CEM

The CE-method is routinely used for estimating failure probabilities for rare events (Homem-de-Mello, 2007) and has been applied to Bayesian posterior inference (Ehre et al., 2023; Engel et al., 2023), Bayesian marginal likelihood estimation (Chan and Eisenstat, 2012) and optimal control problems (Kappen and Ruiz, 2016; Zhang et al., 2014).

more lit. review CEM

Importance sampling is well known to exhibit the curse of dimensionality (COD) (Bengtsson, Bickel, and B. Li, 2008), i.e. the phenomenon that in many problems, unless  $N$  grows exponentially with the dimension of  $\mathcal{X}$ , the weights collapse to a single particle, i.e.  $W^{(N)} \rightarrow 1$  as the dimension of  $\mathcal{X}$  goes to  $\infty$ . As the CE-method employs importance sampling to obtain  $\hat{\psi}_{\text{CE}}$ , it too is affected by this phenomenon, see also Section 3.7. The screening method (Rubinstein and Glynn, 2009) deals with the COD by keeping components of  $\psi^l$  that vary too much from iteration to iteration fixed, in essence reducing the dimension of  $\Psi$ . Alternatively, the improved cross-entropy method (Chan and Kroese, 2012) suggests generating approximately independent samples from  $\mathbf{P}$  by, e.g., MCMC-methods, and replacing the importance sampling version of  $\hat{\mathbf{P}}_N$  in Equation (3.13) by the actual empirical distribution. Still, in high dimensions both of these approaches may be difficult to implement: the screening method may not move far from the initial proposal and MCMC-methods are expensive in high dimensions.

As stated in (Chan and Kroese, 2012) there may be two reasons as to why the CE-method fails: either the parametric family is not rich enough to give a good approximation to  $\mathbf{P}$ , i.e.  $\mathcal{D}_{\text{KL}}(\mathbf{P} \parallel \mathbf{G}_{\psi_{\text{CE}}})$  is still large, or the estimate  $\hat{\psi}_{\text{CE}}$  fails to be close to  $\psi_{\text{CE}}$ . As our simulation studies Section 3.7

suggest, the reason for the degeneracy seems to be the latter. It will thus be beneficial to investigate the asymptotic behavior of  $\hat{\psi}_{\text{CE}}$ .

In the remainder of this section, we will derive novel results on the performance of the estimator  $\hat{\psi}_{\text{CE}}$  of  $\psi_{\text{CE}}$ . In particular, we will investigate under which conditions  $\hat{\psi}_{\text{CE}}$  is consistent and asymptotically normal. To focus on the asymptotic behavior, we will only perform a single iteration of the basic CE-method algorithm (Algorithm 4). While we restrict ourselves here to the setting of  $k$ -dimensional natural exponential families, these results should generalize to other classes of distributions as well. The advantage that this class of families has is that due to the structure of the densities, they provide straightforward (regularity) conditions for the asymptotic results to hold. As the target functions are concave, these conditions are rather liberal. We start with proving the consistency of  $\hat{\psi}_{\text{CE}}$ .

**Theorem 3.4** (consistency of  $\hat{\psi}_{\text{CE}}$ ). *Adopt the same assumptions as in Proposition 3.2. Furthermore, let  $\mathbf{G} \gg \mathbf{P}$  be a proposal distribution and assume that*

- (i)  $\psi_{\text{CE}}$  is the unique maximizer of Equation (3.12),
- (ii)  $\psi_{\text{CE}}$  is in the interior of the convex parameter space  $\Psi$ .

Then  $\hat{\psi}_{\text{CE}}$  is a strongly consistent estimator of  $\psi_{\text{CE}}$ .

The proof is based on the following theorem of Haberman.

**Theorem 3.5** ((Haberman, 1989, Theorem 5.1)<sup>3</sup>). *Let  $\Psi \subseteq \mathbf{R}^k$ ,  $\mathcal{X}$  a separable, complete metric space and  $b_{\mathcal{X}} : \mathcal{X} \times \mathbf{R}^k \rightarrow [-\infty, \infty)$  such that for every  $x \in \mathcal{X}$  the function*

$$b(x, \cdot) : \mathbf{R}^k \rightarrow [-\infty, \infty), \psi \mapsto b(x, \psi)$$

*is concave. Let  $\mathbf{P}$  be a probability measure on  $\mathcal{X}$  such that  $\mathbf{P}[b(\cdot, \psi)] < \infty$  for all  $\psi \in \mathbf{R}^k$ . Assume that  $\psi^* \in \Psi$  is the unique maximizer of*

$$b_{\Psi} : \Psi \rightarrow [-\infty, \infty), \psi \mapsto \mathbf{P}[b(\cdot, \psi)].$$

*Let  $(X^i)_{i \in \mathbf{N}} \stackrel{i.i.d.}{\sim} \mathbf{P}$  be a sequence of i.i.d. random variables with distribution  $\mathbf{P}$  and let for  $N \in \mathbf{N}$  let*

$$\hat{\mathbf{P}}_N = \frac{1}{N} \sum_{i=1}^N \delta_{X^i}$$

*be their empirical distribution. Let  $(\hat{\psi}_N)_{N \in \mathbf{N}}$  be a sequence of  $M$ -estimators, i.e. a sequence of maximizers of*

$$\hat{b}_{\Psi} : \Psi \rightarrow [-\infty, \infty), \psi \mapsto \hat{\mathbf{P}}_N[b(\cdot, \psi)].$$

*Assume that the following conditions hold:*

- (C1) *For some closed set  $V$ ,  $\psi^*$  is in the interior of  $V$  and  $\Psi \cap V$  is closed.*
- (C2)  *$\psi^*$  is the unique maximizer of*

$$b_{\text{cl}(\Psi)} : \text{cl}(\Psi) \rightarrow [-\infty, \infty), \psi \mapsto \mathbf{P}[b(\cdot, \psi)],$$

*where  $\text{cl}$  denotes the closure of  $\Psi$  in  $\mathbf{R}^k$ .*

- (C3)  *$\Psi$  is convex and  $b_{\Psi}$  is finite on a nonempty open set.*

*Then*

$$\hat{\psi}_N \xrightarrow{N \rightarrow \infty} \psi^*$$

*$\mathbf{P}$ -almost surely, so  $\hat{\psi}_N$  is strongly consistent.*

---

<sup>3</sup>Note that while the actual theorem assumes conditions 1,2,5 and 6 in the paper, C3 as stated here implies conditions 5 and 6, see also the discussion in Sections 2.3 and 2.4 in (Haberman, 1989).

The assumptions of this theorem ensure that the unique optimum is in the interior of  $\Psi$  and „well-separated“ from its boundary, so there are no additional maximizers on the boundary. In this case, concavity of  $b(x, \psi)$  together with the law of large numbers yield uniform convergence of  $\hat{\mathbf{P}}_N[b(\cdot, \psi)] \rightarrow \mathbf{P}[b(\cdot, \psi)]$  on compacta and thus also for  $\hat{\psi}_N$ , see (Haberman, 1989, pp. 1652).

To apply this theorem to our setting, let us begin by extending it to incorporate importance sampling.

**Proposition 3.3.** *Assume that the conditions of Theorem 3.5 are fulfilled and let  $\mathbf{G} \gg \mathbf{P}$  be another probability measure with Radon-Nikodym derivative  $w(x) = \frac{d\mathbf{P}}{d\mathbf{G}}(x)$ . Let  $(X^i)_{i \in \mathbf{N}} \stackrel{i.i.d.}{\sim} \mathbf{G}$  and consider the particle approximations*

$$\begin{aligned}\tilde{\mathbf{P}}_N &= \frac{1}{N} \sum_{i=1}^N w(X^i) \delta_{X^i}, \\ \hat{\mathbf{P}}_N &= \sum_{i=1}^N W^i \delta_{X^i},\end{aligned}$$

and suppose for every  $N \in \mathbf{N}$  there exist M-estimators

$$\begin{aligned}\tilde{\psi}_N &\in \operatorname{argmax}_{\psi \in \Psi} \tilde{\mathbf{P}}_N[b(\cdot, \psi)], \\ \hat{\psi}_N &\in \operatorname{argmax}_{\psi \in \Psi} \hat{\mathbf{P}}_N[b(\cdot, \psi)].\end{aligned}$$

Then both  $\tilde{\psi}_N$  and  $\hat{\psi}_N$  are strongly consistent estimators of  $\psi^*$ .

*Proof.* Define a new objective function  $\tilde{b} : \mathcal{X} \times \mathbf{R}^k \rightarrow [-\infty, \infty)$  by

$$\tilde{b}(x, \psi) = w(x)b(x, \psi).$$

Then  $\mathbf{G}[\tilde{b}(\cdot, \psi)] = \mathbf{P}[b(\cdot, \psi)]$  for all  $\psi \in \Psi$ , and so  $\psi^*$  is the unique global maximum of

$$\psi \mapsto \mathbf{G}[\tilde{b}(\cdot, \psi)].$$

As  $\mathbf{G}[\tilde{b}(\cdot, \psi)] = \mathbf{P}[b(\cdot, \psi)] < \infty$  and for fixed  $x \in \mathcal{X}$   $\tilde{b}(x, \cdot) = w(x)b(x, \cdot)$  is concave, we may directly apply Theorem 3.5 to  $\tilde{\psi}_N$ , showing its strong consistency.

For  $\hat{\psi}_N$ , notice that for a fixed sample  $X^1, \dots, X^N \stackrel{i.i.d.}{\sim} \mathbf{G}$  and any function  $f : \mathcal{X} \rightarrow [-\infty, \infty)$  we have, a.s.,

$$\hat{\mathbf{P}}_N[f] = \sum_{i=1}^N W^i f(X^i) = \frac{\mathbf{G}[\tilde{w}]}{\sum_{i=1}^N \tilde{w}(X^i)} \sum_{i=1}^N \frac{\tilde{w}(X^i)}{\mathbf{G}[\tilde{w}]} f(X^i) = \frac{\mathbf{G}[\tilde{w}]}{\sum_{i=1}^N \tilde{w}(X^i)} \tilde{\mathbf{P}}_N[f] \propto \tilde{\mathbf{P}}_N[f],$$

where  $\tilde{w}$  are the unnormalized weights, i.e.  $\frac{\tilde{w}(x)}{\mathbf{G}[\tilde{w}]} = w(x)$ ,  $x \in \mathcal{X}$ . Thus  $\hat{\psi}_N$  maximizes  $\tilde{\mathbf{P}}_N[b(\cdot, \psi)]$  as well, and the result follows from the consistency of  $\tilde{\psi}_N$ .  $\square$

Let us now prove the promised consistency of the CE-method.

*Proof (Theorem 3.4).* We show that the assumptions of Theorem 3.5 are fulfilled. Let

$$b : \mathbf{R}^p \times \mathbf{R}^k \rightarrow [-\infty, \infty) \quad b(x, \psi) = \begin{cases} \log g_\psi(x) & \psi \in \Psi, \\ -\infty & \text{else.} \end{cases}$$

As  $\Psi$  is convex and  $g_\psi(x)$  is log-concave (see Lemma 3.3),  $b(x, \cdot)$  is concave. Let  $X^1, \dots, X^N \stackrel{i.i.d.}{\sim} \mathbf{P}$  and let  $\hat{\mathbf{P}}_N = \frac{1}{N} \sum_{i=1}^N \delta_{X^i}$ . For  $\psi \in \Psi$  we have

$$\mathbf{P}[b(\cdot, \psi)] = \psi^T \mathbf{P}[T] - \log Z(\psi) < \infty,$$

as  $T \in L^1(\mathbf{P})$ , while for  $\psi \notin \Psi$  this integral is  $-\infty$ . Thus we only have to check that (C1)-(C3) are fulfilled.

For condition (C1) note that, as  $\psi_{CE}$  is in the interior of  $\Psi$ , we may choose  $\varepsilon > 0$  such that the closed  $\varepsilon$  ball around  $\psi_{CE}$ ,  $\bar{B}_\varepsilon(\psi_{CE})$  is completely contained in  $\Psi$ , so letting  $V = \bar{B}_\varepsilon(\psi_{CE})$  implies the condition. Condition (C2) is fulfilled by the definition of  $b$  and condition (C3) is fulfilled by considering the neighborhood of  $\psi_{CE}$  that is assumed to be contained in  $\Psi$ . Finally, by Proposition 3.3,  $\hat{\psi}_{CE}$  is strongly consistent.  $\square$

The assumptions on  $\psi_{CE}$  and  $\Psi$  in Theorem 3.4 could be somewhat looser, as the concavity of the target function is a rather strong property. In natural exponential families,

$$\Psi = \{\psi \in \mathbf{R}^k : Z(\psi) < \infty\}$$

is always convex so this is not a strong restriction. In regular exponential families,  $\Psi$  is open and so only the existence and uniqueness of  $\psi_{CE}$  are required. Uniqueness may be attained, e.g., by Lemma 3.8. It will also hold if the exponential family considered is minimal (Brown, 1986, Corollary 2.5). Existence is a matter of correctly specifying the exponential family. For example, in Section 3.5.2 we will exploit the Markov structure of targets to restrict ourselves to Gaussian Markov processes for  $(\mathbf{G}_\psi)_{\psi \in \Psi}$ .

Not only is  $\log g_\psi$  concave, but it also possesses derivatives of any order, at least on the interior of  $\Psi$ . Indeed, its Hessian is given by the inverse of the Fisher-information matrix  $I(\psi)^{-1}$ :

$$H_\psi \log g_\psi = -H_\psi \log Z(\psi) = -\text{Cov}_{\mathbf{G}_\psi}(T) = -I(\psi)^{-1}.$$

These rather strong properties enable us to derive a central limit theorem for the CE-method with natural exponential family proposals under quite liberal conditions.

**Theorem 3.6** (CLT for  $\hat{\psi}_{CE}$ ). *Adopt the same assumptions as in Proposition 3.2. Furthermore, let  $\mathbf{G} \gg \mathbf{P}$  be a proposal distribution with weights  $w = \frac{d\mathbf{P}}{d\mathbf{G}}$  and assume that*

- (i)  $\psi_{CE} \in \Psi$  is the unique maximizer of Equation (3.12) which lies in the interior of the convex parameter space  $\Psi$ ,
- (ii) the Fisher information matrix  $I(\psi_{CE})$  exists and is positive definite,
- (iii)  $w, wT \in L^2(\mathbf{G})$ , and
- (iv)  $T \in L^2(\mathbf{P})$ .

Then

$$\sqrt{N}(\hat{\psi}_{CE} - \psi_{CE}) \xrightarrow{\mathcal{D}} \mathcal{N}(0, BMB)$$

where  $B = I(\psi_{CE}) = \text{Cov}_{\mathbf{G}_{\psi_{CE}}}(T)^{-1}$  and

$$M = \text{Cov}_{\mathbf{G}}(wT) = \mathbf{G} [w^2(T - \mathbf{P}[T])(T - \mathbf{P}[T])^T] = \mathbf{P} [w(T - \mathbf{P}[T])(T - \mathbf{P}[T])^T].$$

To prove Theorem 3.6, let us start again with a general version of a central limit theorem for M-estimators based on concave objective functions.

**Theorem 3.7** ((Haberman, 1989, Theorem 6.1)<sup>4</sup>). *Consider the same setting as in Theorem 3.5.*

*Assume further that  $\psi^*$  lies in the interior of  $\Psi$  and that the following conditions hold:*

- (C7) *The Hessian  $H_\psi \mathbf{P}[b(\cdot, \psi^*)]$  exists and is non-singular.*
- (C10) *For  $X \sim \mathbf{P}$  and some neighborhood  $V$  of  $\psi^*$*

$$\sigma^2(\psi, \xi) = \mathbb{E}(b'(X, \psi, \xi))^2 < \infty \quad \psi \in V, \xi \in \mathbf{R}^k,$$

where  $b'(x, \psi, \xi) = \lim_{a \downarrow 0} a^{-1}(b(x, \psi + a\xi) - b(x, \psi))$  is the directional derivative. Note that if  $b$  is differentiable for all  $\psi \in V$ ,  $b'(x, \psi, \xi) = \xi^T \nabla_\psi b(x, \psi)$  and it suffices to assume  $(\nabla_\psi b(x, \psi))_i (\nabla_\psi b(x, \psi))_j \in L^1(\mathbf{P})$  for all  $\psi \in V$  and  $i, j = 1, \dots, k$ .

<sup>4</sup>Note, again, that the original theorem is based on conditions 7,8,9 in the paper. However, under (C7), condition (C10) implies conditions 8 and 9 in the paper. See the discussion in Section 3.1 in (Haberman, 1989).

Let  $M = \text{Cov}(\nabla_\psi b(X, \psi))$  and let  $B = -(H_\psi \mathbf{P}[b(\cdot, \psi)])^{-1}$ . Then

$$\sqrt{N}(\hat{\psi}_N - \psi) \xrightarrow{\mathcal{D}} \mathcal{N}(0, BMB). \quad (3.14)$$

Similar to the consistency result above (Proposition 3.3), we need to extend this CLT to account for importance sampling.

**Proposition 3.4.** *Assume that the conditions of Theorem 3.7 are fulfilled and use the same notation as in Proposition 3.3. Furthermore, assume that*

- (i)  $w(\cdot)b'(\cdot, \psi, \xi) \in L^2(\mathbf{G})$  in a neighborhood  $N$  of  $\psi^*$  for all  $\xi \in \mathbf{R}^k$ .

Then

$$\sqrt{N}(\tilde{\psi}_N - \psi^*) \xrightarrow{\mathcal{D}} \mathcal{N}(0, BMB), \quad (3.15)$$

where  $M = \text{Cov}(w(X)\nabla_\psi b(X, \psi^*))$  for  $X \sim \mathbf{G}$  and  $B = -(H_\psi \mathbf{P}[b(\cdot, \psi^*)])^{-1}$  is as in Theorem 3.7. Additionally

$$\sqrt{N}(\hat{\psi}_N - \psi^*) \xrightarrow{\mathcal{D}} \mathcal{N}(0, BMB). \quad (3.16)$$

*Proof.* Similar to the proof of Proposition 3.3, define the new objective function  $\tilde{b} : \mathcal{X} \times \mathbf{R}^k \rightarrow [-\infty, \infty)$  by

$$\tilde{b}(x, \psi) = w(x)b(x, \psi),$$

and notice that  $\mathbf{G}[\tilde{b}(\cdot, \psi)] = \mathbf{P}[b(\cdot, \psi)]$ . Let us verify the conditions of Theorem 3.7 for  $\tilde{b}$  and the probability measure  $\mathbf{G}$ .

For condition (C7), as  $H_\psi \mathbf{P}[b(\cdot, \psi)]$  exists and is non-singular, so does

$$H_\psi \mathbf{G}[\tilde{b}(\cdot, \psi)] = H_\psi \mathbf{P}[b(\cdot, \psi)]$$

exist and is non-singular. Similarly, it is easy to see that  $\tilde{b}'(x, \psi, \xi) = w(x)b'(x, \psi, \xi)$  and so for  $X \sim \mathbf{G}$

$$\sigma_b^2(\psi, \xi) = \mathbb{E} \left( \tilde{b}'(X, \psi, \xi) \right)^2 = \mathbb{E} w^2(X) b'(X, \psi, \xi)^2 < \infty$$

by assumption (i), showing condition (C10). Thus we may apply Theorem 3.7 to  $\tilde{b}$  and  $\mathbf{G}$ , finishing the proof.  $\square$

Interestingly, importance sampling only affects the  $M$  component of the asymptotic variance. The reason for this is that  $M$  is a quadratic function of the weights  $w$ , while  $B$  only depends linearly on  $w$ , allowing to switch integrators from  $\mathbf{G}$  to  $\mathbf{P}$ . We now have all the tools at our disposal to proof Theorem 3.6.

*Proof of Theorem 3.6.* We show that the assumptions and conditions of Theorem 3.7 for the objective function  $b : \mathcal{X} \times \mathbf{R}^k \rightarrow [-\infty, \infty)$

$$b(x, \psi) = \begin{cases} \log g_\psi(x) & x \in \Psi \\ -\infty & \text{else,} \end{cases}$$

are fulfilled, which, together with Proposition 3.4 will show the claim.

The Hessian of the objective function is, for  $\psi \in \text{int } \Psi$

$$H_\psi \mathbf{P}[b(\cdot, \psi)] = H_\psi \mathbf{P}[\psi^T T - \log Z(\psi)] = -H_\psi \log Z(\psi) = -I(\psi),$$

as the cumulant generating function is smooth on  $\text{int } \Psi$  (Theorem 3.1). Thus the Hessian is non-singular by assumption (ii), showing that condition (C7) is fulfilled.

For condition (C10), note that for  $\psi \in \text{int } \Psi$ ,  $b$  is differentiable with gradient

$$\nabla_\psi b(x, \psi) = T(x) - \nabla_\psi \log Z(\psi) = T(x) - \mathbf{G}_\psi[T].$$

By assumption (iv),  $\nabla_\psi b(x, \psi) \in L^2(\mathbf{P})$ , showing that condition (C10) is fulfilled.

To show that the central limit theorem applies to  $\hat{\psi}_{\text{CE}}$ , we additionally show that assumption (i) in Proposition 3.4 is fulfilled, which will finish the proof. To this end, note that

$$w(x)b'(x, \psi, \xi) = w(x)\xi^T \nabla_\psi b(x, \psi) = \xi^T (w(x)(T(x) - \mathbf{G}_\psi[T])) \in L^2(\mathbf{G})$$

by assumption (iii).

Finally, to show the representation of  $M$ , note that by Proposition 3.4 we have for  $X \sim \mathbf{G}$

$$M = \text{Cov}(w(X)(T(X) - G_{\psi_{\text{CE}}}[T])),$$

and  $\mathbb{E}w(X)(T(X) - \mathbf{G}_{\psi_{\text{CE}}}[T]) = 0$  as  $\mathbf{G}_{\psi_{\text{CE}}}[T] = \mathbf{P}[T]$ .  $\square$

The form of the asymptotic covariance matrix is that of the sandwich estimator (White, 1982), corrected for the importance sampling with  $\mathbf{G}$ . This is not surprising: the CE-method essentially performs maximum likelihood estimation of  $\psi$  where the data comes from the misspecified  $\mathbf{P}$ . Additionally, we have to correct the variance for performing importance sampling with  $\mathbf{G}$ , instead of sampling directly from  $\mathbf{P}$ .

The assumptions of Theorem 3.6 are minimal to facilitate the proof. The existence and positive definiteness of the Fisher information matrix are easily checked for the exponential family proposal and hold for minimal regular exponential families. Additionally, we have two moment constraints that involve the weights  $w$  and the sufficient statistic  $T$ . That  $wT \in L^2(\mathbf{G})$  may be seen as a generalization of the existence of the second moment  $\rho = \mathbf{G}[w^2]$ , adapted to the exponential family setting. As such it is a natural requirement. That  $T \in L^2(\mathbf{P})$  is required for the application of Theorem 3.7, and, as mentioned before, should not be problematic in practice, except for heavy-tailed distributions.

For our application, we will choose  $(\mathbf{G}_\psi)_{\psi \in \Psi}$  to consist of Gaussian distributions with natural parameter  $\psi = (\Sigma^{-1}\mu, -\frac{1}{2}\Sigma^{-1})$  and sufficient statistic  $T(x) = (x, xx^T)$ . Thus  $T \in L^2(\mathbf{P})$  is equivalent to  $\mathbf{P}$  having fourth order moments, which is reasonable if the target is not heavy-tailed.

If  $(\mathbf{G}_\psi)_{\psi \in \Psi}$  do not form an exponential family,  $\hat{\psi}_{\text{CE}}$  will still be consistent and asymptotically normal, provided the usual regularity conditions for M-estimators apply. These usually include conditions to ensure the maximum is well-separated and the target is sufficiently smooth such that a Taylor expansion around the maximum is feasible. To extend our results to more involved settings, we refer the reader to (Van der Vaart, 2000) for an empirical process treatment of M- and related Z-estimators, (Haberman, 1989) for asymptotics when the objective function is concave, but the maximum may lie on the border of the parameter space and (Liang and Zeger, 1995) for a review of estimators based on estimating equations.

However, these conditions will become more intricate than the ones we have provided here, as the concavity of the log densities is a rather strong property. As a result, we expect that assessing whether these conditions are satisfied in practice be more difficult.

### 3.3.3 Efficient Importance Sampling (EIS)

Efficient Importance Sampling (EIS) (Richard and Zhang, 2007) provides an alternative to the CE-method. Instead of minimizing the KL-divergence between the target  $\mathbf{P}$  and proposal  $\mathbf{G}_\psi$ ,  $\psi \in \Psi$ , EIS aims at minimizing the variance of the logarithm of importance sampling weights. Our discussion of (Chatterjee and Diaconis, 2018), Theorem 3.2, especially Lemma 3.5, suggests that this is worthwhile. Thus, EIS finds  $\psi_{\text{EIS}}$  which is a feasible solution to the following optimization problem

$$\min_{\psi \in \Psi} \text{Var}_{\mathbf{P}} [\log w_\psi] = \min_{\psi \in \Psi} \mathbf{P} [\log w_\psi - \mathbf{P} \log w_\psi]^2, \quad (3.17)$$

where, as in the last section,  $\log w_\psi = \log p - \log g_\psi$ .

Two problems arise:  $\mathbf{P}[\log w_\psi] = \mathcal{D}_{\text{KL}}(\mathbf{P}||\mathbf{G}_\psi)$  is usually intractable and we usually only have access to the unnormalized weights  $\frac{\bar{w}_\psi}{\mathbf{G}_\psi[w_\psi]} = w_\psi$ , with unknown integration constant  $\mathbf{G}_\psi[w_\psi]$ .

Both can be dealt with by introducing the nuisance parameter  $\lambda = \mathbf{P}[\log \tilde{w}_\psi]$ , utilizing the fact that the mean is the minimizer of the squared distance functional with the minimum value equal to the variance, should it exist. Indeed

$$\log w_\psi - \mathbf{P}[\log w_\psi] = \log \tilde{w}_\psi - \log \mathbf{G}_\psi[\tilde{w}_\psi] - \mathbf{P}[\log \tilde{w}_\psi] + \log \mathbf{G}_\psi[\tilde{w}_\psi] = \log \tilde{w}_\psi - \mathbf{P}[\log \tilde{w}_\psi],$$

so

$$\min_{\psi \in \Psi} \mathbf{P}[\log w_\psi - \mathbf{P}[\log w_\psi]]^2 = \min_{\psi \in \Psi, \lambda \in \mathbf{R}} \mathbf{P}[\log \tilde{w}_\psi - \lambda]^2,$$

where  $\psi \in \Psi$  is a minimizer of the left-hand side if, and only if,  $(\psi, \lambda) \in \Psi \times \mathbf{R}$  with  $\lambda = \mathbf{P}[\log \tilde{w}_\psi]$  is a minimizer of the right-hand side.

Similar to the CE-method we restrict our in-depth analysis to natural exponential family proposals where

$$\log g_\psi(x) = \psi^T T(x) - \log Z(\psi).$$

In this case the optimization problem is reduced to

$$\min_{\psi \in \Psi, \lambda \in \mathbf{R}} \mathbf{P}[\log p - \psi^T T - \lambda]^2, \quad (3.18)$$

a weighted linear least squares problem. As we consider unnormalized weights  $\tilde{w}$ , we are additionally able to get rid of the potentially non-linear term  $\log Z(\psi)$ . Noticing that this is a convex objective function in  $\psi$  which, similar to the CE-method, will be very useful to derive asymptotics later on. For now, we begin with studying the existence and uniqueness of  $\psi_{\text{EIS}}$  similar to Proposition 3.2.

**Lemma 3.9** (EIS for exponential families). *Let  $(\mathbf{G}_\psi)_{\psi \in \Psi}$  form a  $k$ -dimensional natural exponential family with log-densities*

$$\log g_\psi(x) = \psi^T T(x) - \log Z(\psi)$$

*for  $\Psi \subseteq \mathbf{R}^k$ . Suppose that  $\log p, T \in L^2(\mathbf{P})$ .*

*If there is a  $\psi_{\text{EIS}} \in \Psi$  with*

$$\text{Cov}_{\mathbf{P}}(T) \psi_{\text{EIS}} = \text{Cov}_{\mathbf{P}}(T, \log p) \quad (3.19)$$

*it is a global minimizer of Equation (3.17). If  $\text{Cov}_{\mathbf{P}}(T)$  is non-singular,*

$$\psi_{\text{EIS}} = \text{Cov}_{\mathbf{P}}(T)^{-1} \text{Cov}_{\mathbf{P}}(T, \log p)$$

*is the unique global minimizer.*

*Proof.* Under the proposed conditions, we may consider Equation (3.18) instead, where the moment conditions on  $\log p$  and  $T$  ensure that the problem is well-posed, i.e. the target is finite for all  $\psi \in \Psi$ . Thus the optimal  $(\psi_{\text{EIS}}, \lambda_{\text{EIS}})$  are given by the best linear unbiased predictor (BLUP) of  $\log p$  by the sufficient statistic  $T$  under  $\mathbf{P}$  for  $\psi_{\text{EIS}}$  and  $\mathbf{P}[\log \tilde{w}_{\psi_{\text{EIS}}}]$  for  $\lambda_{\text{EIS}}$ . Standard results from multivariate regression theory imply that the BLUP is given by any solution of

$$\text{Cov}_{\mathbf{P}}(T) \psi_{\text{EIS}} = \text{Cov}_{\mathbf{P}}(T, \log p),$$

i.e.  $\psi_{\text{EIS}}$  as stated in the lemma. Furthermore, if  $\text{Cov}_{\mathbf{P}}(T)$  is non-singular, the solution to this equation is unique. □

As the optimal  $\psi_{\text{EIS}}$  depends on several unknown quantities, EIS proceeds like the CE-method and employs importance sampling with a proposal  $\mathbf{G}$ , estimating  $\psi_{\text{EIS}}$  by

$$(\hat{\lambda}, \hat{\psi}_{\text{EIS}}) = \text{argmin}_{\lambda, \psi} \hat{\mathbf{P}}_N[\log \tilde{w}_\psi - \lambda]$$



where  $X^1, \dots, X^N \stackrel{\text{i.i.d.}}{\sim} \mathbf{G}$ . Again, if  $\mathbf{G}_\psi, \psi \in \Psi$  form an exponential family with natural parameter  $\psi$ , this optimization problem turns into a weighted least squares problem, so we can estimate  $\psi_{\text{EIS}}$  with the standard weighted least squares estimator

$$(\hat{\lambda}', \hat{\psi}_{\text{EIS}}) = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} y$$

where the random design matrix  $\mathbf{X}$ <sup>5</sup> and diagonal weights matrix  $\mathbf{W}$  are given by

$$\mathbf{X} = \begin{pmatrix} 1 & T(X^1)^T \\ \dots & \dots \\ 1 & T(X^N)^T \end{pmatrix}$$

and

$$\mathbf{W} = \text{diag}(W^1, \dots, W^N),$$

and the observations are

$$y = (\log p(X^1), \dots, \log p(X^N))^T \in \mathbf{R}^N.$$

Alternatively, replacing  $\mathbf{P}$  by  $\hat{\mathbf{P}}_N$  in Equation (3.19), we obtain the equivalent formulation

$$\hat{\psi}_{\text{EIS}} = \text{Cov}_{\hat{\mathbf{P}}_N}(T)^{-1} \text{Cov}_{\hat{\mathbf{P}}_N}(T, \log p), \quad (3.20)$$

as long as  $\text{Cov}_{\hat{\mathbf{P}}_N} T$  is non-singular.

An attractive feature of EIS is that if the target  $\mathbf{P}$  is a member of the exponential family of proposals, i.e. there is a  $\psi_{\mathbf{P}} \in \Psi$  such that  $\mathbf{P} = \mathbf{G}_{\psi_{\mathbf{P}}}$ , then EIS finds the optimal  $\psi_{\text{EIS}} = \psi_{\mathbf{P}}$  a.s. for a finite number of samples.

**Proposition 3.5** (Finite sample convergence of EIS). *Suppose  $\mathbf{G}_\psi, \psi \in \Psi \subseteq \mathbf{R}^k$  for a natural exponential family w.r.t. Lebesgue measure, where the support of the sufficient statistic  $\text{supp } T$  is open in  $\mathbf{R}^k$ . Furthermore let  $\mathbf{G}$  be a probability measure on  $\mathbf{R}^m$  that is equivalent to  $\mathbf{P}$ , i.e.  $\mathbf{G} \ll \mathbf{P}$  and  $\mathbf{P} \ll \mathbf{G}$ .*

*If there is a  $\psi_{\mathbf{P}} \in \Psi$  such that  $\mathbf{P} = \mathbf{G}_{\psi_{\mathbf{P}}}$ , then  $\hat{\psi}_{\text{EIS}} = \psi_{\mathbf{P}}$  a.s. for  $N \geq k$ .*

*Proof.* As  $\mathbf{P}$  stems from the same exponential family as  $\mathbf{G}_\psi$ , the pseudo-observations are

$$\log p = \psi_{\mathbf{P}}^T T - \log Z(\psi_{\mathbf{P}}).$$

Thus  $\text{Cov}_{\hat{\mathbf{P}}_N}(T, \log p) = \text{Cov}_{\hat{\mathbf{P}}_N}(T) \psi_{\mathbf{P}}$ . If we can show that  $\text{Cov}_{\hat{\mathbf{P}}_N} T$  is non-singular, Equation (3.20) implies that  $\hat{\psi}_{\text{EIS}} = \psi_{\mathbf{P}}$  a.s..

If  $\text{Cov}_{\hat{\mathbf{P}}_N} T$  were singular, there would exist a  $\psi \in \mathbf{R}^k$  such that

$$\psi^T \text{Cov}_{\hat{\mathbf{P}}_N}(T) \psi = \text{Cov}_{\hat{\mathbf{P}}_N}(\psi^T T) = 0.$$

In this case the a.s. non-zero  $W^i(X^i)T(X^i)$  would lie in the orthogonal complement  $\psi^\perp$  for all  $i = 1, \dots, N$ . As the weights are a.s. positive by the assumed equivalence of  $\mathbf{G}$  and  $\mathbf{P}$ , the same holds true for  $T(X^i), i = 1, \dots, N$ . If  $N$  is bigger than  $k$ , the probability that this happens is 0, as  $\text{supp } T$  is open. Thus  $\text{Cov}_{\hat{\mathbf{P}}_N} T$  is non-singular almost surely and the result is shown.  $\square$

Note that if in the above proposition only  $\mathbf{G}_\psi \gg \mathbf{P}$  holds, we obtain, by a similar argument, that

$$\mathbb{P}(\hat{\psi}_{\text{EIS}} = \psi_{\mathbf{P}}) \xrightarrow{N \rightarrow \infty} 1.$$

<sup>5</sup>if  $\mathbf{X} \mathbf{W} \mathbf{X}$  is not invertible, replace the inverse by the Moore-Penrose pseudoinverse



Additionally, we then have to take care of the event  $\{w(X) = 0\}$ , whose probability is now potentially positive.

We now turn to deriving asymptotics for  $\hat{\psi}_{\text{EIS}}$ . As for the CE-method, we start with proving that  $\hat{\psi}_{\text{EIS}}$  consistently estimates  $\psi_{\text{EIS}}$ . For this we need to ensure that  $\psi_{\text{EIS}}$  is the unique solution to Equation (3.17), as otherwise, consistent estimators of  $\psi_{\text{EIS}}$  cannot exist. As Equation (3.18) is a linear least squares problem, the objective function is convex, and so we can apply Theorem 3.5 and Proposition 3.3.

**Theorem 3.8** (consistency of  $\hat{\psi}_{\text{EIS}}$ ). *Let  $(\mathbf{G}_\psi)_{\psi \in \Psi}$  form a  $k$ -dimensional natural exponential family with log-densities*

$$\log g_\psi(x) = \psi^T T(x) - \log Z(\psi)$$

*for convex  $\Psi \subseteq \mathbf{R}^k$ . Let  $\mathbf{G} \gg \mathbf{P}$  be a proposal and suppose that*

- (i)  $\log p, T \in L^2(\mathbf{P})$  and
- (ii)  $\text{Cov}_{\mathbf{P}}(T)$  is non-singular,
- (iii)  $\psi_{\text{EIS}} \in \text{int } \Psi$ .

*Then*

$$\hat{\psi}_{\text{EIS}} \xrightarrow{N \rightarrow \infty} \psi_{\text{EIS}}$$

*almost surely.*

*Proof.* We follow the same strategy as in the proof of Theorem 3.4. Let

$$b : \mathbf{R}^p \times \mathbf{R}^{k+1} \rightarrow [-\infty, \infty) \quad b(x, \psi') = \begin{cases} -\frac{1}{2} (\log p(x) - \psi'^T T(x) - \lambda)^2 & \psi' \in \Psi \\ -\infty & \text{else,} \end{cases}$$

where  $\psi' = (\psi, \lambda) \in \mathbf{R}^{k+1}$ . For fixed  $x$  this function is concave, as its Hessian is negative semi-definite:

$$H_{\psi'} b(x, \psi') = - \begin{pmatrix} 1 & T(x)^T \\ T(x) & T(x)T(x)^T \end{pmatrix} = - \begin{pmatrix} 1 & T(x)^T \end{pmatrix} \begin{pmatrix} 1 & T(x)^T \end{pmatrix}^T,$$

if  $\psi' \in \Psi$ . Let  $X^1, \dots, X^N \stackrel{\text{i.i.d.}}{\sim} \mathbf{P}$  and let  $\tilde{\mathbf{P}}_N$  be their empirical distribution. For  $\psi \in \Psi, \lambda \in \mathbf{R}$  we have

$$\mathbf{P}[b(\cdot, \psi')] = -\frac{1}{2} \mathbf{P}[(\log p - \psi'^T T - \lambda)^2] < \infty,$$

as  $\log p, T \in L^2(\mathbf{P})$ . Let us now check that conditions (C1) - (C3) are fulfilled.

(C1) is fulfilled, as we assumed  $\psi_{\text{EIS}} \in \text{int } \Psi$ . (C2) holds, as  $\psi_{\text{EIS}}$  is the unique global maximizer by Lemma 3.9, as  $\text{Cov}(T)$  is non-singular. (C3) obviously holds.

Thus  $\hat{\psi}_{\text{EIS}}$  is strongly consistent if  $\mathbf{G} = \mathbf{P}$ . If  $\mathbf{G}$  is different from  $\mathbf{P}$ , we can apply Proposition 3.3, where the existence of M-estimators is ensured by Equation (3.20), using the Moore-Penrose inverse if  $\text{Cov}_{\tilde{\mathbf{P}}_N}(T)$  is singular.  $\square$

As Equation (3.20) expresses  $\hat{\psi}_{\text{EIS}}$  in terms of empirical covariances, we could alternatively prove consistency by ensuring that the empirical covariances are consistent as well, for which we would need to ensure that fourth-order moments of  $\log p$  and  $T$  w.r.t.  $\mathbf{P}$  exist. This strategy may be fruitful if  $\psi_{\text{EIS}}$  does not lie in the interior of  $\Psi$ , although the more sophisticated treatment of (Haberman, 1989) may also be applicable under these circumstances.

Additionally, if fourth-order moments exist, we can derive a central limit theorem, similar to Theorem 3.6, for EIS.

**Theorem 3.9** (CLT for  $\hat{\psi}_{\text{EIS}}$ ). *Let  $(\mathbf{G}_\psi)_{\psi \in \Psi}$  form a  $k$ -dimensional natural exponential family with log-densities*

$$\log g_\psi(x) = \psi^T T(x) - \log Z(\psi),$$

*and convex parameter space  $\Psi \subseteq \mathbf{R}^k$ . Let  $\mathbf{G} \gg \mathbf{P}$  be a proposal with weights  $w = \frac{d\mathbf{P}}{d\mathbf{G}}$ .*

- (i)  $wT_iT_j, w(\log p)^2 \in L^2(\mathbf{G})$  for  $i, j = 1, \dots, k$ ,
- (ii)  $\log p, T_i \in L^4(\mathbf{P})$  for all  $i = 1, \dots, k$
- (iii)  $\text{Cov}_{\mathbf{P}}(T)$  is non-singular and  $\psi_{\text{EIS}} \in \text{int } \Psi$ .

Then

$$\sqrt{N}(\hat{\psi}_{\text{EIS}} - \psi_{\text{EIS}}) \xrightarrow{\mathcal{D}} \mathcal{N}(0, BMB)$$

where  $B = \text{Cov}_{\mathbf{P}}(T)^{-1}$  and

$$M = \text{Cov}_{\mathbf{G}} \left( w \left( \log p - \psi_{\text{EIS}}^T T - \lambda_{\text{EIS}} - \mathbf{P}[T] \right) T \right).$$

*Proof.* Similar to the proof of Theorem 3.6, we combine Theorem 3.7 and Proposition 3.4. Let

$$b : \mathcal{X} \times \mathbf{R}^{k+1} \rightarrow [-\infty, \infty) \quad b(x, \psi') = \begin{cases} -\frac{1}{2} (\log p(x) - \psi'^T T'(x)) & x \in \Psi \\ -\infty & \text{else,} \end{cases}$$

where  $\psi' = (\psi, \lambda) \in \Psi \times \mathbf{R}$  and  $T'(x) = (T(x) \ 1)$ . For  $\psi \in \Psi$  the map  $(\psi, \lambda) \rightarrow \mathbf{P}[b(\cdot, (\psi, \lambda))]$  is differentiable with gradient

$$\nabla_{\psi'} \mathbf{P}[b(\cdot, \psi')] = -\mathbf{P}[(\log p - \psi'^T T') T'] = \begin{pmatrix} -\mathbf{P}[T' \log p - T'^T T' \psi'] \\ -\mathbf{P}[\log p - \psi'^T T'] \end{pmatrix}$$

and Hessian

$$H_{\psi'} \mathbf{P}[b(\cdot, \psi')] = -\mathbf{P}[T' T'^T] = -\begin{pmatrix} \mathbf{P}[T T^T] & \mathbf{P}[T^T] \\ \mathbf{P}[T] & 1 \end{pmatrix}.$$

The Hessian is negative definite, as for all  $\psi \in \mathbf{R}^k, \lambda \in \mathbf{R}$  we have

$$\begin{aligned} (\psi^T \ \lambda) H_{\psi'} \mathbf{P}[b(\cdot, \psi')] (\psi^T \ \lambda)^T &= -(\psi^T \text{Cov}_{\mathbf{P}}(T) \psi + \psi^T \mathbf{P}[T] \mathbf{P}[T]^T \psi + 2\psi^T \mathbf{P}[T] \lambda + \lambda^2) \\ &= -(\psi^T \text{Cov}_{\mathbf{P}}(T) \psi + (\lambda + \psi^T \mathbf{P}[T])^2) \leq 0, \end{aligned}$$

with equality if, and only if, both  $\lambda$  and  $\psi$  are 0, as  $\text{Cov}_{\mathbf{P}}(T)$  is assumed to be positive definite. Thus condition (C7) is fulfilled.

For condition (C10), we can verify that for all  $i, j = 1, \dots, k+1$

$$(\nabla_{\psi'} b(\cdot, \psi'))_i (\nabla_{\psi'} b(\cdot, \psi'))_j = (\log p - \psi'^T T')^2 T'_i T'_j$$

is in  $L^1(\mathbf{P})$  by assumption (ii) and the Hölder inequality.

To apply Proposition 3.4 we need to show that  $w(\cdot) b'(\cdot, \psi', \xi') \in L^2(\mathbf{G})$  for all  $\xi' \in \mathbf{R}^{k+1}$  and all  $\psi'$  in a neighborhood of  $\psi_{\text{EIS}}$ , for this it suffices that we show

$$w^2 (\nabla_{\psi'} b(\cdot, \psi'))_i (\nabla_{\psi'} b(\cdot, \psi'))_j = w^2 (\log p - \psi'^T T')^2 T'_i T'_j$$

is in  $L^1(\mathbf{G})$ , which holds, again, by assumption Item (i) and the Hölder inequality.

We have thus shown a central limit theorem for  $\hat{\psi}'_{\text{EIS}} = (\hat{\psi}_{\text{EIS}}, \hat{\lambda}_{\text{EIS}})$ , i.e.

$$\sqrt{N}(\hat{\psi}'_{\text{EIS}} - \psi_{\text{EIS}}) \rightarrow \mathcal{N}(0, M' B' M')$$

with  $B' = -(H_{\psi'_{\text{EIS}}} \mathbf{P}[b(\cdot, \psi'_{\text{EIS}})])^{-1}$  and  $M' = \text{Cov}(w(X) \nabla_{\psi'_{\text{EIS}}} b(X, \psi'_{\text{EIS}}))$  for  $X \sim \mathbf{G}$ . By using the inversion formula for block matrices, we obtain

$$\begin{aligned} B' &= \begin{pmatrix} \mathbf{P}[T T^T] & \mathbf{P}[T^T] \\ \mathbf{P}[T] & 1 \end{pmatrix}^{-1} = \begin{pmatrix} \Sigma + \mu \mu^T & \mu^T \\ \mu & 1 \end{pmatrix}^{-1} \\ &= \begin{pmatrix} (\Sigma + \mu \mu^T - \mu \mu^T)^{-1} & 0 \\ 0 & 1 - \mu^T (\Sigma + \mu \mu^T)^{-1} \mu \end{pmatrix} \begin{pmatrix} I_k & -\mu^T \\ -\mu (\Sigma + \mu \mu^T)^{-1} & 1 \end{pmatrix} \\ &= \begin{pmatrix} \Sigma^{-1} & -\mu^T \Sigma^{-1} \\ -\Sigma^{-1} \mu & 1 - \mu^T (\Sigma + \mu \mu^T)^{-1} \mu \end{pmatrix} \end{aligned}$$

where  $\Sigma = \text{Cov}_{\mathbf{P}}(T)$  and  $\mu = \mathbf{P}[T]$ . Similarly,

$$M' = \begin{pmatrix} \text{Cov}_{\mathbf{G}}(wW_{\psi_{\text{EIS}}}T) & \text{Cov}_{\mathbf{G}}(wW_{\psi_{\text{EIS}}}T, wW_{\psi_{\text{EIS}}}) \\ \text{Cov}_{\mathbf{G}}(wW_{\psi_{\text{EIS}}}, wW_{\psi_{\text{EIS}}}T) & \text{Cov}_{\mathbf{G}}(wW_{\psi_{\text{EIS}}}) \end{pmatrix},$$

where  $W_{\psi_{\text{EIS}}} = \log p - \psi'_{\text{EIS}}T$ .

If  $\mu \neq 0$ , we may change the sufficient statistic of the exponential family such that this holds, i.e. let  $\tilde{T} = T - \mathbf{P}[T]$ , then

$$\log g_{\psi}(x) = \psi^T T(x) - \log Z(\psi) = \psi^T \tilde{T}(x) - \log \tilde{Z}(\psi)$$

where  $\tilde{Z}(\psi) = \log Z(\psi) + \mathbf{P}[T]$ . As  $\psi_{\text{EIS}}$ , Equation (3.19), only depends on  $T - \mathbf{P}[T]$  under  $\mathbf{P}$ , this does not change  $\psi_{\text{EIS}}$ . Similarly,  $\hat{\psi}_{\text{EIS}}$ , Equation (3.20), is unaffected by subtracting a constant from  $T$ . Only

$$\tilde{\lambda}_{\text{EIS}} = \lambda_{\text{EIS}} + \mathbf{P}[T]$$

and similarly  $\hat{\lambda}_{\text{EIS}}$  are changed.

Thus, without loss of generality, we may assume that  $\mathbf{P}[T] = 0$ . Then

$$B' = \begin{pmatrix} \Sigma^{-1} & 0 \\ 0 & 1 \end{pmatrix}$$

is a diagonal matrix. Taking the  $\psi_{\text{EIS}}$  marginal of the asymptotic normal distribution, we arrive at

$$\sqrt{N}(\hat{\psi}_{\text{EIS}} - \psi_{\text{EIS}}) \rightarrow \mathcal{N}(0, BMB)$$

with  $B = \text{Cov}_{\mathbf{P}}(T)$  and  $M = \text{Cov}_{\mathbf{G}}\left(w\left(\log p - \psi_{\text{EIS}}^T T - \lambda_{\text{EIS}} - \mathbf{P}[T]\right)T\right)$ , as promised.  $\square$

A discussion of the assumptions of Theorems 3.8 and 3.9 is in order. We start with the consistency result Theorem 3.8. The integrability condition, i.e that  $\log p, T \in L^2(\mathbf{P})$  is necessary to ensure existence of  $\psi_{\text{EIS}}$  and the existence of  $\text{Cov}_{\mathbf{P}}(T)$  as well as  $\psi_{\text{EIS}} \in \text{int } \Psi$  ensure uniqueness, see also Lemma 3.9.

Regarding the central limit theorem Theorem 3.9, requiring the existence of higher order moments is natural. Unfortunately, there is no direct interpretation of these requirements as generalizations of the existence of  $\rho$ , as was the case for the CE-method.

The only integrability conditions related to the proposal  $\mathbf{G}$  are those for  $wT_i T_j$  and  $w\log(p)^2$ . Choosing  $(\mathbf{G}_{\psi})_{\psi \in \Psi}$  to consist of Gaussian distributions, the conditions on  $T$  translate to the existence of certain polynomial moments of  $w^2$  w.r.t. the proposal  $\mathbf{G}$  (or  $w$  w.r.t.  $\mathbf{P}$ ). This technical condition, is not easily interpreted, as assuming existence of moments of the target distribution seem more natural than those involving the extra weighting term  $w = \frac{p}{g}$ , which depends on the proposal  $\mathbf{G}$  as well.

### 3.4 Interim discussion

Before we apply EIS and the CE-method in the SSM context, let us consolidate what we have achieved by the asymptotic analysis in the preceding two subsections and reason which of the two methods should be used in which circumstances.

We start with a discussion of the optimal values  $\psi_{\text{CE}}$  and  $\psi_{\text{EIS}}$ . Notice that  $\psi_{\text{EIS}}$  depends on second-order moments of the sufficient statistic  $T$ , as well as the shape of  $\log p$ , whereas the optimal parameter for the CE-method  $\psi_{\text{CE}}$  depends only on the first-order moments of  $T$ . This dependence on higher-order moments may be beneficial for the EIS method, for example, if the covariance of  $T$  under  $\mathbf{P}$  is very different from that under  $\mathbf{G}_{\psi}$ .

should have an example for this later

The two methods differ concerning the assumptions that are required for uniqueness, consistency and the central limit theorem to hold if the proposals come from an exponential family. For uniqueness, Proposition 3.2 and lemma 3.9, both methods require that the covariance of  $T$  is non-singular, however, the measures under which the covariance are considered differ: for the CE-method we need  $\text{Cov}_{\mathbf{G}_{\psi_{\text{CE}}}}(T)$  to be non-singular, while for EIS the same has to hold for  $\text{Cov}_{\mathbf{P}}(T)$ . While the former is easy to ensure, the latter depends on the intractable target  $\mathbf{P}$  and may be more difficult to verify in practice, depending on  $T$ .

Regarding the consistency results, Theorems 3.4 and 3.8 as well as the central limit theorems, Theorems 3.6 and 3.9, EIS requires that the sufficient statistic be twice as often  $\mathbf{P}$ -integrable as the CE-method. Additionally, the EIS results assume that  $\log p$  is sufficiently often  $\mathbf{P}$ -integrable. Therefore, EIS is, at first glance, more restrictive than the CE-method. However, our application will perform importance sampling with Gaussian proposals where  $T(x) = \begin{pmatrix} x \\ xx^T \end{pmatrix}$ . For importance sampling to be consistent in this setting, we have to assume that the target has thinner tails than the Gaussian proposal, which implies that all polynomial moments of the target, and thus of  $T$  exist. A similar argument can be made for  $\log p$ , and so the assumptions are likely to be fulfilled when Gaussian importance sampling is consistent.

To compare the asymptotic covariance matrices of both methods, note that both covariance matrices have the same „bread-meat-bread“ factorization, as they are asymptotic covariance matrices of M-estimators<sup>6</sup>. We see that both  $B_{\text{CE}} = I(\psi) = \text{Cov}_{G_{\psi_{\text{CE}}}}(T)^{-1}$  and  $B_{\text{EIS}} = \text{Cov}_{\mathbf{P}}(T)^{-1}$  are precision matrices of the sufficient statistic  $T$ , one with respect to the optimal CE-method proposal and one with respect to the target. Thus, if  $\mathbf{P}$  is well approximated by  $\mathbf{G}_{\psi_{\text{CE}}}$ , we would expect these two components to be close to one another. For  $M_{\text{CE}} = \text{Cov}_{\mathbf{G}}(wT)$  and  $M_{\text{EIS}} = \text{Cov}_{\mathbf{G}}(w(\log p - \psi_{\text{EIS}}^T T - \lambda_{\text{EIS}} - \mathbf{P}[T])T)$ , there is a more notable difference, i.e. the presence of the  $\log p - \psi_{\text{EIS}}^T T - \lambda_{\text{EIS}}$  term. If the EIS approximation performs well, we can expect this term to be small, as it is the prediction error of the least squares approximation of  $\log g_{\psi}$  to  $\log p$ . Therefore, we expect that EIS outperforms the CE-method in terms of asymptotic variance in these settings. In agreement with Proposition 3.5,  $M_{\text{EIS}} = 0$  if  $\log p = \log g_{\psi_{\mathbf{P}}}$  so that  $\psi_{\text{EIS}} = \psi_{\mathbf{P}}$ .

Additionally, both  $M_{\text{CE}}$  and  $M_{\text{EIS}}$  depend on the proposal  $\mathbf{G}$ , and indicate how one might tailor the initial proposal  $\mathbf{G}$  to produce low-variance estimates. For the CE-method we might choose  $\mathbf{G}$  such that the trace determinant of  $\mathbf{G}[w^2 T T^T]$  becomes small. This is not necessarily achieved by the CE-method proposal  $\mathbf{G}_{\hat{\psi}_{\text{CE}}}$ , and so it may be worthwhile to investigate using two types of proposals in the CE-method, one that makes  $M_{\text{CE}}$  small and  $\mathbf{G}_{\psi_{\text{CE}}}$ . This is especially relevant as our simulation studies, Section 3.7, suggest that the asymptotic covariance of the CE-method is usually larger than that of EIS. For EIS, a similar approach might be fruitful, but is not as urgent as that for the CE-method, as the asymptotic covariance of EIS is usually small enough to be feasible in practice.

Finally, let us stress that these asymptotic considerations are, to the author’s knowledge, novel results and should be straightforward to extend if the proposals  $(\mathbf{G}_{\psi})_{\psi \in \Psi}$  do not form a natural exponential family. As any minimal exponential family may be reduced to a natural exponential family by reparametrization, see (Brown, 1986, Theorem 1.9), the delta method can be used to derive CLTs in this case as well, as Proposition 3.2 and lemma 3.9 still apply. If the family is not minimal the optimal values  $\psi_{\text{EIS}}$  and  $\psi_{\text{CE}}$  may be non-unique, so we cannot hope to estimate them consistently. In this case the user should choose a minimal parametrization, see again (Brown, 1986, Theorem 1.9). For non-exponential family proposals our results should also carry over, provided the usual regularity conditions ensuring uniqueness, consistency and asymptotic normality for M-estimators hold. If the objective functions are not concave as they are in our setting one usually requires uniformly bounded third-order derivatives of the objective function to exist.

Furthermore, our results can also be extended to the so-called Variance-Minimization method

<sup>6</sup>These are sometimes called sandwich estimators.

(VM-method) which determines an optimal proposal by solving the following optimization problem:

$$\min_{\psi \in \Psi} \text{Var}_{\mathbf{G}_\psi}(w_\psi) = \min_{\psi \in \Psi} \mathbf{G}_\psi[w_\psi^2] = \min_{\psi \in \Psi} \mathbf{P}[w_\psi],$$

where the first equality holds as  $\mathbf{G}_\psi[w_\psi] = 1$  for all  $\psi$ . Thus the VM-method chooses  $\psi$  such that the second moment of importance sampling weights,  $\rho$ , becomes small. Again, this is sensible by the discussion surrounding  $\rho$  and the ESS. Similar to the CE-method and EIS, one uses importance sampling with a proposal  $\mathbf{G}$  to approximate  $\mathbf{P}[w_\psi]$  by  $\hat{\mathbf{P}}_N[w_\psi]$ , and solves this noisy version of the problem. Unfortunately, there is no closed form for the optimal  $\psi_{\text{VM}}$  or  $\hat{\psi}_{\text{VM}}$ , even if the proposals form a natural exponential family. Still, as  $x \mapsto w_\psi(x)$  is convex, so is  $\psi \mapsto \mathbf{P}[w_\psi]$ , and we can apply Theorems 3.5 and 3.7 in combination with Propositions 3.3 and 3.4 to show, under suitable regularity conditions, the consistency and asymptotic normality of the method.

Now that we have gained theoretical insight into optimal importance sampling, let us apply these insights to the SSMs that we are interested in.

### 3.5 Gaussian importance sampling for state space models

For the types of models considered in this thesis, importance sampling is used to infer the posterior distribution. Given a state space model of the form (3.1) and observations  $Y = Y_{:n}$ , let  $\mathbf{P}$  be the distribution of the states  $X = X_{:n}$ , conditional on  $Y$  and  $f$  be a function of interest. The task at hand is now to find a suitable proposal  $\mathbf{G}$ , using the methods presented in the last section. If  $n$  is large, the posterior distribution lives in a high dimensional state of dimension  $m \cdot n$ , so to obtain  $\mathbf{G}$  efficiently, we should exploit the available structure. Additionally, we want  $\mathbf{G}$  to be tractable, so simulating from it is possible and evaluating the weights  $w$  up to a constant is possible.

The multivariate Gaussian distribution is a good candidate in this setting, as simulating from it is straightforward and its density can be evaluated analytically. However, naively performing the optimal importance sampling methods from the previous section for all multivariate Gaussians is computationally inefficient as the family of distributions has  $\mathcal{O}((n \cdot m)^2)$  many parameters. We can, however, exploit the available structure of the SSM to find parameterizations with fewer parameters by either using smoothing distributions of GLSSMs (Section 3.5.1) or approximating with a Gaussian discrete-time Markov process (Section 3.5.2).

Using Gaussian proposals, while computationally efficient, also comes with some drawbacks. The whole procedure hinges on the assumption that there is a Gaussian that is close to the target distribution. In the setting of SSMs this is not guaranteed, as the targets may contain multiple modes or heavy tails, features that may, in the worst case, lead to inconsistent importance sampling estimates. Additionally, even if there is a Gaussian distribution that facilitates consistent importance sampling, finding it in practice may be complicated, as the proposals generated by the LA, CE-method and EIS have deteriorating performance for fixed sample size  $N$  (in terms of ESS and convergence) with increasing dimension, see Section 3.7.4.

Using a GLSSM as an importance sampling proposal for non-Gaussian state space models was introduced by (Durbin and Koopman, 1997) to facilitate maximum likelihood estimation using the LA as a proposal. Concurrently, (Shephard and Pitt, 1997) established a similar result in the context of MCMC analysis of SSMs.

#### 3.5.1 The GLSSM-approach

The first approach is motivated by the fact that the target posterior is again a Markov process, as are posteriors in GLSSMs. Additionally, the posterior distribution in GLSSMs is again Gaussian, and straightforward to simulate from by, e.g., the FFBS algorithm (Algorithm 3) or the simulation smoother (Durbin and Koopman, 2002). Thus parameterizing the proposals  $\mathbf{G}$  by the posterior of a suitably chosen GLSSM may be a fruitful approach. For the models we consider in this thesis, the distribution of states is already Gaussian and the observations are conditionally independent given the states. Thus a natural GLSSM to use as a proposal consists of keeping the prior distribution of states and replacing the distribution of observations with conditionally independent

Gaussian distributions and the actual observations by synthetic ones. By the assumed conditional independence, this model only needs  $2p \cdot (n + 1)$  many parameters,  $p \cdot (n + 1)$  for the synthetic observations and  $p \cdot (n + 1)$  for their variances. We term this approach the **GLSSM-approach** to importance sampling.

In total, the GLSSM-approach considers parametric proposals  $\mathbf{G}_\psi$  of the form

$$\begin{aligned}\mathbf{G}_\psi &= \mathcal{L}(X|Z = z), \\ Z_t &= B_t X_t + \eta_t, \\ \eta_t &\sim \mathcal{N}(0, \Omega_t), \\ \Omega_t &= \text{diag}(\omega_t^2) = \text{diag}(\omega_{t,1}^2, \dots, \omega_{t,p}^2).\end{aligned}\tag{3.21}$$

where the distribution of  $X$  is given by (3.4),  $\psi = (z, \omega^2)$  for  $z = (z_0, \dots, z_n) \in \mathbf{R}^{(n+1) \times m}$  and  $\omega^2 = (\omega_0^2, \dots, \omega_n^2) \in \mathbf{R}^{(n+1) \times m}$ . Alternatively the natural parametrization

$$\psi = (z \oslash \omega^2, -1 \oslash (2\omega^2))\tag{3.22}$$

may also be used, where  $\oslash$  is the Hadamard, i.e. entry-wise, division. Simulation from  $\mathbf{G}_\psi$  may be efficiently implemented by the FFBS algorithm, as  $\mathbf{G}_\psi$  is the smoothing distribution of a GLSSM.

In this setting, the importance sampling weights are given by

$$w(x) = \frac{p(x|y)}{g(x|z)} = \frac{p(y|x)p(x)}{g(z|x)p(x)} \frac{g(z)}{p(y)} \propto \prod_{t=0}^n \frac{p(y_t|x_t)}{g(z_t|x_t)},$$

so they can be computed efficiently. Additionally, for a EGSSM with linear signals,  $p(y_t|x_t)$  and  $g(z_t|x_t)$  depend on  $x_t$  only through the signal  $s_t = B_t x_t$ , and we have

$$w(x) \propto \prod_{t=0}^n \frac{p(y_t|s_t)}{g(z_t|s_t)},\tag{3.23}$$

which implies that auto-normalized weights may be calculated by using the signal smoother (Jungbacker and Koopman, 2007, Theorem 2). As (Durbin and Koopman, 2012) (Durbin and Koopman, 2012, Section 4.5.3) argue, it is often computationally more efficient to treat only on the signals  $S_{:,n}$  instead of the states  $X_{:,n}$ , the idea being that the dimension of  $S_t$ ,  $p$ , is usually much smaller than that of  $X_t$ ,  $m$ .

As the joint distribution of  $(X, S)$  is a Gaussian distribution, by Lemma 3.1  $X|S = s$  is again Gaussian, with known conditional mean and covariance matrix and density  $p(x|s) = g(x|s)$ . If  $(\tilde{X}_t)_{t=0, \dots, n}$  is a draw from this conditional distribution a quick calculation reveals that a.s.  $B_t \tilde{X}_t = S_t$ , and so, as expected, the weights  $w(\tilde{X}_t)$  are a.s. constant and given by (up to the integration constant) Equation (3.23). Producing a draw from this conditional distribution can be achieved by the FFBS algorithm (Algorithm 3), as  $(X, S)$  form a GLSSM with degenerate observation covariance matrices  $\Omega_t = 0$ .

By the assumed conditional independence of observations given signals, we have

$$p(x, s|y) \propto p(x|s)p(s|y),$$

and so if one is interested in the states, rather than the signals, importance sampling with the proposal Equation (3.21) can be achieved in a two-step procedure: first sample from  $g(s|z)$ , then run the FFBS algorithm to sample from  $g(x|s) = p(x|s)$  using the same weights for MC-integration.

The GLSSM-approach is the standard approach for finding the LA in EGSSM (Durbin and Koopman, 2012; Durbin and Koopman, 1997) and may even be applied when the observation densities are not log-concave as (Jungbacker and Koopman, 2007) show. The approach also leads to efficient implementation for EIS (Koopman, Lit, and Nguyen, 2019). However, as will become apparent in the later part of this section, it is infeasible for the CE-method if  $n$  is large.

We now give a concise overview over how to perform the LA and EIS for EGSSM, but refer the reader for more details to the respective literature. (Danielsson and Richard, 1993) were the first to propose minimizing the variance of log-weights, which developed into the EIS method (Liesenfeld and Richard, 2003; Richard and Zhang, 2007). However, as formulated in these earlier works, EIS requires careful tracking of integration constants. For PGSSMs, the modified EIS of (Koopman, Lit, and Nguyen, 2019) provides a more straightforward approach to determine the proposal distribution, by noticing that it may be written as the posterior of an appropriately chosen GLSSM, whose distribution of states coincides with that of the original PGSSM. Instead of approximating the integrals in EIS using MC-integration, (Koopman, Lucas, and Scharth, 2015) suggest to do so using numerical integration, to reduce sampling error.

The LA for PGSSM can be obtained efficiently, by noticing that the Newton-Raphson scheme to obtain the posterior mode (Equation (3.11)) can be implemented using the Kalman smoother, see (Durbin and Koopman, 2012, Chapter 10). For EGSSM with a linear signal, the LA is particularly easy to implement, as the conditional independence of individual observations translates to an approximating GLSSM with independent observations as well — the resulting algorithm is presented in Algorithm 5.

---

**Algorithm 5** The LA for EGSSM

---

**Require:** EGSSM (Definition 3.5) with linear signal and natural parameters  $s_t$ ,  $t = 0, \dots, n$ , observations  $y_0, \dots, y_n$ , initial values  $\psi^0 = (z^0, (\omega^2)^0)$

- 1: Set  $l = 0$ .
- 2: **repeat**
- 3:   Run the Kalman smoother Algorithm 2 for the model (3.21) to obtain  $\hat{s} = \mathbb{E}_{\mathbf{G}_{\psi^l}}(S|Z = z^l)$ .
- 4:   **for**  $t = 0, \dots, n$  **do**
- 5:     Set  $(\omega^{l+1})_{t,i}^2 = H_{s_t^i} Z_t(s_t^i)$  ▷ Hessian evaluated at  $\hat{s}_t$
- 6:     Set  $\Omega_t^{l+1} = \text{diag}((\omega^2)_{t,1}^{l+1}, \dots, (\omega^2)_{t,p}^{l+1})$ .
- 7:     Set  $z_t^{l+1} = z_t^l - (\Omega_t^{l+1})^{-1} (\nabla \log Z_t(s_t^1), \dots, \nabla \log Z_t(s_t^p))$ . ▷ Gradient evaluated at  $\hat{s}_t$
- 8:   **end for**
- 9:   Set  $\psi^{l+1} = (z^{l+1} \odot (\omega^2)^{l+1}, -1 \odot (\omega^2)^{l+1})$ .
- 10:   Set  $l = l + 1$ .
- 11: **until**  $\psi^l$  has converged.

---

As we will show in Section 3.7, the LA may provide poor performance as a proposal, when the dimension of the PGSSM grows. In this case, the EIS proposal may perform better. Recall from Section 3.3.3, that EIS aims at minimizing the mean squared error between  $s \mapsto \log p(s|y)$ , the target log-density and  $s \mapsto \log g_\psi(s|z)$ , the proposals log-density with respect to the target. Thus in the context of PGSSMs, EIS aims to minimize

$$(z, \omega^2) \mapsto \mathbb{P}^{S|Y=y} \left[ (\log p(\cdot|y) - \log g(\cdot|z))^2 \right].$$

By reformulating the integrand to  $(\log p(y|s) - \log g(z|s) - \log p(y) + \log g(z))^2$  and following the discussion surrounding Equation (3.18), we may instead minimize

$$\mathbb{P}^{S|Y=y} \left[ (\log p(y|\cdot) - \log g(z|\cdot) - \lambda)^2 \right] = \mathbb{P}^{S|Y=y} \left[ \left( \sum_{t=0}^n \log p(y_t|\cdot_t) - \log g(z_t|\cdot_t) - \lambda \right)^2 \right].$$

over the parameters of interest,  $(z, \omega^2)$ , and the nuisance parameter  $\lambda$ .

However, the dimension of  $\psi$  is quite high ( $2p(n+1)$ ), so one resorts to solving the lower dimensional problems

$$\mathbb{P}^{S|Y=y} \left[ (\log p(y_t|\cdot_t) - \log g(z_t|\cdot_t) - \lambda_t)^2 \right] = \mathbb{P}^{X_t|Y=y} \left[ (\log p(y_t|\cdot_t) - \log g(z_t|\cdot_t) - \lambda_t)^2 \right],$$

which only depends on the marginal of  $X_t$ .



EIS for PGSSM then proceeds as in Section 3.3.3, using importance sampling to obtain a particle approximation  $\hat{\mathbf{P}}_N$  to  $\mathbb{P}^{S|Y=y}$  and solving the resulting approximate problem. As

$$\log g(z_t|s_t) = -\frac{1}{2}(z_t - s_t)^T \Omega_t^{-1}(z_t - s_t) - \frac{p}{2} \log(2\pi) - \frac{1}{2} \log \det \Omega_t,$$

affine in the natural parameters  $(z \oslash \omega^2, -1 \oslash (2\omega^2))$ , the problem reduces to weighted linear least squares problem, which can be solved analytically. If the model at hand is an EGSSM with linear signal<sup>7</sup>, we can further exploit the independence of marginals, to solve only for the univariate marginals, i.e. minimize

$$(z_{t,i}, \omega_{t,i}^2, \lambda_{t,i}) \mapsto \mathbb{P}^{S|Y=y} \left[ (\log p(y_{t,i}|s_{t,i}) - \log g(z_{t,i}|s_{t,i}) - \lambda_{t,i})^2 \right], \quad (3.24)$$

or an importance sampling version of it.

We present the resulting algorithm in Algorithm 6, restricted to the case where the model is an EGSSM with linear signal. As starting values  $\psi$  we may take those obtained by the LA.

---

**Algorithm 6** EIS for EGSSMs with linear signal

---

**Require:** EGSSM (Definition 3.5) with linear signal and natural parameters  $s_t$ ,  $t = 0, \dots, n$ , observations  $y_0, \dots, y_n$ , initial values  $\psi^0 = (z^0, (\omega^2)^0)$ , number of samples  $N$ .

- 1: Set  $l = 0$ .
  - 2: **repeat**
  - 3:   Run the FFBS (Algorithm 3) to obtain samples of the signals  $(S^i)_{i=1, \dots, N}$ .
  - 4:   **for**  $t = 0, \dots, n$  **do**
  - 5:     **for**  $i = 1, \dots, p$  **do**
  - 6:       Solve Equation (3.24) for  $z_{t,i}^{l+1}$  and  $(\omega^2)_{t,i}^{l+1}$  ▷ Using  $\hat{P}_N$
  - 7:     **end for**
  - 8:   **end for**
  - 9:   Set  $\psi^{l+1} = (z^{l+1} \oslash (\omega^2)^{l+1}, -1 \oslash (\omega^2)^{l+1})$ .
  - 10:   Set  $l = l + 1$ .
  - 11: **until**  $\psi^l$  has converged.
- 

For the CE-method, using the GLSSM-approach turns out to be difficult numerically. For a high-level argument of why this is true, let us ignore the Markov structure of the model for the moment. As the CE-method matches moments of the target and proposal, applying it to fit model (3.21) amounts to matching the moments of  $\mathbf{G}_\psi$  to those of the target posterior  $\mathcal{L}(X|Y=y)$  in the SSM. Unfortunately, the covariance of  $\mathbf{G}_\psi$  is given by  $(\Sigma^{-1} + B^T \Omega^{-1} B)^{-1}$ , where  $\Sigma$  is the covariance of all states,  $B = \text{block-diag}(B_0, \dots, B_n)$  and  $\Omega = \text{block-diag}(\Omega_0, \dots, \Omega_n)$ . Choosing the diagonal matrix  $\Omega$  such that the covariance of  $\mathbf{G}_\psi$  matches this expression is numerically expensive: we either need to invert the large (dimension  $(n+1)m \times (n+1)m$ ) covariance matrix, or solve numerically for the  $(n+1)p$  parameters. The problem at hand is that we cannot decouple this into  $(n+1)$  equations of dimension  $p$  (or even  $(n+1)p$  equations) as we did for EIS, because all entries of  $(\Sigma^{-1} + B^T \Omega^{-1} B)^{-1}$  depend on all entries of  $\Omega$ .

To make matters more concrete, the CE-method finds  $\psi = (z, \omega^2)$  such that model (3.21) maximizes the cross entropy with the target  $\mathbf{P}^{X|Y=y}$ . For simplicity, let us assume that  $m = p$ ,  $B$  is the identity and we only observe a single  $y$ . Using Lemma 3.1, we see that when  $X \sim \mathcal{N}(\mu, \Sigma)$ , the conditional distribution of  $X$  given  $Z = z$ ,  $\mathbf{G}_\psi$ , is a Gaussian distribution with mean  $\tilde{\mu} = \mu + \Sigma(\Sigma + \Omega)^{-1}(z - \mu)$  and covariance matrix  $\tilde{\Sigma} = (\Sigma^{-1} + \Omega^{-1})^{-1}$  for  $\Omega = \text{diag}(\omega^2)$ , where  $\omega^2 > 0$ . Assuming that  $\Sigma$  is

---

<sup>7</sup>Actually, we do not require the model to be an EGSSM, but only that the univariate marginals are conditionally independent.



non-singular, we can reparameterize the objective function of the CE-method by  $\tilde{\mu}$ ,

$$\begin{aligned} \max_{z, \omega^2} \int p(s|y) \log g_\psi(s|z) dx &= \max_{\tilde{\mu}, \omega^2} \int p(s|y) \left( -\frac{1}{2}(s - \tilde{\mu})^T \tilde{\Sigma}^{-1} (s - \tilde{\mu}) - \frac{1}{2} \log \det \tilde{\Sigma} \right) dx \\ &= \max_{\tilde{\mu}, \omega^2} -\frac{1}{2}(\gamma - \tilde{\mu})^T \tilde{\Sigma}^{-1}(\gamma - \tilde{\mu}) - \frac{1}{2} \text{trace} \left( \tilde{\Sigma}^{-1} \Gamma \right) - \frac{1}{2} \log \det \tilde{\Sigma}, \end{aligned} \quad (3.25)$$

where  $\gamma = \mathbb{E}(X|Y = y)$  and  $\Gamma = \text{Cov}(X|Y = y)$ . Thus the optimal  $\tilde{\mu}$  is  $\gamma$  and to find the optimal  $\omega^2$  we have to minimize

$$\text{trace} \left( (\Sigma^{-1} + \Omega^{-1}) \Gamma \right) - \log \det (\Sigma^{-1} + \Omega^{-1}).$$

Taking the derivative w.r.t.  $\frac{1}{\omega^2}$ , we see that

$$\Gamma_{i,i} = \left( \left( \Sigma^{-1} + \text{diag} \left( \frac{1}{\omega_1}, \dots, \frac{1}{\omega_p} \right) \right)^{-1} \right)_{i,i} = \left( \Sigma - \Sigma (\Sigma + \Omega)^{-1} \Sigma \right)_{i,i} \quad (3.26)$$

has to hold for all  $i = 1, \dots, (p \times (n+1))$ , i.e. we have to choose  $\omega^2$  such that the posterior marginal variances  $\Gamma_{i,i}$  coincide with the marginal variances of  $\mathbf{G}_\psi$ .

Several problems arise: First of all, Equation (3.26) is not guaranteed to have a solution. For the  $i$ -th unit-vector  $e_i \in \mathbf{R}^p$  we can reformulate Equation (3.26) to

$$\Sigma_{i,i} - \Gamma_{i,i} = e_i^T \Sigma^T (\Sigma + \Omega)^{-1} \Sigma e_i > 0$$

and so we require  $\Gamma_{i,i} < \Sigma_{i,i}$ . While the law of total covariance asserts that

$$\Sigma = \mathbb{E} \text{Cov}(X|Y) + \text{Cov}(\mathbb{E}(X|Y)),$$

it does not guarantee  $\Gamma \prec \Sigma$ , as  $\text{Cov}(X|Y) = \Gamma$  may not hold in the non-Gaussian case.

Second, even if there is an analytical solution  $\Omega$  to Equation (3.26), in the CE-method we replace  $\Gamma_{i,i}$  by the observed marginal variances  $\hat{\Gamma}_{i,i}$  obtained by importance sampling. The variation introduced by simulation can then lead to situations where  $\hat{\Gamma}_{i,i} > \Sigma_{i,i}$ . As an example take  $X \sim \mathcal{N}(0, 1)$ , and  $Y = X + \eta$  for  $\eta \sim \mathcal{N}(0, \omega^2)$ . Then the conditional variance of  $X$  given  $Y = y$  is  $\Gamma = 1 - \frac{1}{1+\omega^2} < 1$ . Given  $N$  i.i.d. samples  $X^1, \dots, X^N$  from this distribution, their empirical variance  $\hat{\Gamma} = \frac{1}{N} \sum_{i=1}^N (X^i - \bar{X})^2$  follows a scaled  $\chi_{N-1}^2$  distribution, i.e.  $\frac{N\hat{\Gamma}}{\Gamma} \sim \chi_{N-1}^2$ . Notice that we use the non-Bessel corrected version of the empirical variance here, as it is the maximum-likelihood estimate.

Then

$$\mathbf{P} \left( \hat{\Gamma} > 1 \right) = \mathbf{P} \left( \frac{N\hat{\Gamma}}{\Gamma} > \frac{N}{\Gamma} \right) = 1 - F_{\chi_{N-1}^2} \left( N \left( 1 + \frac{1}{\omega^2} \right) \right)$$

is the probability that Equation (3.26) has no solution  $\omega^2 \in \mathbf{R}_{\geq 0}$ . Here  $F_{\chi_{N-1}^2}$  is the cumulative distribution function of the  $\chi_{N-1}^2$  distribution. As  $\omega^2$  goes to  $\infty$ , this probability approaches  $1 - F_{\chi_{N-1}^2}(N)$  which, for large  $N$ , is approximately  $1 - F_{\chi_{N-1}^2}(N-1) \approx \frac{1}{2}$ , as  $\chi_{N-1}^2 \approx \mathcal{N}(N-1, 2(N-1))$  (Johnson, Kotz, and Balakrishnan, 1994, Section 18.5). Thus, even in the basic univariate Gaussian setting, for every  $N$  there is an  $\omega^2$  such that the CE-method fails for Equation (3.21) with practically relevant probability.

In higher-dimensional settings, e.g. when applying the CE-method to SSMS, we can expect this phenomenon to occur even more often. In the extreme case of independent marginals, i.e. when  $\Sigma$  is a diagonal matrix, Equation (3.26) reduces to  $(n+1)p$  many decoupled equations, where  $\hat{\Gamma}_{i,i}$ ,  $i = 1, \dots, (n+1)p$  are independent. If all  $q_i = \mathbf{P}(\Gamma_{i,i} > \Sigma_{i,i})$  are identical to  $q \in (0, 1)$ , e.g. because  $\Sigma$  and  $\Omega$  are multiples of the identity, the number of failures follows a  $\text{Binom}((n+1)p, q)$  distribution, so that even small  $q$  may lead to a non-negligible number of failures if the number of observations is high.

Finally, in the multivariate setting, the system (3.26) has no analytical solution. Instead, we have to resort to numerical methods to find a solution  $\Omega$ . Unfortunately, even evaluating the right-hand side of (3.26) requires  $\mathcal{O}(m^3)$  operations, as we have to invert  $\Sigma + \Omega$ . Additionally, we cannot hope to reuse a singular-value, LR, or eigenvalue-decomposition for further evaluations, as  $\Sigma$  and  $\Omega$  are not guaranteed to be jointly diagonalizable. In the SSM context we may use the Kalman-smoother to compute the marginal variances, but have to re-run the smoother for every evaluation, which is expensive.

If we admit noise variance  $\infty$  in the univariate setting, then  $\Gamma > 1$  implies that the CE-method chooses this as the estimate, i.e.  $\mathbf{G}_{\hat{\psi}_{\text{CE}}}$  is  $\mathcal{N}(0, 1)$ , which is equal to the prior. We can interpret this as having a missing observation, which, going back to the SSM context, the Kalman-filter (Algorithm 1) can handle with only simple modifications, see e.g. (Durbin and Koopman, 2012, Section 4.10). However, if there are a lot of failures, the optimally chosen  $\mathbf{G}_{\hat{\psi}_{\text{CE}}}$  will be close to the prior distribution of states  $X$ , and importance sampling is unlikely to be effective. Hence, we turn to another approach that allows us to apply the CE-method to SSMs.

### 3.5.2 The Markov-approach

An alternative family of Gaussian proposals is given by directly modeling a Gaussian Markov process on the states  $X_{:n}$ . Again, this is sensible given the Markov structure of the target. This parametrization is more flexible than using the posterior of a GLSSM with fixed prior as the proposal. This flexibility, however, comes at the cost of requiring a larger number of parameters. Here we propose with  $\mathbf{G}_{\psi}$  where

$$\begin{aligned} \mathbf{G}_{\psi} &= \mathcal{L}(U + v), \\ v &\in \mathbf{R}^{(n+1)m}, \\ U_0 &\sim \mathcal{N}(0, R_0 R_0^T), \\ U_t &= C_t U_{t-1} + R_t \nu_t, \\ C_t &\in \mathbf{R}^{m \times m}, \\ \nu_t &\sim \mathcal{N}(0, I_m), \\ R_t &\in \mathbf{R}^{m \times m} \text{ lower triangular with positive diagonal} \end{aligned} \tag{3.27}$$

for  $t = 1, \dots, n$ , with  $U_0$  and  $\nu_1, \dots, \nu_n$  independent. The number of parameters in

$$\psi = (v, C_1, \dots, C_n, R_0, \dots, R_n)$$

is  $(n+1) \cdot m$  for the mean  $v$ ,  $n \cdot m^2$  for the transition matrices  $C_t$  and  $(n+1) \frac{m(m-1)}{2}$  for the Cholesky roots of innovation covariances, totaling  $\mathcal{O}(n \cdot m^2)$  many parameters. While these are considerably more parameters than for the GLSSM-approach for large state dimension  $m$ , we will see in the later part of this section that finding the optimal parameters for the CE-method can be done analytically, depending only on the posterior first and second moments.

This approach, which we term the **Markov-approach**, was originally proposed by (Richard and Zhang, 2007) for general unnormalized transition kernels as EIS proposals. However, because of its lower number of parameters, one should favor the GLSSM-approach for EIS that operates on the signals, see (Koopman, Lit, and Nguyen, 2019).

To perform importance sampling with  $\mathbf{G}_{\psi}$  in model (3.27) we not only need to simulate from  $\mathbf{G}_{\psi}$  but also evaluate the unnormalized importance sampling weights  $w(x) = \frac{p(x|y)}{g_{\psi}(x)}$ . Simulation from  $\mathbf{G}_{\psi}$  is achieved by a simple recursion. For the weights note that

$$w(x) \propto \frac{p(y|x)p(x)}{g_{\psi}(x)} = \prod_{t=0}^n \frac{p(y_t|x_t)p(x_t|x_{t-1})}{g_{\psi}(x_t|x_{t-1})}, \tag{3.28}$$

where  $p(x_0|x_{-1}) = p(x_0)$  and  $g_{\psi}(x_0|x_{-1}) = g_{\psi}(x_0)$ .

The Markov structure of model (3.27) implies that the precision matrix of  $\mathbf{G}_{\psi}$  is sparse, i.e. it has a block-tridiagonal form. This is a well-known property of the precision matrix of Gaussian random

vectors, as the following two classical lemmas show. We show their proofs here for completeness. For a general treatment, we refer the reader to (Lauritzen, 1996, Chapters 3 and 5).

**Lemma 3.10.** *Let  $(X, Y)$  be jointly Gaussian with distribution  $\mathcal{N}(\mu, \Sigma)$  where*

$$\mu = (\mu_X, \mu_Y)$$

and

$$\Sigma = \begin{pmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{pmatrix},$$

are partitioned according to the dimensions of  $X$  and  $Y$  and  $\Sigma$  is non-singular. If

$$P = \Sigma^{-1} = \begin{pmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{pmatrix}^{-1} = \begin{pmatrix} P_{XX} & P_{XY} \\ P_{YX} & P_{YY} \end{pmatrix}$$

is the precision matrix of  $(X, Y)$ , partitioned as is  $\Sigma$ , then  $\text{Cov}(X|Y) = P_{XX}^{-1}$ .

*Proof.* Without loss of generality, assume that both  $X$  and  $Y$  are centered. The conditional density  $p(x|y)$  is proportional (in  $x$ ) to the joint density  $p(x, y)$  with

$$\log p(x, y) = -\frac{1}{2} \begin{pmatrix} x & y \end{pmatrix} P \begin{pmatrix} x \\ y \end{pmatrix} + C = -\frac{1}{2} (x^T P_{XX} x + 2x^T P_{XY} y) + C',$$

for constants  $C, C'$  that do not depend on  $x$ . As the conditional distribution of  $X$  given  $Y = y$  is Gaussian (by Lemma 3.1), its covariance matrix is  $P_{XX}^{-1}$ .  $\square$

**Lemma 3.11.** *Let  $(X, Y, Z) \sim \mathcal{N}(\mu, \Sigma)$  be jointly Gaussian with non-singular  $\Sigma$ . Then  $X \perp Y|Z$  if, and only if, the sub-matrix of the precision matrix  $P = \Sigma^{-1}$  whose rows correspond to the entries of  $X$  and columns correspond to the entries of  $Y$  is the 0 matrix.*

*Proof.* Partition the conditional covariance matrix into

$$\text{Cov}((X, Y)|Z) = \begin{pmatrix} \Sigma_{XX|Z} & \Sigma_{XY|Z} \\ \Sigma_{YX|Z} & \Sigma_{YY|Z} \end{pmatrix}.$$

As all distributions involved are Gaussian,  $X \perp Y|Z$  is equivalent to  $\text{Cov}((X, Y)|Z)$  being a block-diagonal matrix with blocks  $\Sigma_{XX|Z}$  and  $\Sigma_{YY|Z}$ , which is equivalent to its inverse being a block-diagonal matrix with blocks  $\Sigma_{XX|Z}^{-1}$  and  $\Sigma_{YY|Z}^{-1}$ . Its inverse is, by Lemma 3.10, the sub-matrix of  $P$  whose rows and columns correspond to  $X$  and  $Y$ .  $\square$

Applying Lemma 3.11 to model (3.27), we see that its precision matrix  $P$  is sparse, i.e. it is a block-tri-diagonal matrix, as  $U_t \perp U_s|U_{-t,-s}$  for  $|t - s| > 1$  and  $U_{-t,-s}$  being the vector of all  $U_0, \dots, U_n$  except for  $U_t, U_s$ . Thus, the only entries of  $P$  that are potentially non-zero are those whose row and column correspond to  $(U_t, U_t)$  for  $t = 0, \dots, n$ ,  $(U_t, U_{t-1})$  and  $(U_{t-1}, U_t)$  for  $t = 1, \dots, n$ . Therefore,  $P$  has the following block-tridiagonal structure:

$$P = \begin{pmatrix} P_{0,0} & P_{0,1} & 0 & \cdots & \cdots & 0 & 0 \\ P_{1,0} & P_{1,1} & P_{1,2} & 0 & \cdots & 0 & 0 \\ 0 & P_{2,1} & P_{2,2} & P_{2,3} & \cdots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 & 0 \\ 0 & 0 & 0 & \cdots & P_{n-1,n-2} & P_{n-1,n-1} & P_{n-1,n} \\ 0 & 0 & 0 & \cdots & 0 & P_{n,n-1} & P_{n,n} \end{pmatrix}. \quad (3.29)$$

As the precision matrix is the natural parameter for the multivariate Gaussian exponential family, we see that model (3.27), parameterized by  $(P^{-1}v, P)$  form a natural exponential family and we can apply Theorem 3.6 to obtain a central limit theorem when applying the CE-method for this model.

The sparsity of  $P$  implies that  $P = LL^T$  has a sparse Cholesky root  $L$ , which will make computations efficient. To see that  $L$  is sparse, we apply the following Theorem, slightly adapted to our notation, from the theory of Gaussian-Markov-Random-fields (GMRF), i.e. Gaussian models whose dependency structure is given by a graph, with edges between nodes indicating non-zero entries in the precision matrix.

**Theorem 3.10** ((Gelfand et al., 2010, Theorem 12.14)). *Let  $X = (X_0, \dots, X_n) \in \mathbf{R}^{(n+1)m}$  be a GMRF wrt to the labeled graph  $G$ , with mean  $\mu$  and symmetric positive-definite precision matrix  $P$ . Let  $L$  be the Cholesky factor of  $P$  and define for  $0 \leq t < s \leq n$  the future of  $t$  except  $s$  as*

$$F(t, s) = \{t+1, \dots, s-1, s+1, n\}.$$

Then

$$X_t \perp X_s | X_{F(t,s)} \Leftrightarrow L_{t,s} = 0.$$

In the preceding theorem  $X_{F(t,s)}$  is the vector of all  $X_u$  for  $u \in F(t, s)$  and  $L_{t,s} \in \mathbf{R}^{m \times m}$  is the sub-matrix of  $L$  whose rows correspond to  $X_t$  and columns to  $X_s$ . From Theorem 3.10 we immediately obtain the following:

**Corollary 3.2** (sparsity of  $L$  in model (3.27)). *Let  $U \sim \mathbf{G}_\psi$  as in Equation (3.27),  $P \succ 0$  be the precision matrix of  $\overleftarrow{U} = (U_n, \dots, U_0)$  and  $L$  be the Cholesky root of  $P$ . Then  $L$  is a lower-block-diagonal matrix, with at most  $nm^2 + (n+1)m\frac{m-1}{2}$  non-zero entries:*

$$L = \begin{pmatrix} L_{n,n} & 0 & \cdots & \cdots & \cdots & 0 & 0 \\ L_{n-1,n} & L_{n-1,n-1} & 0 & \cdots & \cdots & 0 & 0 \\ 0 & L_{n-2,n-1} & L_{n-2,n-2} & 0 & \cdots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 & 0 \\ 0 & 0 & 0 & \cdots & L_{1,2} & L_{1,1} & 0 \\ 0 & 0 & 0 & \cdots & 0 & L_{0,1} & L_{0,0} \end{pmatrix}, \quad (3.30)$$

where  $L_{t,t} \in \mathbf{R}^{m \times m}$ ,  $t = 0, \dots, n$  are lower triangular matrices with positive diagonal entries and  $L_{t-1,t} \in \mathbf{R}^{m \times m}$ ,  $t = 1, \dots, n$  are square matrices.

From  $L$  in Corollary 3.2 we obtain an iterative method of sampling from  $\mathbf{G}_\psi$ : If  $v + U \sim \mathbf{G}_\psi$ , then, as  $\text{Cov} U = (LL^T)^{-1} = L^{-T}L^{-1}$ , it holds that  $L^T U \sim \mathcal{N}(0, I)$  follows a standard normal distribution. Thus to simulate from  $\mathbf{G}_\psi$  we may solve

$$L^T U = \overleftarrow{Z}$$

where  $\overleftarrow{Z} = (Z_n, \dots, Z_0) \sim \mathcal{N}(0, I)$ . Using the structure available in  $L$ , we see that this is equivalent to first solving

$$L_{0,0}^T U_0 = Z_0$$

and then recursively solving for  $t = 1, \dots, n$

$$L_{t,t}^T U_t + L_{t-1,t}^T U_{t-1} = Z_{t-1}.$$

Rearranging terms, provided  $L_{t,t}$  is non-singular, we end up with the Markov-process

$$U_t = L_{t,t}^{-T} L_{t-1,t}^T U_{t-1} + L_{t,t}^{-T} Z_t, \quad (3.31)$$

where  $Z_t$  is, by construction, independent of  $U_{t-1}$ . Thus for model (3.27), we obtain

$$\begin{aligned} R_t &= L_{t,t}^{-T} \text{ for } t = 0, \dots, n, \\ C_t &= L_{t,t}^{-T} L_{t-1,t}^T \text{ for } t = 1, \dots, n. \end{aligned} \quad (3.32)$$

Here we see why we chose to use  $\overleftarrow{U}$  in Corollary 3.2: had we applied Theorem 3.10 to  $U$  directly we would have ended up with a Markov process in reverse time.

We now turn our attention to applying the CE-method to model (3.27). Following a similar argument as in the discussion surrounding Equation (3.25), we see that we may match the mean  $v$  to that of  $\mathbf{P}$  and it suffices to choose  $P$ , the precision matrix of  $U$ , such that it minimizes

$$\frac{1}{2} \text{trace} \left( P \hat{\Gamma} \right) - \frac{1}{2} \log \det P \quad (3.33)$$

where  $\hat{\Gamma}$  is the importance sampling estimate of the joint covariance matrix of all states  $X$ . This is equivalent to minimizing

$$\mathcal{D}_{\text{KL}} \left( \mathcal{N}(0, \hat{\Gamma}) \parallel \mathcal{N}(0, P^{-1}) \right).$$

Here  $P$  is restricted to precision matrices that may arise in model (3.27), i.e., by Corollary 3.2,  $P = LL^T$  where  $L$  possess structure as in (3.30). At first glance, this problem seems more involved than solving Equation (3.26): after all, the optimal  $P$  depends on the whole covariance matrix  $\hat{\Gamma}$ . However, it turns out that the sparsity we enforce in  $L$  allows us to compute analytically the optimal  $\hat{L}$  that minimizes Equation (3.33). Additionally, due to the Markov-structure of our proposal,  $\hat{L}$  depends only on the block-tri-diagonal component of  $\hat{\Gamma}$ , i.e. only the covariances  $\text{Cov}(X_t, X_{t-1})$  and  $\text{Cov}(X_0)$  are required. This is sensible - all information about the Markov transitions is encoded in these covariances if we assume that  $X$  is a Gaussian Markov process.

To make this argument rigorous, let us apply the following result (stated in our notation).

**Theorem 3.11** ((Schäfer, Katzfuss, and Owhadi, 2021, Theorem 2.1)). *Let  $\Gamma$  be a positive-definite matrix of size  $n \times n$ . Given a lower-triangular sparsity set  $S \subset \{1, \dots, n\}^2$ , i.e.  $i \geq j$  for all  $(i, j) \in S$ , let*

$$\hat{L} = \underset{L \in S}{\text{argmin}} \mathcal{D}_{\text{KL}} \left( \mathcal{N}(0, \Gamma) \parallel \mathcal{N}(0, (LL^T)^{-1}) \right)$$

*be the Cholesky root of the closest Gaussian (wrt. the KL-divergence) with sparsity  $S = \{A \in \mathbf{R}^{n \times n} : A_{i,j} \neq 0 \Rightarrow (i, j) \in S\}$ .*

*Then the following holds: The nonzero entries of the  $i$ -th column of  $\hat{L}$  are given by*

$$L_{s_i, i} = \frac{\Gamma_{s_i, s_i}^{-1} e_1}{\sqrt{e_1^T \Gamma_{s_i, s_i}^{-1} e_1}}, \quad (3.34)$$

*where  $s_i = \{j : (i, j) \in S\}$ ,  $\Gamma_{s_i, s_i}$  is the restriction of  $\Gamma$  to the set of indices  $s_i$  and  $e_1 \in \mathbf{R}^{|s_i|}$  is the first unit vector.*

Exploiting the Markov structure of our proposals, we immediately obtain the following:

**Corollary 3.3.** *Let  $S$  be the sparsity set of a Gaussian Markov process of the form Equation (3.27), i.e.*

$$S = \left\{ ((t, i), (s, j)) \in (\{0, \dots, n\} \times \{1, \dots, m\})^2 \mid (t = s \text{ and } i \geq j) \text{ or } t = s + 1 \right\},$$

*see also Equation (3.30), and let  $\Gamma$  be a positive definite matrix of size  $((n+1)m) \times (n+1)m$  with blocks*

$$\Gamma_{s,t} = (\Gamma_{(s,i),(t,j)})_{i,j=1,\dots,m}.$$

*Then  $\hat{L}$  in Theorem 3.11 depends only on the block-diagonal entries  $\Gamma_{t,t}$ ,  $t = 0, \dots, n$  and block off-diagonal entries  $\Gamma_{t,t+1}$ ,  $t = 0, \dots, n$ .*

*If, in particular,  $\Gamma$  is the covariance matrix of Gaussian Markov process,  $\hat{L} = \text{chol}(\Gamma^{-1})$ .*

We have thus shown the following: The covariance matrix of the KL-optimal Gaussian Markov process for the positive definite covariance matrix  $\Gamma$  with  $\mathcal{O}(n^2 m^2)$  entries only depends on  $\mathcal{O}(nm^2)$  many entries, the marginal covariances. In particular, if we can find a centered Gaussian Markov process  $(X_t)_{t=0,\dots,n}$  whose marginal covariances fulfill

$$\begin{aligned} \text{Cov}(X_t) &= \Gamma_t & t &= 0, \dots, n \\ \text{Cov}(X_t, X_{t+1}) &= \Gamma_{t,t+1} & t &= 0, \dots, n, \end{aligned}$$

then its law  $\mathcal{L}(X)$  is the one we seek. The following proposition puts all the pieces together.

**Theorem 3.12** (the CE-method for the Markov proposal). *Let  $\mathbf{P}$  be a probability measure on  $\mathbf{R}^{(n+1) \times m}$  with mean  $\mu$  and positive definite covariance matrix  $\Gamma$ , partitioned into blocks*

$$\Gamma_{s,t} = (\Gamma_{(s,i),(t,j)})_{i,j=1,\dots,m}.$$

Let

$$\begin{pmatrix} J_{t,t} & 0 \\ J_{t+1,t} & Z_{t+1,t+1} \end{pmatrix} = \text{chol} \begin{pmatrix} \Gamma_{t,t} & \Gamma_{t,t+1} \\ \Gamma_{t+1,t} & \Gamma_{t+1,t+1} \end{pmatrix}.$$

Then the optimal cross-entropy parameter

$$\psi_{CE} = \operatorname{argmin}_{\psi=(v,C_1,\dots,C_n,R_0,\dots,R_n)} \mathcal{D}_{KL}(\mathbf{P} \parallel \mathbf{G}_\psi)$$

for the Markov proposal  $\mathbf{G}_\psi$  from model (3.27) exists and is unique. The components of  $\psi_{CE}$  are given by

$$\begin{aligned} v &= \mu \\ R_0 &= \text{chol}(\Gamma_{0,0}) \end{aligned}$$

and for  $t = 1, \dots, n$

$$\begin{aligned} C_t &= J_{t+1,t} J_{t,t}^{-1} \\ R_t &= Z_{t+1,t+1} \end{aligned}$$

Thus, given  $\nu$  and  $\Gamma$ ,  $\psi_{CE}$  can be obtained in  $\mathcal{O}(nm^3)$  many operations.

*Proof.* It only remains to show the uniqueness and existence of  $\psi_{CE}$ , as well as its representation. The discussion surrounding Equation (3.33) shows that  $v = \mu$  has to hold, so we may assume that  $\mathbf{P}$  and the proposal are both centered. As  $\Gamma$  is positive definite, so are all of its sub-matrices, and we may apply Corollary 3.3. Therefore, if we can show that there is a unique Gaussian Markovian probability measure whose covariance matrix matches  $\Gamma$  as in that corollary we are done.

Let  $(U_t)_{t=0,\dots,n} \sim \mathbf{G}_{\psi_{CE}}$ . Then

$$\text{Cov}(U_0) = R_0 R_0^T = \Gamma_{0,0},$$

and from the Cholesky decomposition we obtain for  $t = 0, \dots, n-1$

$$\begin{pmatrix} J_{t,t} J_{t,t}^T & J_{t,t} J_{t+1,t}^T \\ J_{t+1,t} J_{t,t}^T & J_{t+1,t} J_{t+1,t}^T + Z_{t+1,t+1} Z_{t+1,t+1}^T \end{pmatrix} = \begin{pmatrix} \Gamma_{t,t} & \Gamma_{t,t+1} \\ \Gamma_{t+1,t} & \Gamma_{t+1,t+1} \end{pmatrix}.$$

As  $Z_{t+1,t+1}$  is a lower triangular matrix with positive diagonal and

$$\Gamma_{t+1,t+1} - \Gamma_{t+1,t} \Gamma_t^{-1} \Gamma_{t,t+1} = Z_{t+1,t+1} Z_{t+1,t+1}^T,$$

it is the Cholesky root of the Schur complement  $\Gamma_{t+1,t+1} - \Gamma_{t+1,t} \Gamma_t^{-1} \Gamma_{t,t+1}$ , which, recalling Lemma 3.1, we can think of as a conditional covariance matrix. Therefore, using induction over  $t = 0, \dots, n-1$ , we obtain

$$\begin{aligned} \text{Cov}(U_{t+1}) &= C_{t+1} \text{Cov}(U_t) C_{t+1}^T + R_{t+1} R_{t+1}^T \\ &= J_{t+1,t} J_{t,t}^{-1} \Gamma_{t,t} J_{t,t}^{-T} J_{t+1,t}^T + \Gamma_{t+1,t+1} - \Gamma_{t+1,t} \Gamma_t^{-1} \Gamma_{t,t+1} \\ &= \Gamma_{t+1,t+1} \end{aligned}$$

and

$$\text{Cov}(U_{t+1}, U_t) = C_t \text{Cov}(U_t) = J_{t+1,t} J_{t,t}^{-1} J_{t,t} J_{t,t}^T = \Gamma_{t+1,t}.$$

This shows the existence. For uniqueness, note that model (3.27) enforces that  $R_t$  is a lower triangular matrix with positive diagonals. As  $R_{t+1} R_{t+1}^T$  is the conditional covariance of  $U_{t+1}$  given  $U_t$  which is, by Lemma 3.1 given by  $\Gamma_{t+1,t+1} - \Gamma_{t+1,t} \Gamma_t^{-1} \Gamma_{t,t+1}$ . Thus the  $R$  matrices are unique as well. As  $\text{Cov}(U_{t+1}, U_t) = C_t \text{Cov}(U_t)$ , we can show that, additionally, also  $C_t$  is unique.  $\square$

When using the CE-method, we do not have access to the mean and covariances necessary to apply this proposition. Instead, we may apply the CE-method to estimate  $\psi$  in model (3.27) by replacing these unknown moments with their importance sampling estimates. Given importance samples  $U^1, \dots, U^N$  for  $\mathcal{L}(X|Y = y)$  and associated auto-normalized weights  $W^1, \dots, W^N$ , we estimate  $v$  by

$$\hat{v} = \sum_{i=1}^N W^i X^i \quad (3.35)$$

and the empirical covariance matrices

$$\begin{aligned} \widehat{\text{Cov}}(X_t, X_{t-1}) &= \sum_{i=1}^N W^i (X_{t:t-1}^i - \hat{v}_{t-1:t}) (X_{t:t-1}^i - \hat{v}_{t-1:t})^T \\ \widehat{\text{Cov}}(X_0) &= \sum_{i=1}^N W^i (X_0^i - \hat{v}_0) (X_0^i - \hat{v}_0)^T \end{aligned} \quad (3.36)$$

These steps are summarized in Algorithm 7.

---

**Algorithm 7** The CE-method for the Markov proposal (3.27)

---

**Require:** EGSSM (Definition 3.5), observations  $Y$ , initial estimate  $\hat{\psi}^0 = (v^0, C^0, R^0)$ , sample size  $N$

- 1: set  $l = 0$
  - 2: **repeat**
  - 3:   sample  $U^1 + v^l, \dots, U^N + v^l \stackrel{\text{i.i.d.}}{\sim} \mathbf{G}_{\hat{\psi}^l}$  with fixed seed ▷ Equation (3.27)
  - 4:   determine auto-normalized weights  $W^1, \dots, W^N$  ▷ Equation (3.28)
  - 5:   estimate  $\hat{v}^{l+1}$  ▷ Equation (3.35)
  - 6:   estimate  $\widehat{\text{Cov}}(U_t, U_{t-1}), t = 1, \dots, n$ , and  $\widehat{\text{Cov}}(U_0)$  ▷ Equation (3.36)
  - 7:   determine  $C^{l+1}$  and  $R^{l+1}$  ▷ Theorem 3.12
  - 8:   set  $\hat{\psi}^{l+1} = (\hat{v}^{l+1}, C^{l+1}, R^{l+1})$
  - 9:   set  $l = l + 1$
  - 10: **until**  $\hat{\psi}^l$  converged
  - 11: **return**  $\hat{\psi}_{\text{CE}} = \hat{\psi}^l$
- 

To run Algorithm 7 we require an initial value for  $\hat{\psi}^0$ . If a suitable  $\hat{\psi}^0$  is not available, we can obtain one from the LA by sampling  $X^1, \dots, X^N$  from the LA and performing steps 5 to 8 from the loop. Alternatively, we could also directly base our initial value on the smoothing distribution of the GLSSM that the LA is based on. The Kalman smoother (Algorithm 2) provides us with the analytically available covariances  $\text{Cov}(X_t, X_{t-1}|Z = z)$  and the marginal covariance  $\text{Cov}(X_0|Z = z)$  can be computed as well.

The convergence criteria in Algorithm 7 is similar to that used for EIS: we stop until the absolute or entry-wise relative difference of  $\hat{\psi}^l$  and  $\hat{\psi}^{l+1}$  is smaller than a predetermined threshold, or a fixed number of iterations has passed. For the matrices involved, we use the Frobenius norm and the Euclidean distance for the mean  $v$ .

In Algorithm 7 we use the standard praxis of CRNs to ensure numerical convergence. This is similar to EIS and the maximum likelihood estimates from Section 3.6.

We give an overview of the time and space complexities of each line in Algorithm 7 in Table 3.1. The total time complexity of a single iteration of Algorithm 7 is  $\mathcal{O}(Nnm^2 + nm^3)$  and its space complexity is  $\mathcal{O}(Nnm + nm^2)$ . Let us elaborate on the complexities of each step:

Algorithm 7 Generate  $N$  i.i.d. samples from model (3.27), where each simulation requires  $\mathcal{O}(n)$  matrix-vector multiplications of dimension  $m$ .

Algorithm 7 To evaluate the weights, Equation (3.28), we have to evaluate for every sample  $\mathcal{O}(n)$ -times the density of a  $m$ -variate Gaussian distribution, while this usually has time-complexity



step	time complexity	space complexity
simulation (Algorithm 7)	$\mathcal{O}(Nnm^2)$	$\mathcal{O}(Nnm)$
weights (Algorithm 7)	$\mathcal{O}(Nnm^2)$	$\mathcal{O}(N)$
estimating $v$ (Algorithm 7)	$\mathcal{O}(Nnm)$	$\mathcal{O}(nm)$
estimating covariances (Algorithm 7)	$\mathcal{O}(Nnm^2)$	$\mathcal{O}(nm)$
determining $C$ and $R$ (Algorithm 7)	$\mathcal{O}(nm^3)$	$\mathcal{O}(nm^2)$

Table 3.1: Time and space complexities of individual steps in Algorithm 7.

$\mathcal{O}(m^3)$ , we have access to the Cholesky root  $R_t$ , so this step has only time-complexity  $\mathcal{O}(m^2)$ . In Equation (3.28) we also need to compute  $p(y_t|x_t)$  and  $p(x_t|x_{t-1})$ . Assuming conditional independence of observations,  $p(y_t|x_t) = \prod_{i=1}^m p(y_t^i|(B_t x_t)^i)$ , evaluating the first term requires only  $\mathcal{O}(m^2)$  operations. For the second term, if we allow pre-computation of the Cholesky roots of innovations off-line (in  $\mathcal{O}(m^3)$  time), this step reduces to  $\mathcal{O}(m^2)$  as well.

Algorithm 7 Calculating the weighted mean  $\hat{v} \in \mathbf{R}^{(n+1)m}$ , Equation (3.35), requires  $\mathcal{O}(Nnm)$  operations.

Algorithm 7 Calculating the weighted covariance matrices, Equation (3.36), requires  $(n+1)$  times multiplying  $N$  many  $m \times 1$  with  $1 \times m$  vectors.

Algorithm 7 For each of the  $\mathcal{O}(n)$  many  $C_t$  and  $R_t$  we have to calculate Cholesky decompositions and invert triangular matrices of dimension  $m$ .

An efficient implementation of Algorithm 7 can improve on the practically relevant computational time. There is no need to calculate the  $C_t$  matrices explicitly, instead we can calculate  $C_t U_{t-1} = J_{t+1,t} J_{t,t}^{-1} U_{t-1}$  efficiently by back-substitution, as  $J_{t,t}$  is a lower triangular matrix.

The main bottleneck for space lies in the  $\mathcal{O}(Nnm)$  simulation part, and we may reduce this by simulating twice from model (3.27) using CRNs, and only storing the samples for a single time step (dimension  $\mathcal{O}(Nm)$ ) in each simulation. In the first pass, we only calculate the weights, and in the second pass, we calculate  $\hat{v}$  and the required covariance matrices. For this, we only need the  $2N$  samples of dimension  $m$  from time  $t$  and  $t+1$ , i.e.  $\mathcal{O}(Nm)$  space. This reduces the total space complexity to  $\mathcal{O}(Nm + nm^2)$ .

We demonstrate these improvements in Algorithm 8. Additionally, we calculate the weights on the log scale for numerical stability.

The advantage of Algorithms 7 and 8 over applying the CE-method to the GLSSM model (3.21) are multiple: First of all, as long as the involved covariance matrices are positive definite, the two algorithms produce valid proposals, i.e. they do not have the degeneracy problem we observed in Section 3.5.1. When matrices are only positive-semi definite, replacing inverses with generalized inverses still yields a valid model. Additionally, determining the optimal parameters  $(v, C, R)$  or  $(v, J, R)$  is numerically stable, involving only inversion of small matrices. Compare this with solving Equation (3.26), where we need to employ a numerical scheme to solve for the diagonal entries of  $\Omega$ .

After having determined  $\hat{\psi}_{\text{CE}}$  for model (3.27), generating  $N$  samples requires only  $\mathcal{O}(Nnm^2)$  operations, whereas sampling from model (3.21) requires  $\mathcal{O}(nm^3 + Nnm^2)$  operations, as we need an initial run of the Kalman filter. Unless  $N < m$ , this difference is negligible, and the case where  $N < m$  is not really of interest, as we would expect importance sampling to require a much larger number of samples, i.e.  $N \gg m$ .

However, the two algorithms presented in this section also come with some drawbacks, especially if the dimension  $m$  of states is large. This affects the algorithms in multiple ways: when  $m$  is large, computation of the Cholesky decomposition in Theorem 3.12 becomes more time-intensive. Additionally, the dimension of the parameter  $\psi$  increases quadratically in  $m$ , so we expect convergence to be slower, requiring a larger sample size  $N$  to find the optimal  $\hat{\psi}_{\text{CE}}$ . For an empirical study in this direction, see Section 3.7.



---

**Algorithm 8** Time and space improved version of Algorithm 7. Instructions involving the free index  $i$  are to be performed for all  $i = 1, \dots, N$  samples. For simplicity of notation we let  $R^l = (R_0^l, \dots, R_n^l)$  and  $J^l = (J_{0,0}^l, J_{1,0}^l, \dots, J_{n-1,n-1}^l, J_{n,n-1}^l)$  for  $l \in \mathbf{N}_0$ .

---

**Require:** EGSSM (Definition 3.5), observations  $Y$ , initial estimate  $\hat{\psi}^0 = (v^0, R^0, J^0)$ , sample size  $N$

```

1: set  $l = 0$ 
2: repeat
3:   simulate  $\nu_0^1, \dots, \nu_0^N \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I)$ 
4:   set  $U_0^i = R_0^l \nu_0^i$ 
5:   set  $X_0^i = v_0^l + U_0^i$ 
6:   set  $\log w^i = \log p(y_0 | X_0^i) + \log p(X_0^i) + \frac{1}{2} \|\nu_0^i\|^2$   $\triangleright \log g(X_0^i) = -\frac{1}{2} \|\nu_0^i\|_2^2 + C$ 
7:   store current RNG state
8:   for  $t \leftarrow 1, \dots, n$  do
9:     simulate  $\nu_t^1, \dots, \nu_t^N \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I)$ 
10:    set  $U_t^i = (J_{t+1,t})^T (J_{t,t})^{-1} U_{t-1}^i + R_t^l \nu_t^i$   $\triangleright$  backsubstitution
11:    set  $X_t^i = v_t^l + U_t^i$ 
12:    set  $\log w^i = \log w^i + \log p(y_t | X_t^i) + \log p(X_t^i | X_{t-1}^i) + \frac{1}{2} \|\nu_t^i\|^2$ 
13:  end for
14:  set  $\log w^i = \log w^i - \max_{i=1, \dots, N} \log w^i$   $\triangleright$  ensure  $\log w^i \leq 0$ 
15:  set  $w^i = \exp(\log w^i)$ 
16:  set  $W^i = \frac{w^i}{\sum_{i=1}^N w^i}$   $\triangleright$  auto-normalized weights
17:  set  $v_0^{l+1} = \sum_{i=1}^N W^i X_0^i$ 
18:  restore RNG state
19:  for  $t \leftarrow 1, \dots, n$  do
20:    simulate  $\nu_t^1, \dots, \nu_t^N \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I)$ 
21:    set  $U_t^i = (J_{t+1,t})^T (J_{t,t})^{-1} U_{t-1}^i + R_t^l \nu_t^i$   $\triangleright$  backsubstitution
22:    set  $X_t^i = v_t^l + U_t^i$ 
23:    calculate  $\hat{v}_t^{l+1}$   $\triangleright$  Equation (3.35)
24:    calculate covariances  $\triangleright$  Equation (3.35)
25:  end for
26:  set  $\hat{\psi}^{l+1} = (\hat{v}^{l+1}, \hat{R}^{l+1}, \hat{J}^{l+1})$ 
27:  set  $l = l + 1$ 
28: until  $\hat{\psi}^l$  converged
29: return  $\hat{\psi}_{\text{CE}} = \hat{\psi}^l$ 

```

---

### 3.6 Inference in PGSSMs

Once we have chosen a suitable PGSSM to model the observations  $(y_t)_{t=0,\dots,n}$ , we are interested in statistical inferences. This is a two-part procedure: first we must estimate the unknown hyperparameters  $\theta$ , which we will do by maximum likelihood estimation. Then we have to obtain a description of the conditional distributions of interest, e.g. the conditional distribution of states given observations or the conditional distribution of future, yet unavailable, observations.

#### 3.6.1 Maximum likelihood estimation

Until now, we have assumed that the SSM under consideration is completely known, i.e. we have access to the true transition and observation kernels. For the models considered in this thesis (Chapter 4), this is unrealistic, as they are not based on concrete physical processes but are rather statistical approximations of the true underlying dynamics. The transition densities of, e.g., Equation (3.4) will depend on the covariance matrix of innovations, of which we have no a priori knowledge and for negative binomially distributed observations the overdispersion parameter  $r$  will be unknown. Let us denote by  $\theta \in \mathbf{R}^l$  the vector of these hyperparameters. To make this dependence explicit, we will introduce subscripts  $\theta$  where appropriate, i.e.  $\mathbf{P}_\theta$  is a target distribution that additionally depends on  $\theta$ ,  $p_\theta$  its density et cetera. This section is loosely based on (Durbin and Koopman, 2012, Chapter 7 & 11) and (Chopin and Papaspiliopoulos, 2020, Chapter 14).

To determine a suitable value of  $\theta$ , multiple options are available. Here, we opt for a frequentist approach, using maximum likelihood estimation to determine an optimal  $\hat{\theta}$ . Therefore, given observations  $y \in \mathbf{R}^{(n+1) \times p}$ ,  $\hat{\theta}$  maximizes the likelihood  $p_\theta(y)$  and can be obtained as the global maximum of the following optimization problem:

$$\max_{\theta \in \Theta} p_\theta(y).$$

For numerical stability, we should maximize the log-likelihood instead, i.e. solve

$$\max_{\theta \in \Theta} \log p_\theta(y). \quad (3.37)$$

Here  $\Theta \subseteq \mathbf{R}^l$  is the parameter space. To solve this optimization problem using gradient ascent algorithms, we need access to both the likelihood and its derivatives. Thus, in the following, we will assume that  $\theta \mapsto \log p_\theta(y)$  is sufficiently smooth, to apply these methods, i.e. it has continuous derivatives of second order.

While the Kalman-filter (Algorithm 1) allows analytical computation of this likelihood GLSSMs, in general SSMs it is numerically intractable. The reason for this is that

$$p_\theta(y) = \int p_\theta(x, y) d\mu(x)$$

is a high-dimensional integral, which is hard to evaluate numerically. Instead, we will use importance sampling to estimate the likelihood. For this, let us regard  $p_\theta(x, y)$  as an unnormalized density in  $x$ . The missing integration constant is then just  $p_\theta(y)$  and the normalized density is  $p_\theta(x|y)$ . If  $\mathbf{G} \gg \mathbf{P}_\theta$  is a proposal distribution whose density  $g$  with respect to  $\mu$  we can evaluate analytically, i.e. not only up to a constant, we see that for the unnormalized weights  $\tilde{w}_\theta(x) = \frac{p_\theta(x, y)}{g(x)}$ , that  $p_\theta(y) = \mathbf{G}[\tilde{w}_\theta]$ . Thus we may estimate the likelihood by

$$\widehat{p_\theta(y)} = \frac{1}{N} \sum_{i=1}^N \tilde{w}_\theta(X^i)$$

for  $X^1, \dots, X^N \stackrel{\text{i.i.d.}}{\sim} \mathbf{G}$  and  $N \in \mathbf{N}$ . To evaluate the gradient, notice that as  $\nabla_\theta p_\theta(x, y) = p_\theta(x, y) \nabla_\theta \log p_\theta(x, y)$ , we have, provided we can exchange integration and differentiation,

$$\begin{aligned} \nabla_\theta p_\theta(y) &= \nabla_\theta \int p_\theta(x, y) d\mu(x) = \int p_\theta(x, y) \nabla_\theta \log p_\theta(x, y) d\mu(x) \\ &= \mathbf{G}[\tilde{w}_\theta \nabla_\theta \log p_\theta(x, y)], \end{aligned}$$

and so we may estimate the gradient by

$$\widehat{\nabla_{\theta} p_{\theta}(y)} = \frac{1}{N} \sum_{i=1}^N \tilde{w}_{\theta}(X^i) \nabla_{\theta} \log p_{\theta}(X^i, y)$$

Similarly, we can estimate the log-likelihood by Plug-In

$$\widehat{\log p_{\theta}(y)} = \log \left( \frac{1}{N} \sum_{i=1}^N \tilde{w}_{\theta}(X^i) \right) \quad (3.38)$$

and its gradient, using the fact that the gradient of  $\log f$  for  $f : \mathbf{R}^l \rightarrow \mathbf{R}$  is  $\frac{1}{f} \nabla_{\theta} f$ , by

$$\begin{aligned} \widehat{\nabla_{\theta} \log p_{\theta}(y)} &= \left( \frac{1}{N} \sum_{i=1}^N \tilde{w}_{\theta}(X^i) \right)^{-1} \left( \frac{1}{N} \sum_{i=1}^N \tilde{w}_{\theta}(X^i) \nabla_{\theta} \log p_{\theta}(X^i, y) \right) \\ &= \sum_{i=1}^N W_{\theta}^i \nabla_{\theta} \log p_{\theta}(X^i, y) \end{aligned}$$

where  $W_{\theta}^i = \frac{\tilde{w}_{\theta}(X^i)}{\sum_{i=1}^N \tilde{w}_{\theta}(X^i)}$  are the auto-normalized weights. Note that, by Jensen's inequality, these estimates are biased.

To solve the optimization problem (3.37) we will again employ CRNs. If the densities involved are twice differentiable, this device ensures that the random objective function  $\theta \mapsto \sum_{i=1}^N \tilde{w}_{\theta}(X^i)$  is twice differentiable, and so we can indeed apply gradient ascent to find a local maximum. This is an advantage of performing global importance sampling over SMC, i.e. particle filter, methods. To avoid collapse to a single particle, SMC methods perform intermediate resampling steps, which make the objective function discontinuous. While particle smoothing methods can mitigate this problem, they are more expensive than standard SMC and, as the importance sampling estimates of the log-likelihood and its gradient are biased, the usual requirements for stochastic approximation methods are not fulfilled. For a more thorough discussion of the challenges maximum likelihood estimation with SMC methods faces, we recommend (Chopin and Papaspiliopoulos, 2020, Chapter 14).

While MLEs have a strong frequentist foundation, let us stress that, for the models that we investigate in Chapter 4, the frequentist properties of the estimates are not of interest. The reason for this is that a frequentist interpretation requires us to imagine, at least hypothetically, an infinite repetition of the data-generating process. For the data at hand, such repetition is nonsensical: the pandemic is a „one-off“ event that will not be replicated under even approximately similar circumstances. Therefore, we will choose to view the estimation procedure more as a hyper-parameter tuning step, rather than true frequentist inference. While we can compute asymptotic confidence intervals for  $\hat{\theta}$ , see, e.g., (Durbin and Koopman, 2012, Chapter 11.6), (Chopin and Papaspiliopoulos, 2020, Chapter 14.8), these are not of practical interest for similar reasons.

As an alternative to modeling  $\theta$  as fixed, but unknown, and performing maximum-likelihood estimation to obtain  $\hat{\theta}$ , one might also model  $\theta$  as random with prior density  $p(\theta)$ , such that the full model becomes  $p(x, y, \theta) = p(x, y|\theta)p(\theta)$ . In this setup, sometimes called the Bayesian treatment of SSMs (Durbin and Koopman, 2012, Section 13.1), the main interest still lies in the posterior density  $p(x, \theta|y)$ , which, depending on the model at hand, can drastically increase the difficulty of the problem: even if  $p(x, y|\theta)$  is an analytically tractable model such as a GLSSM, unless the prior is chosen to be conjugate, one has to resort to, e.g., MCMC-methods.

By the structure of the model, Equation (3.2), the log density and its gradient can be computed efficiently by

$$\begin{aligned} \log p_{\theta}(x, y) &= \log p_{\theta}(x_0) + \sum_{t=1}^n \log p_{\theta}(x_t|x_{t-1}) + \log p_{\theta}(y_t|x_t, y_{t-1}) \\ \nabla_{\theta} \log p_{\theta}(x, y) &= \nabla_{\theta} \log p_{\theta}(x_0) + \sum_{t=1}^n \nabla_{\theta} \log p_{\theta}(x_t|x_{t-1}) + \nabla_{\theta} \log p_{\theta}(y_t|x_t, y_{t-1}), \end{aligned}$$

respectively.

Similarly, when proposing with a GLSSM or Markov-proposal for a PGSSM, the weights have similar structure, see Equations (3.23) and (3.28), which makes calculation of  $\tilde{w}$  efficient.

For the remainder of this section, let us consider the GLSSM-proposal obtained by the LA or EIS for a PGSSM with linear signal, as this is the main setting of Chapter 4. For this we obtain

$$\tilde{w}_\theta(x) = \tilde{w}_\theta(s)g(z)\frac{p_\theta(y|s)}{g(z|s)} = g(z)\prod_{t=0}^n\frac{p_\theta(y_t|s_t)}{g(z_t|s_t)},$$

where  $s_t = B_tx_t$ ,  $t = 0, \dots, n$ , is the signal, and so the log-likelihood is given by

$$\log p_\theta(y) = \log g_\theta(z) + \log \mathbb{E}(w_\theta(S)|Y = y) \quad (3.39)$$

and can be estimated by

$$\widehat{\log p_\theta(y)} = \log g_\theta(z) + \log \left( \frac{1}{N} \sum_{i=1}^N \prod_{t=0}^n \frac{p_\theta(y_t|S_t^i)}{g(z_t|S_t^i)} \right). \quad (3.40)$$

Notice that  $\log g_\theta(z)$  is the likelihood in a GLSSM, which can be computed efficiently by the standard Kalman filter (Algorithm 1). As in the GLSSM-approach we propose with an GLSSM whose state density  $g(x)$  and observation matrices  $B_t$ ,  $t = 0, \dots, n$  are equal to those of the target, the log-likelihood  $\log g_\theta(z)$  also depends on  $\theta$ . The estimated gradient of the log-likelihood is

$$\widehat{\nabla_\theta \log p_\theta(y)} = \nabla_\theta \log g_\theta(z) + \sum_{i=1}^N W_\theta^i \sum_{t=0}^n \nabla_\theta \log p_\theta(y_t|S_t^i).$$

The gradient of the GLSSM log-likelihood can be obtained either numerically or analytically by employing the Kalman filter and smoother (Koopman and Shephard, 1992), however, numerical evaluation may be faster if the dimension of  $\theta$  is small compared to the length of the time series, as evaluating the likelihood only requires a single application of the Kalman filter.

As the observation densities  $g(z_t|s_t)$  do not depend on  $\theta$ , their derivatives do not appear in the above estimate. However, when using EIS to determine an optimal proposal, the parameter  $\psi = (z, \omega)$  implicitly depends on  $\theta$ . Accounting for this yields the gradient

$$\widehat{\nabla_\theta \log p_\theta(y)} = \nabla_\theta \log g_\theta(z) + \sum_{i=1}^N W_\theta^i \left( \sum_{t=0}^n \nabla_\theta \log p_\theta(y_t|S_t^i) - \nabla_\theta \log g_\theta(z_t|S_t^i) \right),$$

as  $\nabla_\theta \frac{1}{g_\theta(z|s)} = -\frac{1}{g_\theta(z|s)} \nabla_\theta \log g_\theta(z|s)$ . The computation of this additional term is much more involved, as the parameters  $z, \Omega$  are found through an iterative numerical scheme. Instead, we favor numerical differentiation of the whole procedure to evaluate the likelihood at  $\theta$ , including the method of finding an optimal importance sampling scheme.

As an alternative one may try keeping proposal  $\mathbf{G}$  fixed, which would avoid calculation of the involved derivatives in the previous equation. However, this makes the calculation of weights more involved, as then  $p_\theta(x) \neq g(x)$ . Additionally, we would expect the target  $p_\theta(x|y)$  to be quite sensitive to small changes in  $\theta$ , as  $\theta$  will likely contain parameters related to the covariance structure of model, leading to fast degeneration of weights. Nevertheless, a combination of analytical and numerical gradient descent steps may improve the performance of the optimization procedure, but is beyond the scope of this thesis.

As a single evaluation of the log-likelihood can become very expensive we want our procedure to be as efficient as possible. To this end, (Durbin and Koopman, 1997) provides several improvements to the basic algorithm if the model is a PGSSM with a linear signal. Their contributions consist of a bias correction for the log-likelihood, the use of antithetic and control variables to reduce Monte-Carlo error for importance sampling and a deterministic initialization procedure. Let us briefly summarize these ideas, adapted to our notation. As the computational gains for control

variates in the presence of antithetic variables seem to be limited, we do not give the same level of detail here, for an in-depth analysis, we refer the reader to the source.

For bias reduction, a second-order Taylor series expansion shows that for  $\tilde{w} = \frac{1}{N} \sum_{i=1}^N \tilde{w}(X^i)$ ,

$$\begin{aligned} \mathbb{E}(\log \tilde{w}) - \log \mathbf{G}\tilde{w} &= \mathbb{E} \log \left( 1 + \frac{\tilde{w} - \mathbf{G}\tilde{w}}{\mathbf{G}\tilde{w}} \right) \\ &= \frac{\tilde{w} - \mathbf{G}\tilde{w}}{\mathbf{G}\tilde{w}} - \frac{1}{2} \left( \frac{\tilde{w} - \mathbf{G}\tilde{w}}{\mathbf{G}\tilde{w}} \right)^2 + \mathcal{O}_p(N^{-\frac{3}{2}}), \end{aligned}$$

provided  $\tilde{w} \in L^3(\mathbf{G})$ . Thus, estimating the second order term by  $-\frac{\hat{\sigma}^2}{2N\tilde{w}}$ , where  $\hat{\sigma}^2$  is the empirical variance of the unnormalized weights, we can perform a bias reduction by estimating

$$\widehat{\log p_\theta(y)} = \log(\tilde{w}) + \log g_\theta(z) + \frac{\hat{\sigma}^2}{2N\tilde{w}}. \quad (3.41)$$

The second improvement of (Durbin and Koopman, 1997) is the use of antithetic variables and control variates, a device to reduce Monte-Carlo variance. The main idea of an antithetic variable is to construct for each sample  $X^i$ ,  $i = 1, \dots, N$ , another sample  $\tilde{X}^i$  that has the same distribution as  $X^i$ , but is negatively correlated with  $X^i$ . This has two effects: first of all, we increase the number of samples used for importance sampling and second, as the new samples are negatively correlated with the old samples, the Monte-Carlo variance is reduced. The computation of these samples is usually much faster than creating new samples, which requires the use of the expensive FFBS or simulation smoother algorithms.

**Definition 3.6** (antithetic variable). Let  $X, \tilde{X} \in \mathbf{R}^k$  be two random variables with the same distribution,  $\mathcal{L}(X) = \mathcal{L}(\tilde{X})$  and  $f : \mathbf{R}^k \rightarrow \mathbf{R}$ . Then  $\tilde{X}$  is called an antithetic variable of  $X$  for  $f$ , if  $\text{Cov}(f(\tilde{X}), f(X)) < 0$ . If  $k = 1$  and  $f$  is the identity, we just say that  $\tilde{X}$  is an antithetic variable of  $X$ .

(Durbin and Koopman, 1997) introduce two antithetic variables: balanced for location and balanced for scale, both of which are tailored to the multivariate normal distribution.

**Definition 3.7** (antithetic variable balanced for location and scale, (Durbin and Koopman, 1997)). Let  $X \sim \mathcal{N}(\mu, \Sigma)$  for  $\mu \in \mathbf{R}^k$  and  $\Sigma \in \mathbf{R}^{k \times k}$  positive definite. We call

$$\tilde{X} = \mu + (\mu - X) \quad (3.42)$$

the (entry-wise) antithetic balanced for location. If  $L \in \mathbf{R}^{k \times k}$  is a Cholesky root of  $\Sigma$  and

$$X = \mu + L\varepsilon$$

with  $\varepsilon \sim \mathcal{N}(0, I)$ , let  $c = \varepsilon^T \varepsilon \sim \chi_k^2$  and  $c' = F_{\chi_k^2}^{-1}(1 - F_{\chi_k^2}(\sqrt{c}))$ . We call

$$\tilde{X} = \mu + \sqrt{\frac{c'}{c}} (X - \mu) \quad (3.43)$$

the antithetic balanced for scale.

**Lemma 3.12.** In the above definition,  $\tilde{X}_i$  is an antithetic variable of  $X$  for the coordinate functions  $f_i : \mathbf{R}^k \rightarrow \mathbf{R}$ ,  $f_i(x) = x_i$ ,  $i = 1, \dots, k$ . Furthermore,  $\tilde{c}$  is an antithetic variable of  $c$ .

*Proof.* It is easy to see that  $\tilde{X}$  has the same distribution as  $X$ . Furthermore

$$\text{Cov}(f_i(X), f_i(\tilde{X})) = \text{Cov}(2\mu_i - X_i, X_i) = -\Sigma_{i,i} < 0.$$

For  $c$  and  $\tilde{c}$ , let  $U = F_{\chi_k^2}(c)$ , then  $U \sim \text{Unif}(0, 1)$  and  $\tilde{U} = 1 - U = F_{\chi_k^2}(\tilde{c})$ . As  $\tilde{U} \sim \text{Unif}(0, 1)$  as well,  $\mathcal{L}(c) = \mathcal{L}(\tilde{c})$ . In (Whitt, 1976, Lemma 2.3) it is shown that for any pair of real-valued random variables  $(Y, W)$  with CDF  $H$  and marginal CDFs  $F, G$ , it holds

$$\text{Cov}(Y, W) = \int_{\mathbf{R}^2} H(y, w) - F(y)G(w) \, dy \, dw,$$

and, furthermore, by (Whitt, 1976, Theorem 2.1 and Lemma 2.4) that the joint CDF of  $(c, \tilde{c})$  is  $(y, w) \mapsto \max\{0, F(y) + G(w) - 1\}$ , where  $F$  is the CDF of  $c$  and  $G$  the CDF of  $\tilde{c}$ . As

$$a + b - 1 = ab + a(1 - b) + b - 1 = ab - (1 - a)(1 - b) < ab$$

for all  $a, b \in (0, 1)$ , we have

$$\begin{aligned} \text{Cov}(c, \tilde{c}) &= \int_{\mathbf{R}^2} H(y, w) - F(y)G(w) \, dy \, dw \\ &= \int_{\mathbf{R}^2} \max\{0, F(y) + G(w) - 1\} - F(y)G(w) \, dy \, dw < 0. \end{aligned}$$

□

Let us mention that, by the properties of the standard multivariate normal distribution,  $c = \|u\|$  and  $\frac{u}{\|u\|}$  are independent. Writing

$$X = \mu + \|u\|L\frac{u}{\|u\|} = \mu + \|u\|\frac{X - \mu}{\sqrt{c}},$$

we see that

$$\check{X} = \mu + \sqrt{\tilde{c}}\frac{X - \mu}{\sqrt{c}}$$

has the same distribution as  $X$ , as  $\tilde{c} \sim \mathcal{L}(\|u\|^2)$  and is independent of  $\frac{X - \mu}{\sqrt{c}}$ .

Given a GLSSM-proposal and samples  $X^1, \dots, X^N$  from it, we can calculate these antithetic variables efficiently: for the location balanced antithetic we can calculate the mean using the Kalman-smoother and for the scale balanced antithetic we can calculate  $c$  and  $c'$  using the inverse CDF of the  $\chi_k^2$  distribution and the standard normal samples used to sample  $X^i$  in the first place, for which fast implementations are readily available. Incidental, we obtain a third antithetic,

$$\check{\check{X}} = \mu - \sqrt{\frac{c'}{c}}(X - \mu) \quad (3.44)$$

for free. We can then estimate the log-likelihood in Equation (3.41) by replacing each occurrence of  $\tilde{w}_\theta(X^i)$  by

$$\frac{1}{4} \left( \tilde{w}_\theta(X^i) + \tilde{w}_\theta(\tilde{X}^i) + \tilde{w}_\theta(\check{X}^i) + \tilde{w}_\theta(\check{\check{X}}^i) \right). \quad (3.45)$$

As the procedure to evaluate the likelihood by importance sampling becomes expensive as the dimension of the model increases, (Durbin and Koopman, 1997) recommend finding an initial value  $\hat{\theta}_0$  by maximizing a deterministic version of Equation (3.38). For this, denote by  $s^*$  the mode of the linear signal, conditional on the pseudo-observations  $z$ . As  $S$  follows a multivariate Gaussian,  $s^*$  is also the mean which can be computed efficiently by the Kalman or signal-smoother. Approximating the conditional expectation in Equation (3.39) by  $w_\theta(s^*)$  then yields

$$\log p_\theta(y) \approx \log g_\theta(z) + \log w_\theta(s^*), \quad (3.46)$$

which can be evaluated without simulation by the LA. A better approximation can be obtained by performing a fourth-order Taylor expansion of  $s \mapsto w_\theta(s)$  around the mode  $s^*$ , which yields

$$\log p_\theta(y) \approx \log g_\theta(z) + \log w_\theta(s^*) + \log \left( 1 + \frac{1}{8} \sum_{t=1}^n \sum_{j=1}^m l_{t,j}^{(4)}(s^*) v_{t,j}^2 \right), \quad (3.47)$$

where  $l^{(4)}$  is the fourth derivative of the log-weights  $s \mapsto \log w_\theta(s)$  and  $v_{t,j}$  is the conditional variance  $\text{Var}(S_{t,j}|Z = z)$  in the proposal. Again, we refer the interested reader to the source for the details.

---

**Algorithm 9** Maximum likelihood estimation in a PGSSM with linear signal using EIS.

---

**Require:** parameterized PGSSM with linear signal, initial  $\theta^0 \in \Theta$ , observations  $y \in \mathbf{R}^{(n+1)p}$ , number of samples  $N$

```

1: function APPROX_LOGLIK( $\theta$ )
2:   obtain LA of the PGSSM for  $\theta$                                 ▷ Algorithm 5
3:   obtain mode  $s^*$  and conditional variances  $v_{t,j}$  from the LA    ▷ Algorithms 1 and 2
4:   return approximate log-likelihood                             ▷ Equation (3.46) or Equation (3.47)
5: end function

6: function ESTIMATE_LOGLIK( $\theta$ )
7:   obtain LA of the PGSSM for  $\theta$                                 ▷ Algorithm 5
8:   obtain EIS proposal  $\mathbf{G}_{(z,\Omega)}$                              ▷ Algorithm 6, LA as initial values
9:   sample  $N$  signals  $S^i$  from  $S|Z = z$  in EIS                     ▷ Algorithm 3 or signal smoother
10:  obtain mode  $s^*$  in EIS proposal                                ▷ Algorithm 2 or signal smoother
11:  calculate antithetic variables  $\tilde{S}^i, \check{S}^i, \breve{S}^i$            ▷ Equations (3.42) to (3.44)
12:  set  $\tilde{w}_\theta^i = \frac{1}{4} \left( \tilde{w}_\theta(X^i) + \tilde{w}_\theta(\tilde{X}^i) + \tilde{w}_\theta(\check{X}^i) + \tilde{w}_\theta(\breve{X}^i) \right)$  Equation (3.45)
13:  set  $\tilde{w}_\cdot = \frac{1}{N} \sum_{i=1}^N \tilde{w}_\theta^i$ 
14:  set  $\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (\tilde{w}_\theta^i - \tilde{w}_\cdot)^2$ 
15:  calculate  $\log g_\theta(z)$                                          ▷ Algorithm 1
16:  return  $\log p_\theta(y)$                                          ▷ Equation (3.41)
17: end function

18: maximize APPROX_LOGLIK with initial value  $\theta^0$               ▷ numerically
19: set  $\theta^0$  to optimal value
20: maximize ESTIMATE_LOGLIK with initial value  $\theta^0$  and CRNs    ▷ numerically
21: set  $\hat{\theta}$  to optimal value
22: return  $\hat{\theta}$ 

```

---

The resulting procedure to find the MLE  $\hat{\theta}$  in a PGSSM with linear signal is summarized in Section 3.6.1. Notice that we use CRNs to ensure numerical convergence. The numerical optimization can be performed using any standard solver such as the BFGS algorithm (Nocedal and Wright, 2006, Chapter 6.1). We cannot give guarantees that this procedure produces the true MLE, i.e. finds the global maximizer. However, as we have discussed earlier, we are not interested in frequentist properties of  $\hat{\theta}$  but see the estimation procedure as a hyperparameter tuning step. Thus, a local maximum may well be sufficient. Nevertheless, checking different starting points and random number seeds should be used to get as close as possible to the global maximum.

Notice that our discussion implies that we cannot reuse a GLSSM proposal used for  $\theta$  at another  $\theta'$ , as  $p_{\theta'}(x) \neq g_{\theta}(x)$ . While we can still calculate the weights using the general Equation (3.38), we presume that the old proposal is not a good choice for the new target. The reason for this is that  $\theta$  will usually contain parameters related to the covariance structure of the innovations and observations, and these parameters usually affect many, if not all states or observations. For example, it is common to model states that perform a random walk with common innovation variance  $\sigma^2$  as an element of  $\theta$ . As the distributions lie in a high-dimensional space, slight misspecification of the covariance structure will drastically deteriorate the performance of importance sampling.

If computations are so involved that we want to avoid running the optimal importance sampling scheme as much as possible, one could try, if the model under investigation allows for it, to split  $\theta$  into  $(\theta_x, \theta_y)$  where  $\theta_x$  only affects the state transitions and  $\theta_y$  only affects the observation densities. Then a coordinate ascent scheme could be employed, where the update step for  $\theta_y$  can reuse the proposal, provided that  $\theta_y$  does not change too much and the observation density  $p_{\theta}(y|x)$  is not too sensitive to changes in  $\theta_y$ , which should imply that the proposal is still close enough to give good importance sampling performance. Then numerical differentiation is only required to update  $\theta_x$ .

### 3.6.2 Posterior inference

Once we have chosen  $\theta$  (by maximum likelihood estimation), and thus a concrete PGSSM with which to perform statistical inference, we are interested in, e.g., the conditional distribution of states  $X$  given observations  $Y$ , or functionals thereof. Here we will assume, for computational reasons, that the PGSSM has a linear signal, otherwise the same arguments can be applied to the states directly as well, at the expense of higher computation cost.

At our disposal we will have, after obtaining a GLSSM proposal using the EIS method, signal samples  $S^i \in \mathbf{R}^{(n+1) \times p}$ ,  $i = 1, \dots, N$  and associated auto-normalized weights  $W^i$ . Let  $\mathfrak{X} \in \mathbf{R}$  be a univariate random variable which is conditionally independent of the observations  $Y$  given the signal  $S$ , i.e.  $\mathfrak{X} \perp Y|S$ , whose conditional expectation and variance given  $Y$  exist, as well as a regular version of this conditional distribution.  $\mathfrak{X}$  can be a marginal of  $X$ , a scalar function of  $X$ , a future or missing observation, or function thereof. We will assume that we can sample from  $\mathfrak{X}|S$ . This is reasonable for all scenarios we are interested in: states and signals are jointly Gaussian, so samples can be obtained using the FFBS (Algorithm 3), and the distribution of missing or future observations, conditional on states, is tractable in the models we consider. The following paragraphs are based on (Durbin and Koopman, 2012, Section 11.5), but stated in more general terms using  $\mathfrak{X}$ .

We are then interested in estimating several quantities: the conditional expectation  $\mathbf{E}(\mathfrak{X}|Y)$ , the conditional variance  $\text{Var}(\mathfrak{X}|Y)$  or  $\alpha$ -quantiles of the conditional distribution  $\mathfrak{X}|Y$ . By the assumed conditional independence, we have

$$\mathbf{E}(\mathfrak{X}|Y) = \mathbf{E}(\mathbf{E}(\mathfrak{X}|S)|Y),$$

and, assuming that  $\mathbf{E}(\mathfrak{X}|S)$  is known analytically, we may estimate the conditional expectation by

$$\sum_{i=1}^N W^i \mathbf{E}(\mathfrak{X}|S = S^i).$$

In the case that  $\mathbf{E}(\mathfrak{X}|S)$  is not known analytically, but we can simulate from the conditional distribution  $\mathfrak{X}|S$ , we can obtain samples  $\mathfrak{X}^i$ ,  $i = 1, \dots, N$  where  $\mathfrak{X}^i$  is a draw from  $\mathfrak{X}|S = S^i$ . By the



conditional independence  $\mathfrak{X} \perp Y|S$ , we have  $p(\mathfrak{x}, s|y) = p(\mathfrak{x}|s)p(s|y)$ , and  $g(\mathfrak{x}, s|z) = p(\mathfrak{x}|s)g(s|z)$ , so

$$\frac{p(\mathfrak{x}, s|y)}{g(\mathfrak{x}, s|z)} = \frac{p(s|y)}{g(s|z)} \propto \frac{p(y|s)}{g(z|s)}$$

and  $(\mathfrak{X}^i, S^i), i = 1, \dots, N$  are draws from a proposal whose auto-normalized weights coincide with  $W^i$ . Thus, we may estimate  $\mathbb{E}(\mathfrak{X}|Y)$  by

$$\sum_{i=1}^N W^i \mathfrak{X}^i.$$

While this produces estimates with slightly larger variance (due to the additional simulation), we can control the simulation error by choosing the sample size large enough.

Similarly, by considering  $\mathfrak{X}^2$ , we can estimate the conditional variance, and by considering  $\mathbf{1}_{\mathfrak{X} \leq x}$ , we may estimate the conditional cumulative distribution function of  $\mathfrak{X}$  given  $Y$  at  $x$ , which is just the empirical CDF of samples  $\mathfrak{X}^i$  with associated weights  $W^i$ ,  $i = 1, \dots, N$ . Consequently, we can estimate the  $\alpha$ -quantile of  $\mathfrak{X}|Y$  by the  $\alpha$ -quantile of this empirical CDF, where we use the standard convention for empirical quantiles of linear interpolation between samples to make quantiles unique, see also (Durbin and Koopman, 2012, Section 11.5.3).

### 3.7 Comparison of Importance Sampling method

We now have three tools to produce Gaussian importance sampling proposals: the LA, the CE-method and EIS. Naturally, we want to choose the optimal tool for the problem at hand. In this section, we investigate under which circumstances which method is to be preferred over the others. To judge the performance of each method, we will discuss the following criteria:

- breakdown of methods,
- time and space complexity of the method,
- speed of stochastic convergence, as indicated by the asymptotic variance, for the CE-method and EIS, and
- performance of the optimal proposal, as measured by the efficiency factor, with a focus on degradation as the dimension of the target grows.

Let us elaborate on these criteria. With a breakdown of the methods, we mean settings in which either the numerical scheme diverges, produces parameters that lead to invalid proposals, i.e. negative variances, or where the proposals fail to produce consistent importance sampling estimates.

Time and space complexity of a single iteration of each method allow us to compare the theoretical computational resources required to apply the method in practice and can be used to inform our choice, provided there is no relevant difference in performance of the achieved importance sampling proposals. The speed of stochastic convergence is relevant for the CE-method and EIS as well: The smaller the asymptotic variance, the smaller we can choose the sample size  $N$  and thus decrease computation time. Similarly, numerical convergence directly affects computation time.

Finally, if one method has vastly better performance at the optimum, we might be willing to spend more time initially to save time later when we use the proposal to perform inference. Of special interest is the performance for long (large  $n$ ) or fat (large  $m$ ) time series, as the models we fit in Chapter 4 usually fall into one of these categories.

#### 3.7.1 Breakdown of methods

It is generally difficult to determine whether the proposals produced by the three methods under consideration are valid, i.e. whether the second moment  $\rho$  is finite, see the discussion surrounding Example 3.2. Nevertheless, by focusing on Gaussian targets, one can employ Lemma 3.6. Let us thus begin with a classical example in which the LA fails to produce consistent importance sampling estimates.

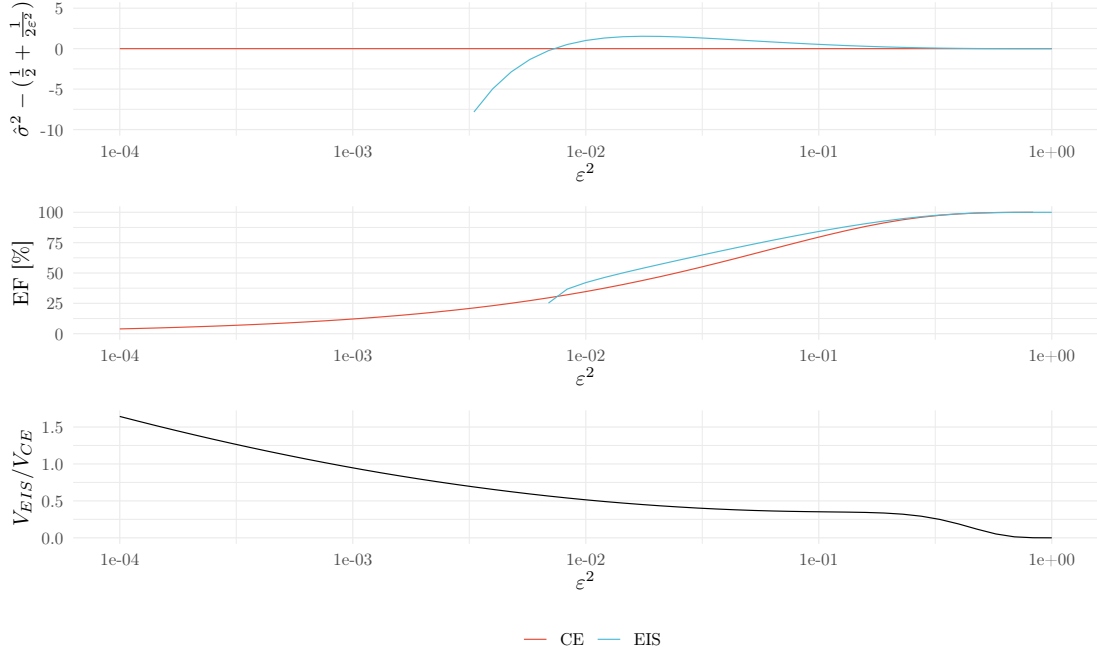


Figure 3.3: Performance of the CE-method and EIS for Example 3.3. The **top** figure shows the excess variance, with positive values corresponding to importance sampling proposals that provide consistent estimates. The **middle** figure shows the degeneration of the efficiency factor as  $\varepsilon^2$  goes to 0. The **bottom** figure shows the asymptotic relative efficiencies of both methods, with EIS outperforming the CE-method for practically relevant values of  $\varepsilon^2$ .

**Example 3.3** (Failure of LA). Consider the Gaussian scale mixture  $\mathbf{P} = \frac{1}{2} (\mathcal{N}(0, 1) + \mathcal{N}(0, \varepsilon^{-2}))$  with mode  $x^* = 0$ , this is the same setup as in Example 3.2. We will perform importance sampling with a normal distribution  $\mathbf{G}_\psi = \mathcal{N}(\mu, \sigma^2)$  for  $\psi = (\mu, \sigma^2)$ .

The LA is  $\mathbf{G}_{\text{LA}} = \mathcal{N}\left(0, \frac{1}{\varepsilon^2 - \varepsilon + 1}\right)$ , whose variance goes to 1 as  $\varepsilon$  goes to 0, so the LA will miss close to  $\frac{1}{2}$  of the total mass. For  $\varepsilon$  small enough, the variance of the LA will be smaller than  $\frac{1}{2\varepsilon^2}$ , whence the second moment of the weights is infinite (Lemma 3.6) and importance sampling with  $\mathbf{G}_{\text{LA}}$  is inconsistent.

The CE-method minimizes the KL-divergence between  $\mathbf{P}$  and  $\mathbf{G}_\psi$  and the optimal proposal is given by  $\mathbf{G}_{\text{CE}} = \mathcal{N}(0, \sigma_{\text{CE}}^2)$ , where  $\sigma_{\text{CE}}^2 = \frac{1}{2} (1 + \varepsilon^{-2})$  is the variance of  $\mathbf{P}$ , as the CE-method matches the first and second moments in this case. As  $\sigma^2 > \frac{1}{2}\varepsilon^{-2}$ , the weights have finite second moment, and importance sampling with  $\mathbf{G}_{\text{CE}}$  is consistent.

As EIS does not yield analytically tractable proposals in this setting, we resort to a simulation study. Using the same setup as described in Example 3.4, we replicate  $M = 100$  times  $\hat{\psi}_{\text{EIS}}$  for varying levels of  $\varepsilon^2$ , averaging over the  $M$  runs. The resulting excess variances, i.e.  $\hat{\sigma}_{\text{EIS}}^2 - \left(\frac{1}{2} + \frac{1}{2\varepsilon^2}\right)$  and  $\sigma_{\text{CE}}^2 - \left(\frac{1}{2} + \frac{1}{2\varepsilon^2}\right)$ , efficiency factors and asymptotic efficiencies are displayed in Figure 3.3. We see that for small  $\varepsilon^2$ , EIS is inconsistent, while the CE-method stays consistent. However, as is to be expected, for small  $\varepsilon^2$ , the efficiency factor becomes very small for both methods, as the tails of both proposals are thinner than those of the target.

Regarding the asymptotic relative efficiencies, we see that  $V_{\text{EIS}} < V_{\text{CE}}$  when  $\varepsilon^2$  is large, and only for very small  $\varepsilon^2$ , i.e. mixtures where one component has very different tail behavior than the other, the CE-method has smaller asymptotic variance.

In the EGSSM setting EIS may produce invalid proposals, as estimates of the variance component

method	single iteration (time)	single iteration (space)	simulation (time)
LA	$\mathcal{O}(np^3)$	$\mathcal{O}(np^2)$	$\mathcal{O}(n(p^3 + m^3 + Nm^2))$
EIS	$\mathcal{O}(n(m^2 + p^3 + Np^2))$	$\mathcal{O}(Np + n(p^2 + m^2))$	$\mathcal{O}(n(p^3 + m^3 + Nm^2))$
CE-method	$\mathcal{O}(n(Nm^2 + m^3))$	$\mathcal{O}(Nm + nm^2)$	$\mathcal{O}(Nnm^2)$

Table 3.2: Computational complexities of importance sampling algorithms for EGSSM with linear signals.

in the weighted least squares regression are not guaranteed to be negative. Thus EIS may produce negative variances. To deal with this, the original EIS paper (Richard and Zhang, 2007, Section 3.2) recommends either inflating the prior or setting the parameters in question to arbitrary fixed values. Alternatively using a more expensive constrained linear least squares solver, such as a conjugate-gradient method (Branch, Coleman, and Y. Li, 1999) or the BVLS (bounded variable least squares) solver (Stark and Parker, 1995) may be appropriate, as is re-running the EIS procedure with a different random seed. Finally, in the EGSSM setting, we could also identify the corresponding observation as missing, similar to the argument presented in Section 3.5.1 for the CE-method.

The CE-method presented in Section 3.5.2 (Algorithm 8) depends on the fact that the covariance matrix of the posterior  $\text{Cov}(X|Y = y)$  is symmetric positive definite (SPD), i.e. non-singular. This might be violated if, e.g., the model contains seasonal components whose associated innovations have variance 0. In this case, the Cholesky roots involved will not be unique. Still Algorithm 8 will, as  $N \rightarrow \infty$  converge a globally optimal solution, though it may not be unique.

### 3.7.2 Computational complexity

Throughout this section, we assume that the model in question is an EGSSM with linear signal (c.f. Definition 3.5) to simplify the treatment. This assumption benefits the LA and EIS approaches, enabling the use of the efficient simulation and signal smoother. If the observation dimension  $p$  is smaller than that of states  $m$ , these algorithms are more efficient, and we will adopt them as well. An overview of computational complexities is presented in Table 3.2. It is important to acknowledge that most operations can be parallelized in one way or the other, e.g. sampling from the proposals. Therefore the time-complexities may not accurately reflect real-world-performance. Nevertheless, they provide theoretical insight into the performance of the three methods considered.

Let us begin with a discussion of the computational complexity associated with determining the optimal parameters,  $\psi_{\text{LA}}$ ,  $\psi_{\text{EIS}}$  and  $\psi_{\text{CE}}$ . In this context we focus on a single iteration and consider the number of iterations as fixed.

As the LA is based on the Kalman-smoother, the time complexity of a single iteration is  $\mathcal{O}(n(m^2 + p^3))$ . The CE-method and EIS need to generate  $N$  samples from the current proposal. For the CE-method this amounts to  $\mathcal{O}(Nnm^2)$  operations (see Section 3.5.2). For EIS, using the simulation smoother (Durbin and Koopman, 2002) requires  $\mathcal{O}(n(m^2 + p^3 + Np^2))$  operations: we need to run the Kalman filter once, while preparing the matrices required for the simulation smoother. Then, assuming Cholesky roots of the innovation covariance matrices  $\Sigma_t$  are already available, only matrix-vector multiplications are necessary for the simulation smoother. Obtaining the EIS model parameters is efficient, requiring only  $\mathcal{O}(n(Np^2 + p^3))$  operations for constructing the  $n \times p$  design matrices and estimating the optimal parameters.

Another concern is the time required to generate  $N$  samples from the fitted model. For both the LA and EIS, this procedure requires using either the simulation smoother or the FFBS algorithm. This necessitates inverting  $p \times p$  matrices in the Kalman filter and  $m \times m$  matrices when simulating the states. Notably, these computational steps can be performed offline, after which the simulation of a single sample requires only  $\mathcal{O}(n)$  matrix-vector multiplications. The CE-method simulation is based on applying Equation (3.27), which requires  $\mathcal{O}(nm^2)$  time per sample.

With respect to space complexity, the LA implementation has to run the Kalman filter which requires  $\mathcal{O}(n(p^2 + m^2))$  space and storage of  $\mathcal{O}(np)$  parameters. EIS has the same space requirement,

yet requires needs additional  $\mathcal{O}(Np)$  storage for the simulated signals — it is sufficient to store just a single set of signals at once, as we can integrate the marginal EIS step into the simulation smoother. As the weights  $w_t$  in EIS depend only on the current signals  $S_t^1, \dots, S_t^N$ , they can be discarded afterwards. See Section 3.5.2 for the derivation of the  $\mathcal{O}(Nm + nm^2)$  space requirement of the CE-method.

The LA has the fastest and most space-efficient iteration of the three methods because it does not require the simulation of  $N$  samples. This makes it an ideal candidate as an initial guess for the other two methods. For  $p \ll m$ , EIS is faster than CE-method as it is based on the signals  $S$  only, thus having access to the efficient simulation and signal smoother algorithms. The same is true for the space complexity. If, however,  $p \approx m$ , there is no linear signal or the observations are not conditionally independent given the states or signals, the speed of a single iteration of EIS and CE-method are comparable. While theoretically, the CE-method performs sampling faster than the other two methods, for large numbers of samples  $N$  the difference is negligible because the additional computations only have to be performed once.

### 3.7.3 Asymptotic variance and relative efficiencies

As we have seen in the previous section, the number of samples  $N$  used to estimate  $\psi_{\text{CE}}$  and  $\psi_{\text{EIS}}$  enter linearly into the computational complexities. Naturally, we want to know how big a sample size we should choose for our procedures to estimate a proposal that is close to the true optimal value and whether one of the two simulation-based procedures requires fewer samples than the other. To answer this question we turn to the two central limit theorems, Theorems 3.6 and 3.9. If  $N$  is large, the asymptotic variances (or rather: the asymptotic standard deviations) tell us how much stochastic variation we should expect around the optimal value, and can thus guide us in choosing  $N$ . We start with two examples in an univariate setting, where both the CE-method and EIS use Gaussian proposals with either fixed variance (Example 3.4) or mean (Example 3.5).

To compare both methods we will determine the asymptotic relative efficiencies, i.e.  $\frac{\text{Var}(\hat{\psi}_{\text{EIS}})}{\text{Var}(\hat{\psi}_{\text{CE}})}$ , with values smaller than 1 indicating that EIS requires (asymptotically) fewer samples for the same precision as the CE-method. Let us note that we are comparing the efficiencies of parameters  $\psi$ , not those of derived parameters such as the standard deviation or the ESS. However, should both methods have the same optimal value, the relative efficiencies are the same for all parameters derived from  $\psi$ , by the delta method. By a continuity argument, the same is approximately true if the optimal values of the CE-method and EIS are close.

To make as much the of the following examples analytically tractable, we will apply the CE-method and EIS in a univariate and single-parameter setting. This allows us to focus on the distinctive properties of the two methods, investigating under which circumstances each method performs well or poorly. In addition, we will use Gaussian proposals where mean (or variance) is fixed and the optimal variance (mean) is obtained by the CE-methods or EIS. As such, we are able to focus on specifying of either the mean or variance and determine which of the two is more crucial to specify accurately. Additionally, the univariate setting allows us, in some cases, to derive analytical expressions of the efficiencies involved.

**Example 3.4** (univariate Gaussian proposal,  $\sigma^2$  fixed). On  $\mathbf{R}$ , consider the probability measure  $\mathbf{P} = p\lambda$  for the Lebesgue measure  $\lambda$  and assume that  $\mathbf{P}$  is symmetric around 0, i.e.  $p(-x) = p(x)$  for  $\lambda$ -a.e.  $x \in \mathbf{R}$  and possesses up to third order moments. Let  $\mathbf{G} = \mathbf{P}$  be a proposal, so  $W \equiv 1$  and let  $\mathbf{G}_\psi = \mathcal{N}(\sigma\psi, \sigma^2)$  be the single parameter natural exponential family of Gaussians with fixed variance  $\sigma^2 > 0$ . Then

$$\log g_\psi(x) = \psi T(x) - \frac{\psi^2}{2} + \log h(x),$$

where  $T(x) = \frac{x}{\sigma}$  and  $h(x)$  is the density of  $\mathcal{N}(0, \sigma^2)$  w.r.t. Lebesgue measure. Note that  $T$  is centered under  $\mathbf{P}$ . To compare the asymptotic behavior of the CE-method and EIS we compute the asymptotic variances arising from their respective central limit theorems (Theorems 3.6 and 3.9).

By symmetry, both  $\psi_{\text{CE}}$  and  $\psi_{\text{EIS}}$  are equal to 0. The Fisher information  $I(\psi)$  is equal to 1 for all

$\psi$ , so

$$V_{\text{CE}} = \text{Cov}_{\mathbf{P}}(T) = \frac{\tau^2}{\sigma^2}, \quad (3.48)$$

where  $\tau^2 = \mathbf{P} \text{id}^2$  is the second moment of  $\mathbf{P}$ .

Additionally,  $B_{\text{EIS}} = (\text{Cov}_{\mathbf{P}}(T))^{-1} = \frac{\sigma^2}{\tau^2}$  and

$$\begin{aligned} M_{\text{EIS}} &= \text{Cov}_{\mathbf{P}} \left( \left( \log \frac{p(x)}{h(x)} - \lambda_{\text{EIS}} \right) T \right) \\ &= \text{Cov}_{\mathbf{P}} \left( (\log p - \log h - \mathbf{P}(\log p - \log h)) T \right) \\ &= \frac{1}{\sigma^2} \int_{-\infty}^{\infty} p(x) x^2 \left( \log p(x) + \frac{x^2}{2\sigma^2} - \mathbf{P} \left( \log p + \frac{\tau^2}{2\sigma^2} \right) \right)^2 dx. \end{aligned}$$

Thus

$$V_{\text{EIS}} = B_{\text{EIS}} M_{\text{EIS}} B_{\text{EIS}} = \sigma^2 \frac{\gamma}{\tau^4},$$

where  $\gamma = \int_{-\infty}^{\infty} p(x) x^2 \left( \log p(x) + \frac{x^2}{2\sigma^2} - \mathbf{P} \left( \log p + \frac{\tau^2}{2\sigma^2} \right) \right)^2 dx$ , and the efficiency of EIS relative to the CE-method is

$$\frac{V_{\text{EIS}}}{V_{\text{CE}}} = \frac{\sigma^4}{\tau^6} \gamma.$$

Let us now consider three exemplary choices of  $\mathbf{P}$  that illustrate a target that is well-behaved (the standard normal), multimodal (a Gaussian location mixture) and has different behavior in the tails than indicated at the mode (a Gaussian scale mixture). For each target, we vary  $\sigma^2$  from  $\frac{1}{2}$  to 3 and obtain relative efficiencies of the CE-method and EIS either analytically or by simulation, the results are shown in the left-hand side of Figure 3.4.

**Normal distribution** If  $\mathbf{P} = \mathcal{N}(0, \tau^2)$  is a normal distribution, this reduces to

$$V_{\text{EIS}} = \frac{5}{2} \left( \frac{\tau^2}{\sigma^2} - 1 \right)^2 \frac{\sigma^2}{\tau^2} = \frac{5}{2} \frac{(V_{\text{CE}} - 1)^2}{V_{\text{CE}}}$$

and so for  $\tau^2 = \sigma^2$  we have  $V_{\text{EIS}} = 0$ , so  $\hat{\psi}_{\text{EIS}}$  might converge faster than the standard  $\mathcal{O}(N^{-\frac{1}{2}})$  rate. Indeed, in this case  $\hat{\psi}_{\text{EIS}} = \psi_{\text{EIS}}$  a.s. for  $N > 1$ , see Proposition 3.5.

**Gaussian location mixture** Consider now the case where  $\mathbf{P} = \frac{1}{2}\mathcal{N}(-1, \omega^2) + \frac{1}{2}\mathcal{N}(1, \omega^2)$  is a Gaussian location mixture. The second moment is  $\tau^2 = 1 + \omega^2 = -\frac{1}{2\psi_{\text{CE}}}$ . Unfortunately, there is no closed-form expression for many of the terms required for the analysis of EIS. Instead, we resort to a simulation study to determine the asymptotic variances and relative efficiencies for three different values of  $\omega^2 \in \{0.1, 0.5, 1.0\}$ .

To this end we draw  $M = 400$  times from the distribution of  $\hat{\psi}_{\text{CE}}$  and  $\hat{\psi}_{\text{EIS}}$ , where we use  $N = 1000$  samples from the tractable  $\mathbf{P}$  as importance samples<sup>8</sup>. We only iterate a single time for both procedures. From individual estimates, we estimate the asymptotic variances  $V_{\text{CE}}$  and  $V_{\text{EIS}}$  by the respective empirical variances, and determine the relative efficiency of EIS over the CE-method as  $\frac{\hat{V}_{\text{EIS}}}{\hat{V}_{\text{CE}}}$ . Again, we vary the fixed variance of the proposals,  $\sigma^2$ , from  $\frac{1}{2}$  to 3. To quantify uncertainty in these asymptotic relative efficiencies, we perform the non-parametric bootstrap with 10 000 samples to estimate their standard errors  $\widehat{\text{se}}_b$ .  $M = 400$  has been chosen to ensure that the relative bootstrap standard error  $\frac{\widehat{\text{se}}_b}{\hat{V}_{\text{EIS}}/\hat{V}_{\text{CE}}}$  is less than 10% across all simulations.

<sup>8</sup>Code for all simulation studies is available in the associated GitHub repository, see Chapter A.

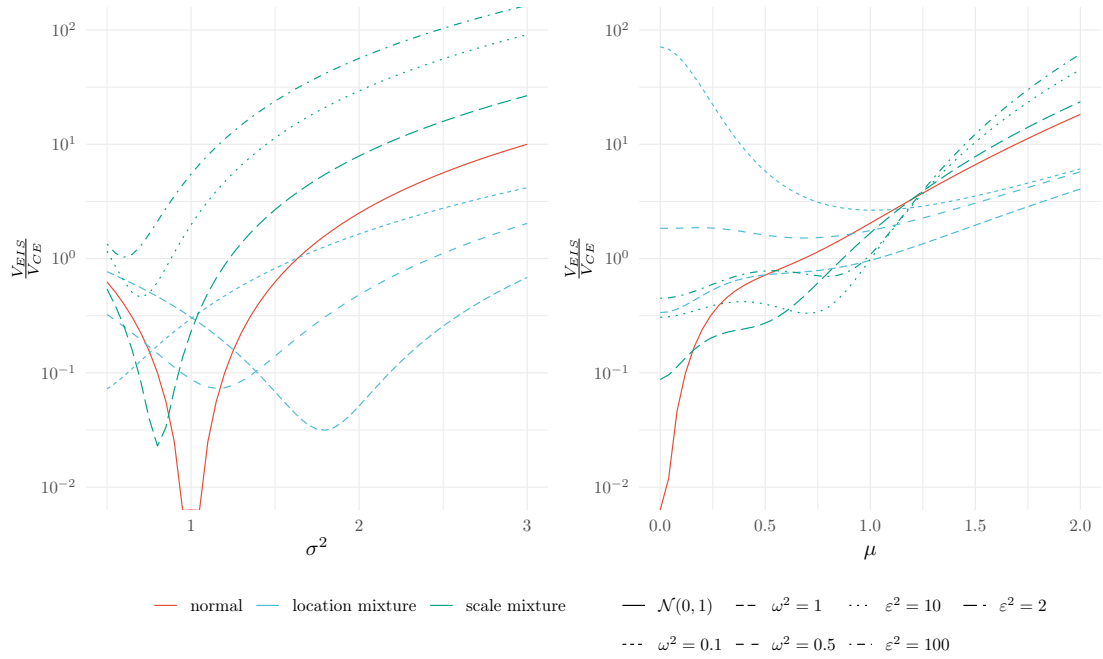


Figure 3.4: Asymptotic relative efficiency  $\frac{V_{EIS}}{V_{CE}}$  for the normal distribution from Example 3.4 (left hand side) and Example 3.5 (right hand side). Here  $\mathbf{P}$  is either the standard normal distribution, a Gaussian location mixture, or a Gaussian scale mixture.  $\mathbf{G}_\psi$  is the normal distribution  $\mathcal{N}(\mu, \sigma^2)$ , where either  $\sigma^2$  is fixed (left) and  $\mu$  determined by the CE-method / EIS, or the other way around (right). Notice the log scale of the y-axis. As  $\mu$  or  $\sigma^2$  get close to their true values, EIS outperforms the CE-method in terms of asymptotic variance, see Proposition 3.5.

**Gaussian scale mixture** Finally we consider  $\mathbf{P} = \frac{1}{2} (\mathcal{N}(0, 1) + \mathcal{N}(0, \varepsilon^{-2}))$  for  $\varepsilon^2 \in \{2, 10, 100\}$ , a scale mixture similar to the one seen in Example 3.3. Contrary to that example, we choose  $\varepsilon$  big, making the  $\mathcal{N}(0, 1)$  component the one with large variance, to make importance sampling with  $\sigma^2$  in the range considered consistent. Here  $\tau^2 = \frac{1}{2} + \frac{1}{2\varepsilon^2}$ . Again, we estimate the asymptotic  $V_{\text{EIS}}$  in the same way as for the Gaussian location mixture, with  $M = 100$  estimates using  $N = 1000$  samples each and obtain a Monte-Carlo standard error for the asymptotic variances of  $1.1 \times 10^{-3}$ .

Note that for fixed  $\sigma^2$  the asymptotic variance of the CE-method  $V_{\text{CE}}$  is the same in all of the examples considered, as we sample directly from the tractable  $\mathbf{P}$ , so  $V_{\text{CE}}$  only depends on  $\mathbf{P}$  through its second moment  $\tau^2$ . The asymptotic variance of EIS however depends on both  $\tau^2$  and  $\gamma$ , which depends on higher order moments of  $\mathbf{P}$ .

From the left-hand side of Figure 3.4 we can observe that in the case of  $\mathbf{P} = \mathcal{N}(0, 1)$  EIS has smaller asymptotic variance compared to the CE-method, as long as  $\sigma^2$  is not heavily misspecified. Indeed, if  $\sigma^2 = 1$  is correctly specified, by Proposition 3.5, EIS has asymptotic variance 0 and converges already for a single sample.

Consider now the case where  $\mathbf{P}$  is a Gaussian location mixture. For  $\omega^2 = 1$ , the location mixture is unimodal with variance 2 and EIS outperforms the CE-method in terms of asymptotic variance in the range considered. For the smaller values of  $\omega^2$  considered here, the location mixture is bimodal. Close to the true variance  $1 + \omega^2$ , EIS still outperforms the CE-method.

For the Gaussian scale mixture, the case is less clear. Here the true variance is  $\frac{1}{2} + \frac{1}{2\varepsilon^2}$ . The location of the minimal relative efficiency is still close to this true variance, however, as  $\varepsilon^2$  grows, the CE-method starts to dominate EIS. Additionally, recall from Example 3.3 that for large  $\varepsilon^2$  EIS becomes inadmissible.

**Example 3.5** (univariate Gaussian,  $\mu$  fixed). Consider the same setup as in Example 3.4, i.e.  $\mathbf{P}$  is symmetric around 0 with second moment  $\tau^2$ , but let  $\mathbf{G}_\psi = \mathcal{N}(\mu, -\frac{1}{2\psi})$  be the single parameter natural exponential family of Gaussians with fixed mean  $\mu$  and variance  $\sigma^2 = -\frac{1}{2\psi}$ .

Then

$$\log g_\psi(x) = \psi T(x) + \frac{1}{2} \log(-2\psi) - \frac{1}{2} \log 2\pi$$

for  $T(x) = (x - \mu)^2$ . Thus  $\mathbf{P}T = \tau^2 + \mu^2$  and  $\text{Cov}_{\mathbf{P}} T = \nu - \tau^4 + 4\tau^2\mu^2$  where  $\nu = \mathbf{P} \text{id}^4$  and  $\tau^2 = \mathbf{P} \text{id}^2$ .

By matching moments, we obtain  $\psi_{\text{CE}} = -\frac{1}{2(\tau^2 + \mu^2)}$  and  $I(\psi_{\text{CE}}) = \frac{1}{2\psi_{\text{CE}}^2} = 2(\tau^2 + \mu^2)^2$ . In total

$$V_{\text{CE}} = \frac{1}{4(\tau^2 + \mu^2)^4} (\nu - \tau^4 + 4\tau^2\mu^2) \quad (3.49)$$

For EIS,

$$\begin{aligned} \psi_{\text{EIS}} &= (\text{Cov}_{\mathbf{P}} T)^{-1} \text{Cov}_{\mathbf{P}} (T, \log p) \\ &= (\nu - \tau^4 + 4\tau^2\mu^2)^{-1} \underbrace{\int p(x)((x - \mu)^2 - \tau^2 - \mu^2)(\log p(x) - \mathbf{P} \log p(x)) dx}_{=\gamma}. \end{aligned}$$

As  $B_{\text{EIS}} = \text{Cov}_{\mathbf{P}} T^{-1} = (\nu - \tau^4 + 4\tau^2\mu^2)^{-1}$ , we have

$$V_{\text{EIS}} = (\nu - \tau^4 + 4\tau^2\mu^2)^{-2} \mathbf{P} \left( (\text{id} - \mu)^4 (\log p - \psi_{\text{EIS}}(\text{id} - \mu)^2 - \mathbf{P} \log p + \psi(\tau^2 + \mu^2))^2 \right),$$

where the last term is equal to  $M_{\text{EIS}}$ .

We now perform the same analysis as in Example 3.4, the resulting ratio of asymptotic variances is displayed in the right-hand side of Figure 3.4. In general, the variances  $\sigma_{\text{CE}}^2 = -\frac{1}{2\psi_{\text{CE}}}$  and

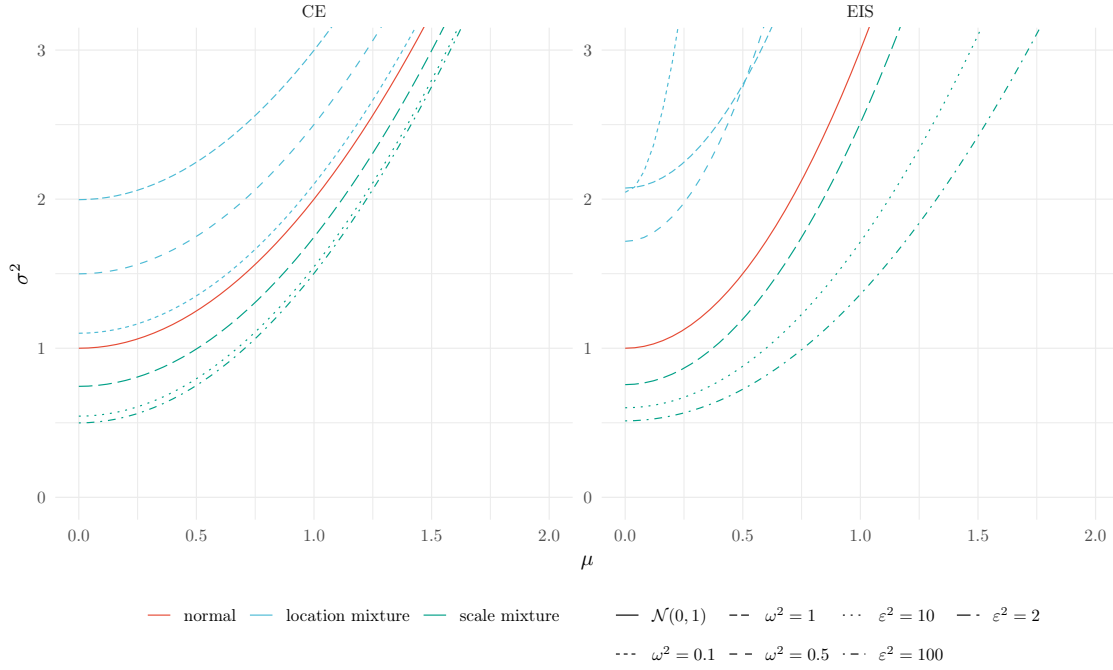


Figure 3.5: Optimal variances  $\sigma_{\text{CE}}^2 / \sigma_{\text{EIS}}^2$  for the CE-method (left) and EIS (right) as a function of the misspecified mean  $\mu$ . The variances produced by EIS tend to be larger than those produced by the CE-method.

$\sigma_{\text{EIS}}^2 = -\frac{1}{2\psi_{\text{EIS}}}$  are different, so the ratio is no longer an asymptotic relative efficiency. However, it is still relevant as a measure of the relative speed of stochastic convergence of both methods. Additionally, we display the resulting optimal variances in Figure 3.5.

**Normal distribution** For the normal distribution  $\mathbf{P} = \mathcal{N}(0, \tau^2)$  where  $\nu = 3\tau^4$  and  $\gamma = -\tau^2$ , so

$$\psi_{\text{EIS}} = \frac{-\tau^2}{2\tau^2(\tau^2 + 2\mu^2)} = \frac{-1}{2(\tau^2 + 2\mu^2)}.$$

Thus the EIS proposal uses variance  $\sigma_{\text{EIS}}^2 = \tau^2 + 2\mu^2$ , which is bigger than the variance of  $\sigma_{\text{CE}}^2 = \tau^2 + \mu^2$  optimal for the CE-method.

In this case the asymptotic variances are

$$V_{\text{CE}} = \frac{\tau^2(\tau^2 + 2\mu^2)}{2(\tau^2 + \mu^2)^4}$$

and

$$V_{\text{EIS}} = \frac{\mu^2(2\mu^6 + 45\mu^4\tau^2 + 15\tau^6)}{4\tau^4(2\mu^2 + \tau^2)^4},$$

see the Appendix for details.

**Gaussian location and scale mixture** To estimate asymptotic relative efficiencies for the Gaussian location and scale mixtures, for the same targets as in Example 3.4, we again perform a simulation study with the same parameters ( $M = 400$  repetitions, 10 000 bootstrap samples estimate the standard error of estimation). Here the choice of  $M$  leads to a relative standard error of at most 11%.



On the left-hand side of Figure 3.4 we see that for  $\mu$  close to the optimal value, EIS has smaller asymptotic variance than the CE-method, except for the two bimodal location measures. Again, due to the finite sample convergence of EIS, Proposition 3.5, the asymptotic variance  $V_{\text{EIS}}$  goes to 0 as  $\mu \rightarrow 0$ . As  $\mu$  is further from the true 0 the ratio of asymptotic variances starts to grow.

In Figure 3.5 we see that, except for the extreme scale mixtures, EIS tends to produce proposals that have a larger variance than those produced by the CE-method. As we will see in the discussion of Figure 3.7, this might be advantageous for EIS as proposals with a small variance run the risk of missing a large part of the probability mass of the target.

In applications, e.g. the model studied in Chapter 4, we are interested in the performance of the importance sampling proposals generated by the LA, CE-method and EIS under more complex circumstances than those discussed in Examples 3.4 and 3.5. In particular, the dimension of  $\psi$  is high ( $\mathcal{O}(n \cdot m)$  or even  $\mathcal{O}(n \cdot m^2)$ ) and proposals may not come from a natural exponential family, so analysis based on Theorems 3.6 and 3.9 is not possible.

really?

Instead, we resort to simulation studies to gain insights into the circumstances when one should prefer one method over the other. As a leading example, we will use the following vector-autoregressive state space model with negative binomial observations. A similar, though more involved, model is studied in Section 4.2 with real data.

**Example 3.6** (Negative Binomial VAR(1) SSM). In this example, we consider a SSM where states  $X_t$  follow a stationary Gaussian VAR(1) process, initialized in its stationary distribution  $\mathcal{N}(0, \Sigma)$  for SPD  $\Sigma$ . For simplicity let the transition matrices be given by a multiple of the identity, i.e.  $A_t = \alpha I_m$  for all  $t$  where  $\alpha \in (-1, 1)$

add I to symbols

. In total, the states are governed by

$$\begin{aligned} X_0 &\sim \mathcal{N}(0, \Sigma) \\ X_t &= \alpha X_{t-1} + \varepsilon_t \\ \varepsilon_t &\stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, (1 - \alpha^2)\Sigma), t = 1, \dots, n \end{aligned}$$

where the  $\varepsilon_1, \dots, n$  and  $X_0$  are jointly independent. The observations follow a conditional negative binomial distribution

$$Y_t^i | X_t \sim \text{NegBinom}(\exp(X_t^i), r), \quad i = 1, \dots, p \quad t = 0, \dots, n$$

and individual observations are conditionally independent given the current state. The parametrization of the negative binomial distribution  $\text{NegBinom}(\mu, r)$  is such that the density is

$$p_{\mu, r}(y) = \binom{y+r-1}{r} \left( \frac{\mu}{r+\mu} \right)^y \left( \frac{r}{r+\mu} \right)^r \propto \mu^y (\mu+r)^{-(r+y)},$$

where proportionality is in  $\mu$ , with expectation  $\mu$ , variance  $\mu + \frac{\mu^2}{r}$  and support  $\mathbf{N}_0$ .

Our first simulation study concerns the non-asymptotic behavior of the CE-method and EIS estimators, i.e. finite sample analogs of Theorems 3.6 and 3.9. To this end, we let  $m = 1$  in Example 3.6 and fix  $n$  to

...

. We then simulate once from the marginal distribution of  $Y$  and perform the LA to a prespecified precision  $\epsilon$  and maximum number of iterations  $n_{\text{iter}}$ , obtaining a proposal distribution  $\mathbf{G}_{\text{LA}}$ . Using a large number of samples  $N_{\text{true}}$  from this proposal we find the optimal  $\mathbf{G}_{\text{CE}}$  and  $\mathbf{G}_{\text{EIS}}$  using the same desired precision and number of iterations as for the LA. For the remainder of this section, we ignore sampling variation in these proposals and treat them as exact.

To determine the non-asymptotic sampling behavior we now perform the above procedure again, using only  $N \ll N_{\text{true}}$  many samples for both procedures, obtaining proposals  $\hat{\mathbf{P}}_{\text{CE}}^N$  and  $\hat{\mathbf{P}}_{\text{EIS}}^N$ . As the full proposals are Gaussian distributions on  $\mathbf{R}^{(n+1) \times m}$ , either given as the posterior of a GLSSM (LA, EIS) or by a Gaussian Markov process (CE-method), see Section 3.5. This procedure is repeated  $M$  times for every sample size  $N$  considered, with different initial random seeds, obtaining  $\hat{\mathbf{P}}_{\text{CE}}^{N,i}$  and  $\hat{\mathbf{P}}_{\text{EIS}}^{N,i}$  for  $i = 1, \dots, M$ .

To assess the speed of convergence of the CE-method and EIS we then estimate the mean squared error of means and variances of the  $(n+1) \times m$  univariate marginals as  $N$ , the number of samples used to obtain  $\hat{\psi}_{\text{CE}}$  or  $\hat{\psi}_{\text{EIS}}$ , grows. For the true value, we take the univariate means and variances of  $\mathbf{G}_{\text{CE}}$  and  $\mathbf{G}_{\text{EIS}}$  respectively. Additionally, we perform a bias-variance decomposition to see where the estimation error originates.

More concretely, fix  $N$  and denote by  $\mu, \sigma^2 \in \mathbf{R}^{(n+1) \cdot m}$  the marginal means and variances of  $\mathbf{G}_{\text{CE}}$  ( $\mathbf{G}_{\text{EIS}}$ ). Let  $\hat{\mu}_i, \hat{\sigma}_i^2 \in \mathbf{R}^{(n+1) \cdot m}$  be the marginal means and variances of  $\mathbf{G}_{\text{CE}}^{N,i}$  ( $\mathbf{G}_{\text{EIS}}^{N,i}$ ) for  $i = 1, \dots, M$ . Now

$$\widehat{\text{aMSE}} = \frac{1}{M} \frac{1}{(n+1)m} \sum_{i=1}^M \|\mu - \hat{\mu}_i\|_2^2 + \|\sigma^2 - \hat{\sigma}_i^2\|_2^2$$

is an estimate of the mean-squared error of  $(\mu, \sigma^2)$ , where we divide by  $(n+1)m$  to make estimates comparable across models of different dimensions.

In Figure 3.6 we show the  $\widehat{\text{aMSE}}$  for both the CE-method and EIS for varying values of  $N$ . As is evident from this Figure, the CE-method consistently has a larger aMSE than EIS, for all values of  $N$ . Thus the CE-method requires several orders of magnitude more samples to obtain the same precision as EIS.

For further investigation, we perform a bias-variance decomposition of the aMSE for both the means  $\mu$  and variances  $\sigma^2$ . Consider the average means and variances over the  $M$  simulations,

$$\bar{\mu} = \frac{1}{M} \sum_{i=1}^M \hat{\mu}_i \quad \bar{\sigma}^2 = \frac{1}{M} \sum_{i=1}^M \hat{\sigma}_i^2,$$

and the state-average squared bias and variance

$$\begin{aligned} \text{aBias}_{\mu}^2 &= \frac{1}{(n+1)m} \|\mu - \bar{\mu}\|_2^2, \\ \text{aVar}_{\mu} &= \frac{1}{M-1} \frac{1}{(n+1)m} \sum_{i=1}^M \|\bar{\mu} - \hat{\mu}_i\|_2^2, \\ \text{aBias}_{\sigma^2}^2 &= \frac{1}{(n+1)m} \|\sigma^2 - \bar{\sigma}^2\|_2^2, \\ \text{aVar}_{\sigma^2} &= \frac{1}{M-1} \frac{1}{(n+1)m} \sum_{i=1}^M \|\bar{\sigma}^2 - \hat{\sigma}_i^2\|_2^2. \end{aligned}$$

These values are depicted in Figure 3.6.

interpretation of Figure 3.6, equal contribution of bias and var, not much to gain from bias correction

is bias of CEM really of this order? would expect bias usually to be of order  $1/n$ , bias squared of order  $1/n$  squared, so negligible compared to  $1/n$  mse?

### 3.7.4 Performance of the optimal proposal

change EF to aEF (asymptotic EF) everywhere in this section

For the performance of importance sampling the efficiency factor  $\text{EF} = \frac{\text{ESS}}{N}$  plays an important role, see Section 3.3. Additionally, it allows a comparison of the effectiveness of importance sampling

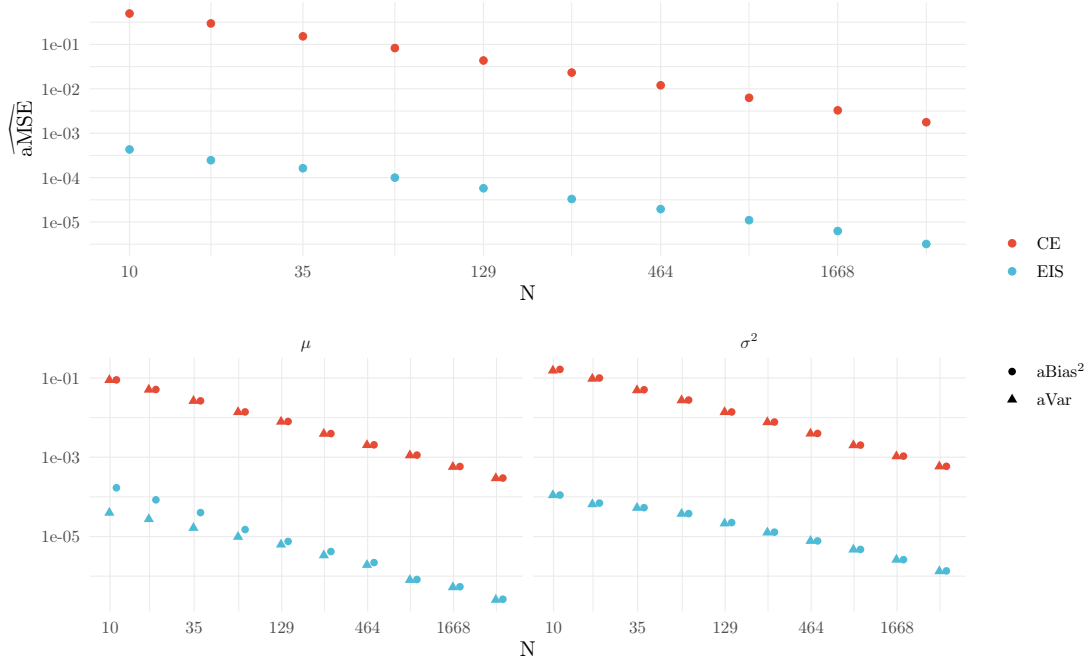


Figure 3.6: TODO

across multiple sample sizes  $N$ , indeed, as  $N \rightarrow \infty$ , EF converges to  $\rho^{-1}$ , where  $\rho$  is the second moment of importance sampling weights,  $\int w^2 d\mathbf{G}$ .

Returning to the distributions studied in Examples 3.4 and 3.5, we now calculate the asymptotic efficiency factor

$$\text{EF} = \frac{1}{\rho} \in (0, 1].$$

As the proposal is always  $\mathcal{N}(\mu, \sigma^2)$  with either  $\mu$  or  $\sigma^2$  fixed, and  $\mathbf{P}$  is a mixture of Gaussians or  $\mathcal{N}(0, 1)$ ,  $\rho$  is analytically available.

For Example 3.4, both EIS and the CE-method have, by symmetry, the same optimal  $\mu = 0$ . Thus the efficiency factor only depends on the fixed  $\sigma^2$ , see Figure 3.7, and is the same for EIS and the CE-method.

For Example 3.5 the two methods have different optimal proposals, thus also different asymptotic efficiency factors. In Figure 3.8, the first two subfigures show how the efficiency factor depends on the misspecified  $\mu$  for both methods. The optimal variances are based on the results from Example 3.5, i.e. based on simulation for EIS. The right-hand subfigure shows the relative efficiency factor, i.e. the ratio of the efficiency factor for the CE-method and EIS. Here values smaller than 1 indicate that EIS has a larger efficiency factor than the CE-method.

In this figure, we can observe that, as expected, misspecification in  $\mu$  almost always results in a smaller efficiency factor, an exception being the scale mixture with  $\varepsilon^2 = 100$  for the CE-method. Compared to Figure 3.7, we see that already small misspecification in  $\mu$  results in a large decline in EF, although we should keep in mind that this is not a fair comparison, as  $\mu$  and  $\sigma^2$  live on different scales. If  $\mu = 0$  is correctly specified, both methods have comparable performance, except for extreme cases of the mixture models, i.e. when  $\omega^2 = 0.1$  or when  $\varepsilon^2 = 100$ . For small misspecification of  $\mu$ , this remains true, but for larger misspecification, the CE-method has a larger efficiency factor, especially for the bimodal location mixture with  $\omega^2 = 0.1$ , where the performance of EIS deteriorates.

stress that cem gives global optimum, eis only approximate

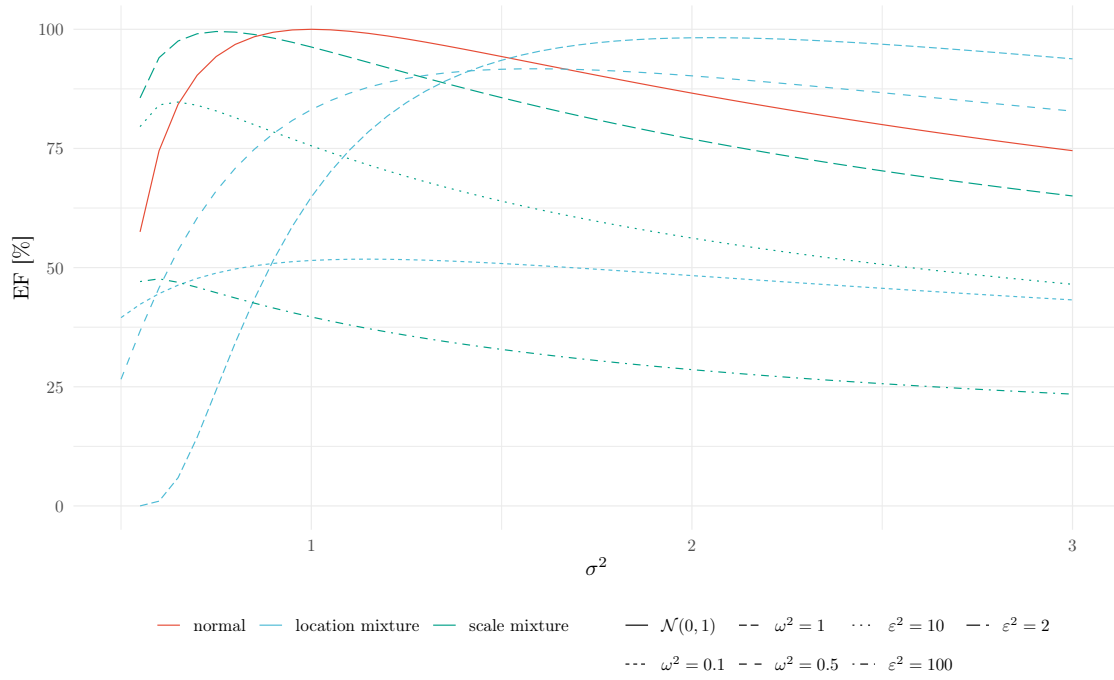


Figure 3.7: TODO

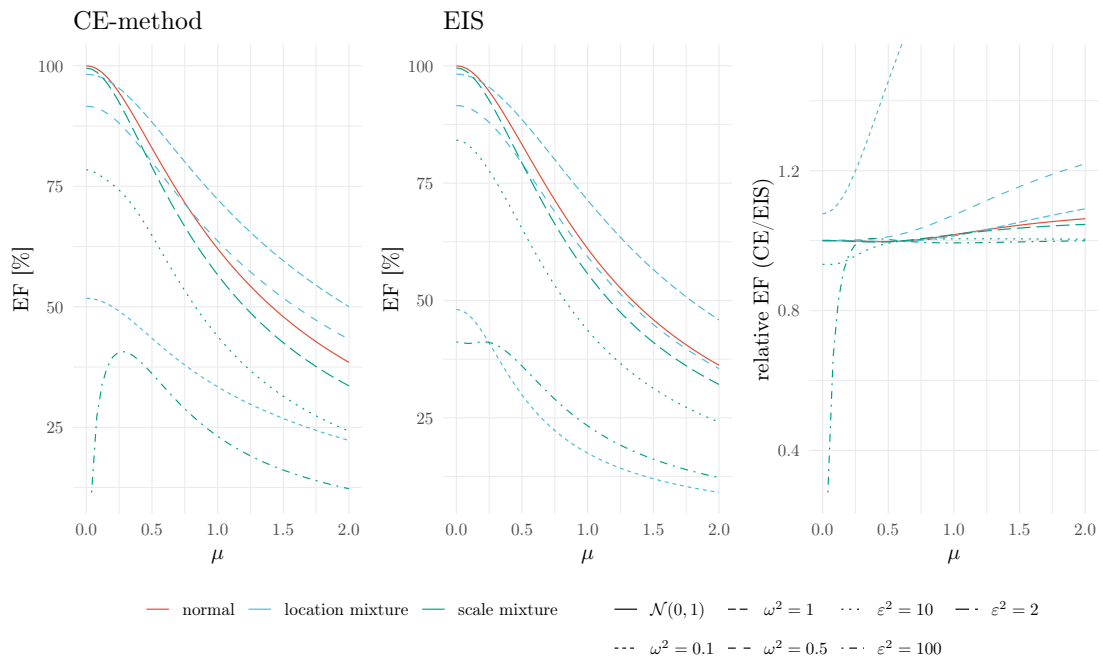


Figure 3.8: TODO

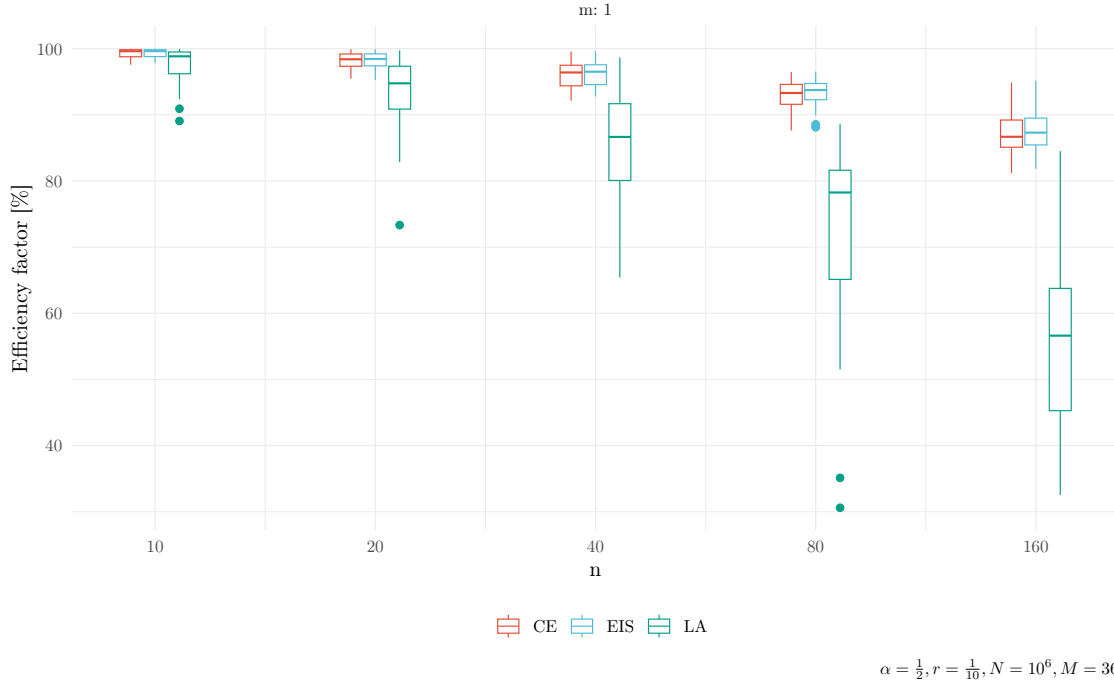


Figure 3.9: The asymptotic efficiency factor degenerates as the number of time steps  $n$  increases. We show the estimated efficiency factor over 100 replications of estimating the optimal parameters for Example 3.6 with the CE-method and EIS with  $N_{\text{true}} = 10^6$  and the resulting estimated efficiency factors at the optimum. Notice the log scale of the x-axis. The performance of the optimal CE-method and EIS parameters is comparable and superior to that of the LA

For the model from Example 3.6 we cannot determine  $\rho$  analytically, so we fall back to a simulation study. Thus, we also estimate EF for each of the  $M$  runs, using the same number of samples  $N = N_{\text{true}}$  as was used to determine the true optimal parameter. We display the resulting efficiency factors in Figure 3.9. The parameters  $\alpha, r, N, M$  may be found in the bottom right corner of the figure. For a low number of time steps  $n$ , all three methods perform comparably. With increasing  $n$ , their performance expectedly worsens, however, more so for the local LA, while the CE-method and EIS perform comparably around their optimal value.

### 3.8 Conclusion

note CEM fails because it cannot work with independent marginals compared to EIS

compare independent components exponential family



## Chapter 4

# Analysis of selected models

### Contributions of this chapter

The main contribution of this chapter is to apply the methods derived in Chapter 3 to selected inference and prediction problems in the context of COVID-19 in Germany.

**Removing reporting delays and weekday effects**

**Regional growth factor model**

**Nowcasting hospitalizations**

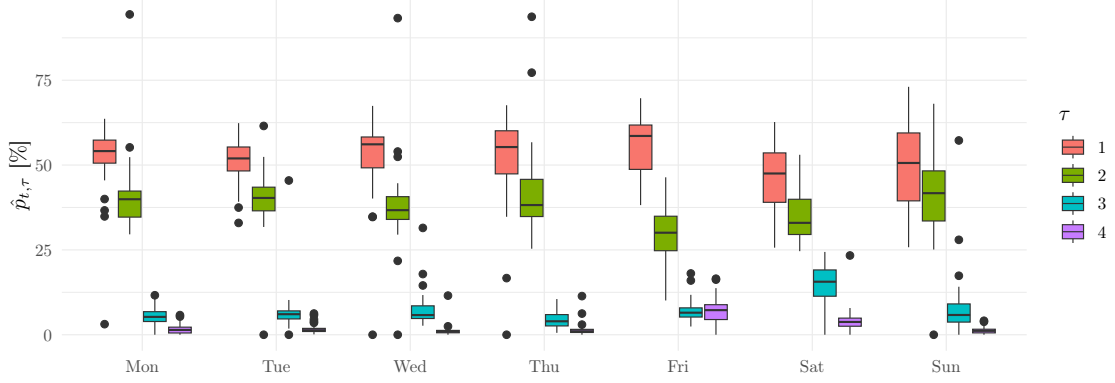


Figure 4.1: Box plots of delay probabilities  $\hat{p}_{t,\tau}$  by weekday of case reporting date  $t$ . As there are systematically fewer cases reported on Sunday, there is a small weekday effect:  $p_{t,1}$  for Saturdays,  $p_{t,2}$  for Fridays,  $p_{t,3}$  for Thursdays and  $p_{t,4}$  for Wednesdays are small compared to other days.

all analyses on Macbook Air M3, if not noted otherwise, include computation time

## 4.1 Removing reporting delays and weekday effects

### 4.1.1 Context

Retrospective analysis of the reported cases is one of the main tasks of epidemiological monitoring, see ???. For this analysis, it is crucial to have the finest temporal resolution possible, as we want, e.g., to link the dates on which NPIs were enforced to the growth factors on, or surrounding these dates. However, as we have observed in ???, the data at our hand are contaminated by the reporting process, most notably the weekday effects, reporting delays, and reporting artifacts, e.g. due to public holidays. While we can make the influence of these effects small by aggregating data to the weekly level, see Section 4.2, modeling on the daily level facilitates better retrospective analyses and as such it is the goal of this section.

Here, we use the RKI case incidence data discussed in ???. As we have seen in ??? A and ??, most delays are shorter than 4 days. Thus, ignoring any cases reported with longer delays, we get for any reporting date  $t$  four observations, say

$$Y_t = (Y_{t,1}, \dots, Y_{t,4}) \in \mathbf{N}_0^4.$$

Here  $Y_{t,\tau}$ ,  $\tau = 1, \dots, 4$ , is the number of newly reported cases for reporting date  $t$  with delay  $\tau$ , such that  $Y_{t,\cdot} = \sum_{\tau=1}^4 Y_{t,\tau}$  is the total number of cases reported for reporting date  $t$  with delay  $\leq 4$ . Let  $\hat{p}_{t,\tau} = \frac{Y_{t,\tau}}{Y_{t,\cdot}}$  be the empirical delay probability for day  $t$  with delay  $\tau$ . We have already observed in ???, that  $Y_{t,\cdot}$  is subject to weekday effects, and similar to hospitalizations (???), there is a small weekday effect for the delay of cases, i.e.  $\hat{p}_{t,\tau}$ , see Figure 4.1.

To produce accurate retrospective analyses of the daily growth factor, we will construct a SSM that allows to account for these delays, as well as the weekday effects and dynamics of the incidences. This model will enable us to better understand the delay process, allow to account for periods of inconsistent reporting, and yield daily growth factors useful for the interpretation of NPI efficacy.

### 4.1.2 Model

To model the development of cases over time, we start with the exponential growth equation ???. Let  $I_t$  be the total number of cases for reporting date  $t$ , unaffected by weekday effects and reporting delays. Ignoring variation around the mean, the exponential growth ansatz gives

$$\log I_{t+1} \approx \log \rho_{t+1} + \log I_t$$



for the growth factor  $\rho_t$  on day  $t$ . It is then sensible to assume that the growth factor  $\rho_t$  performs a random walk on the log-scale, as we would expect large day-to-day variation of  $\rho_t$  for large values, and small variation for small values, i.e. multiplicative, rather than additive, day-to-day changes. Thus, we assume that

$$\log \rho_{t+1} = \log \rho_t + \varepsilon_{t+1,\rho}$$

for  $\varepsilon_{t+1,\rho} \sim \mathcal{N}(0, \sigma_\rho^2)$ . To incorporate weekday effects, consider a weekly seasonal component

$$\log W_{t+1} = - \sum_{s=0}^5 \log W_{t-s} + \varepsilon_{t+1,W},$$

for  $\varepsilon_{t+1,W} \sim \mathcal{N}(0, \sigma_W^2)$ , as described in ??

explicit this here

. Finally, to model the reporting delay probabilities  $p_{t,\tau}$ ,  $\tau = 1, 2, 3, 4$ , we parameterize them by log ratios

$$q_{t,\tau} = \log \frac{p_{t,\tau}}{p_{t,4}} \quad \tau = 1, 2, 3,$$

which also perform a random walk in time:

$$q_{t+1,\tau} = q_{t,\tau} + \varepsilon_{t+1,q,\tau},$$

with  $\varepsilon_{t+1,q,\tau} \sim \mathcal{N}(0, \sigma_q^2)$  whose variance does not depend on the delay  $\tau$ . To account for the weekday effect visible in Figure 4.1, we introduce three further weekday effects, for  $\tau = 1, 2, 3$  let

$$\log W_{t+1}^{q,\tau} = - \sum_{s=0}^5 \log W_{t-s}^{q,\tau} + \varepsilon_{t+1,W^{q,\tau}},$$

with  $\varepsilon_{t+1,W^{q,\tau}} \sim \mathcal{N}(0, \sigma_{W_q}^2)$  and shared variance  $\sigma_{W_q}^2$ . We can recover the delay probabilities  $p_{t,\tau}$  from the log-ratios by

$$\begin{aligned} p_{t,4} &= \frac{1}{1 + \sum_{\tau=1}^3 \exp(q_{t,\tau} + \log W_t^{q,\tau})}, \\ p_{t,\tau} &= \exp(q_{t,\tau} + \log W_t^{q,\tau}) p_{t,4}, \end{aligned} \quad (4.1)$$

for  $\tau = 1, 2, 3$ .

Finally, there are reporting artifacts and other effects that we have not yet considered in our model contribute to the dirtiness of the data. To account for these effects, we model daily, multiplicative, „muck“  $M_t$ , for date  $t$ , such that the total expected number of reported cases on this date is  $M_t I_t$  instead of  $I_t$ . We assume that  $(\log M_t)_{t=0,\dots,n} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(-\frac{1}{2}\sigma_M^2, \sigma_M^2)$ , independent of all other states. Thus,  $M_t$  follows a log-normal distribution with mean 1.

With these components at our disposal, we can model the observed incidences  $Y_{t,\tau}$  by

$$Y_{t,\tau} | \log I_t, \log W_t, q_t, \log M_t \sim \text{Pois}(p_{t,\tau} \exp(\log I_t + \log W_t + \log M_t)), \quad (4.2)$$

conditionally independent for fixed  $t$ . Thus,  $W_t$  acts as a multiplicative factor that modulates the observed cases depending on the day of the week, and the delay probabilities distribute the total expected number of cases  $M_t W_t I_t$  onto the delays. In this model,  $Y_t = \sum_{\tau=1}^4 Y_{t,\tau}$  has conditional expectation

$$\mathbb{E}(Y_t | \log I_t, \log W_t, q_t, \log M_t) = M_t W_t I_t$$

As it is sensible to model the conditional distribution of  $Y_t$  by a Poisson distribution (see ??), we can view Equation (4.2) as a multinomial thinning of this distribution. Notice that including  $M_t$  introduces overdispersion in this Poisson distribution, similar to modeling with a negative binomial distribution.

Letting  $X_t = \left( \log I_t, \log \rho_{t+1}, \log W_t, \dots, \log W_{t-5}, q_{t,1}, q_{t,2}, q_{t,3}, \log W_t^{q,1}, \dots, \log W_{t-5}^{q,3} \right)^T$ , assuming that

$$\varepsilon_{t+1} = \begin{pmatrix} \varepsilon_{t+1,\rho} \\ \varepsilon_{t+1,W} \\ \varepsilon_{t+1,q,1} \\ \varepsilon_{t+1,q,2} \\ \varepsilon_{t+1,q,3} \end{pmatrix}$$

has independent marginals, and fixing an initial distribution of  $X_0$  fully specifies a PGSSM for the joint distribution of  $(X, Y)$ . For the initial distribution we use

$$X_0 \sim \mathcal{N}(u_0, \Sigma_0)$$

where  $u_0$  is 0 for all elements, except to the third entry (corresponding to  $M_0$ ), which we set to  $-\frac{1}{2}\sigma_M^2$ . For the initial covariance we use a diagonal matrix

$$\Sigma_0 = \text{diag} \left( \underbrace{25}_{\log I}, \underbrace{0.2^2}_{\log \rho}, \underbrace{s_M^2}_M, \underbrace{1}_{\log W_0}, \dots, \underbrace{1}_{\log W_{-5}}, \underbrace{1}_{q_{0,1}}, \underbrace{1}_{q_{0,2}}, \underbrace{1}_{q_{0,3}}, \underbrace{1}_{\log W_0^{q,1}}, \dots, \underbrace{1}_{\log W_{-5}^{q,3}} \right).$$

The model has a linear signal

$$S_t = \begin{pmatrix} \log I_t + \log W_t \\ q_{t,1} \\ q_{t,2} \\ q_{t,3} \end{pmatrix},$$

but due to the non-linear dependence of  $p_{t,\tau}$  on  $q_{t,\tau}$ ,  $Y_{t,\tau}$  depends not just on  $S_{t,\tau}$  but on the whole of  $S_t$ . Fortunately, this is not a problem for either the LA or EIS. For the LA (Algorithm 5), notice that the covariance matrix  $\Omega_t$  is given by the inverse of the negative Hessian of  $s_t \mapsto \log p(y_t|s_t)$ , which is now non-diagonal. While it is not guaranteed that  $\Omega_t$  is positive semi-definite during the Newton-Raphson iteration, we can still employ the Kalman filter and signal smoother to perform the iteration efficiently, see (Jungbacker and Koopman, 2007) and the discussion in Section 3.5.1. Furthermore, at the global optimum, the Hessian is negative semi-definite, so  $\Omega_t$  is positive semi-definite, specifying a valid GLSSM proposal. Similarly, we may extend EIS to account for non-diagonal  $\Omega_t$ . Recall from Section 3.3.3, that EIS minimizes for a given  $t$

$$\sum_{i=1}^N \left( \log p(y_t|S_t^i) + \langle \Omega_t^{-1} z_t, S_t^i \rangle - \frac{1}{2} \text{tr}(\Omega_t^{-1} S_t^i (S_t^i)^T) - \lambda_t \right)^2$$

over  $z_t, \Omega_t, \lambda_t$ . Noticing that  $(A, B) \mapsto \text{tr}(A^T B)$  is the Frobenius inner-product, we see that this optimization problem is still a weighted linear least squares problem for  $\Omega_t^{-1} z_t, \Omega_t^{-1}, \lambda_t$ , when we let  $\Omega_t^{-1}$  take values in the symmetric matrices in  $\mathbf{R}^{p \times p}$ . As the dimension of this vector space is  $\frac{p(p+1)}{2}$ , we may still perform the computationally efficient weighted linear least squares routine, but at an increased cost: the number of parameters increases from  $2p+1$  ( $\Omega_t$  diagonal) to  $p + \frac{p(p+1)}{2} + 1$  ( $\Omega_t$  symmetric).

The parameters of the model are  $\theta = \left( \log \sigma_\rho^2, \log \sigma_W^2, \log \sigma_q^2, \log \sigma_M^2, \log \sigma_{W_q}^2 \right)$ , which we model on the log-scale to avoid having to take care of constraints. Given observations  $Y = (Y_0, \dots, Y_n)$  we perform maximum likelihood estimation as described in Section 3.6.1. As tuning parameters in this procedure we use 20 iterations for the LA and EIS, with relative tolerance of convergence set to  $10^{-5}$ . For the EIS proposals we also use 1000 samples and all four antithetic variables, i.e. we use Equation (3.45). At the MLE we again determine the EIS proposal using the same parameters and perform inference for the conditional distribution using 10000 samples, applying the method described in Section 3.6.2 to obtain estimates of the posterior mean, standard deviation and prediction intervals.

method	$\hat{\sigma}_\rho$	$\hat{\sigma}_W$	$\hat{\sigma}_q$	$\hat{\sigma}_M$	$\hat{\sigma}_{W_q}$
manual	0.001	0.100	0.50	0.01	0.10
initial	0.015	0.024	0.12	0.14	0.81
MLE	0.015	0.024	0.12	0.14	0.81

Table 4.1: Standard deviations for the models' showcase determined either by hand, by the initial search or by maximum likelihood estimation described in Section 3.6. The difference between the initial search and the MLE is negligible and is not visible for the precision shown here.

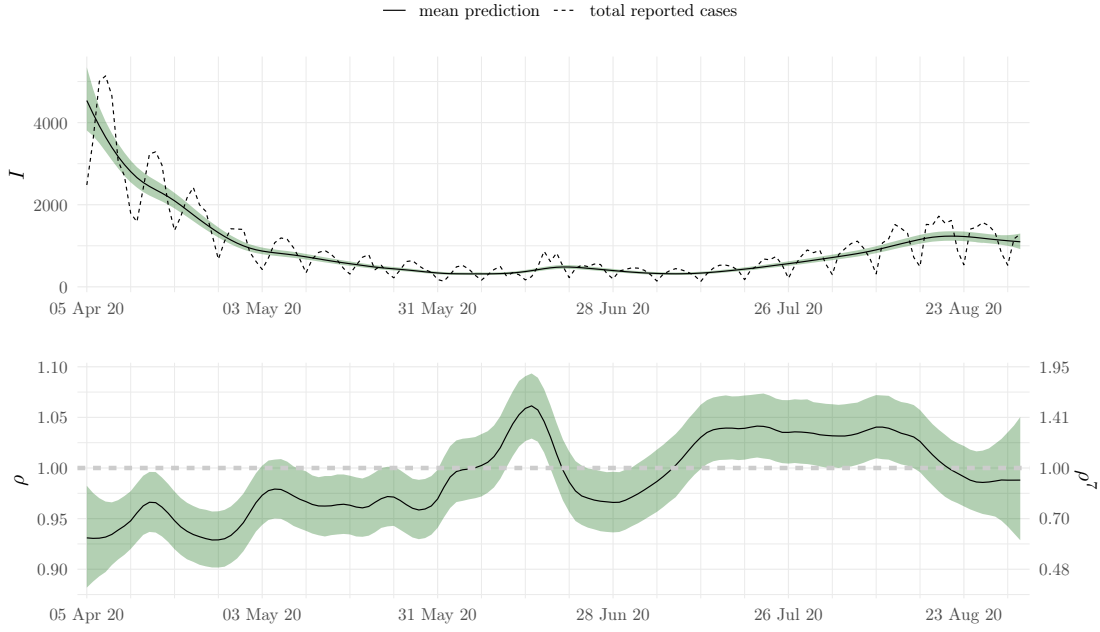


Figure 4.2: Monte-Carlo estimates of mean (black lines) and 95% prediction intervals (shaded green regions) for smoothed incidences  $I$  and daily growth factors  $\rho$ . The total reported cases with delay at most 4 days,  $Y_t$ , is shown as a dotted line. The secondary axis for the daily growth factor  $\rho$  indicates the corresponding weekly growth factors  $\rho^7$  which are easier to interpret. The gray dashed line indicates the threshold for growth  $\rho = 1$ .

### 4.1.3 Results

We start by a showcase of the models' capability, fitting it to the reported case date in the period from April 5th to September 1st 2020, starting from the first day when 4 delays are available in the dataset to the initial period of exponential growth in the fall of 2020. We estimate the parameters  $\theta = (\log \sigma_\rho^2, \log \sigma_W^2, \log \sigma_q^2, \log \sigma_M^2, \log \sigma_{W_q}^2)$  by maximum-likelihood estimation, yielding the parameters displayed in Section 4.1.3. There, we see that  $\log \rho_t$ ,  $\log W_t$ ,  $q_{t,1}$ ,  $q_{t,2}$  and  $q_{t,3}$  vary slowly over time, compared to the faster varying  $W^{q,1}$ ,  $W^{q,2}$ ,  $W^{q,3}$ .

We show importance sampling estimates of the mean and 95% prediction intervals of the conditional distribution of  $I$  and  $\rho$  (Figure 4.2) as well as  $W$ ,  $M$  and  $p$  (Figure 4.3), based on the procedure described in Section 3.6.2. For  $I$  we additionally show the total number of reported cases with delay at most 4 days,  $Y_t = Y_{t,1} + \dots + Y_{t,4}$ , as a sanity check. Indeed,  $I$  is a smoothed version of  $Y$ , which removes weekday-effects and small discrepancies in reporting, as these effects are captured by the  $W$  and  $M$  terms.

For the daily growth factor  $\rho$  we additionally display the corresponding weekly growth factors  $\rho^7$  on the secondary axis. We see that uncertainty for  $\rho$  is roughly constant over time, except close

to the beginning and end of the time period considered here. We see that until June 2020  $\rho$  is below 1, followed by a short skip above 1 during the local outbreak highlighted in ?? and a return to  $\rho < 1$  until beginning of July 2020. We will deal with this sudden increase and the following decrease more extensively in the following section. From the middle of July 2020 to the middle of August 2020,  $\rho$  is consistently above 1, with a slight dip at the end of August, before rising above 1 again. That cases are, or will be, rising exponentially is easier to infer from  $\rho$  compared to  $I$ , as  $\rho$ , or  $\rho^7$  for that matter, directly quantifies the increase in cases. Thus, this sustained period of exponential growth could have been a warning sign to policymakers of the buildup of infections in the population, which only became noticeable in the cases starting in October 2020.

For the muck term  $M$ , we see that is centered around 1 and allows capturing variation of the reported cases that is not captured by other terms in the model. As  $M$  follows a log Normal distribution, its variance is  $(\exp(\sigma_M^2) - 1) \exp(2(-\frac{1}{2}\sigma_M^2) + \sigma_M^2) = \exp(\sigma_M^2 - 1) \approx 0.02$ , so  $M$  has standard deviation  $\approx 0.12$  for the MLE from Section 4.1.3, consistent with Figure 4.3. As such, we expect the reported cases to vary around  $\pm 24\%$  on any given day, due to residual effects not captured by the weekday effect. We also investigated qq-plots of the mean predictions of  $M$ , which indicate that there might be some outliers, e.g. those around the local outbreak in June 2020, present. To improve the fit, we could replace the distribution of  $M$  by, e.g., a t-distribution with a low degree of freedom, allowing for heavier tails. The LA for such a model can still be found efficiently, see (Durbin and Koopman, 2012, Section 11.7.2), so the methods of this section are still applicable. However, we deem such a modification to be outside the scope of this thesis.

The weekday effect  $W$  exhibits the expected seasonal pattern: on Sundays, which are marked by the minor breaks in the figures' grid,  $W$  is below 1, while it is high for Tuesdays, Wednesdays and Thursdays. Over the period considered, this pattern is quite stable, with only slight changes over time:  $W$  is slightly larger for Mondays and Fridays at the end of the period compared to the beginning. By construction, we have  $\overline{\log W_t} = \frac{1}{7} \sum_{\tau=-3}^3 \log W_{t-\tau} \approx 0$  for all  $t$ , so Jensen's inequality suggests  $\bar{W}_t = \frac{1}{7} \sum_{\tau=-3}^3 W_{t-\tau} \gtrsim 1$ . However, the practical difference is small:

$$\frac{1}{7} \sum_{\tau=-3}^3 \mathbb{E}(W_{t-\tau}|Y) \approx 1.05,$$

for  $t = 3, \dots, n - 3$ , with small standard deviation. Consequently, we could correct  $I_t$  for the bias introduced by  $W$  by an increase of 5% (or, more precisely, consider  $I_t \bar{W}_t$ ).

Finally, for the delay probabilities, we compute both the signals probabilities, given by Equation (4.1), and a smoothed version, obtained by setting  $\log W_t^{q,\tau}$  to 0 in Equation (4.1). From Figure 4.3, we see that starting in the middle of April, reporting became faster, with a larger share of cases being reported with a delay of only a single day. While this seems to reverse at the end of the considered period, this is likely due to the reporting artifacts at the end of August, indicated by the large spike in  $p_{t,2}$ .

Now that we have seen an application of the model, we use it to demonstrate how easily we can incorporate missing or faulty observations. Recall from ?? the problem of reporting artifacts during the 2020 Christmas season. In Figure 4.4 we show undesirable effects of directly applying our model to the data in this period. In this figure, the red lines correspond to inferences made using all available observations, while turquoise lines correspond to inferences made where we remove all observations from December 19th 2020 until January 17th 2021, marked by the gray background in the figure.

As outlined in

ref correct section

, we can fit both models using the same methods, as we only have to replace the observation matrices  $B_t$  for missing dates  $t$  by zero matrices and the conditional distribution of  $Y_{t,\tau}|S_t$  by  $\delta_0$ , while replacing  $Y_{t,\tau}$  by 0 for  $\tau = 1, \dots, 4$ . In the approximating LAs and EIS proposals we set  $z_t$  and  $\Omega_t$  to the zero vector and matrix, respectively.

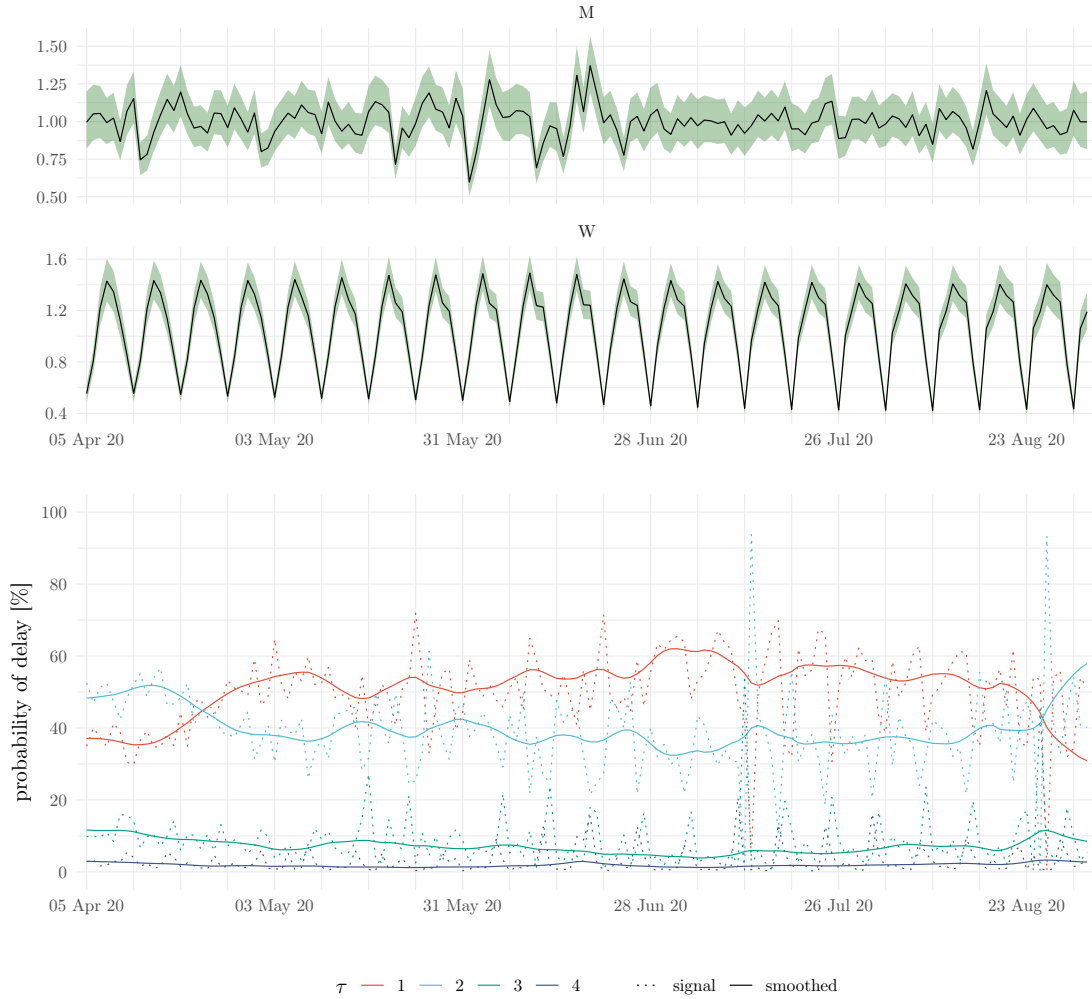


Figure 4.3: Importance sampling estimates of mean (black lines) and 95% prediction intervals (green ribbons) for weekday effect, „muck“ and delay probabilities in the showcase model, based on the method described in Section 3.6.2. We omit the small prediction intervals for delay probabilities for better readability. Note that all variables are not included directly in the model, but may be written as a function of states, either taking the exponential or converting from log-ratios to probabilities. The minor breaks in the x-axis grid indicate Sundays.

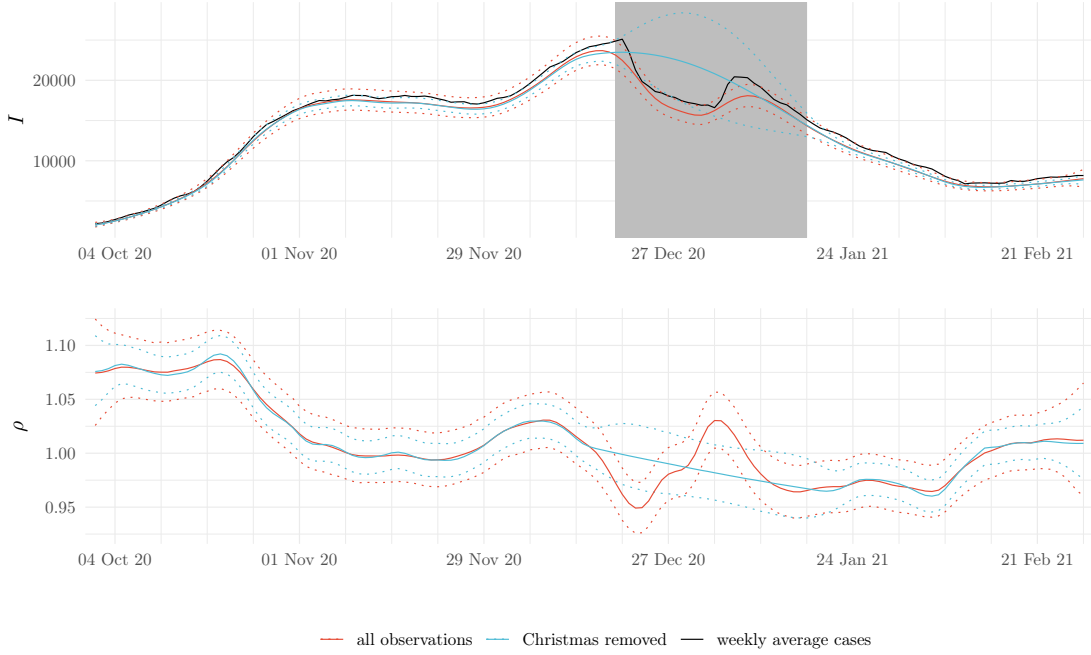


Figure 4.4: Importance sampling estimates of 95% prediction intervals and means of the conditional distribution of  $I$  and  $\rho$  given reported cases for the reporting delay model applied to the period of October 1st 2020 until February 28th 2021. For  $I$  we additionally show weekly average reported cases as in ??.

For the model using all available data, we see that the reporting artifacts affect both the incidences  $I$  and growth factors  $\rho$ , with a sharp decrease in  $\rho$  during the holidays, followed by a sharp increase in the new year. For the model that has the flawed observations removed, we see that both  $I$  and  $\rho$  behave more smoothly, as the estimated standard deviations, displayed in Table 4.2, are also smaller. The price we pay for this smoother transition is larger uncertainty where observations are now missing, i.e. the 95% prediction intervals are larger in this period than those for the model with all data available. However, when data are available, the prediction intervals for the second model are smaller, as its estimated standard deviations are smaller. The means, however, tend to agree rather well.

In Figure 4.5 we additionally show the expected smoothed delay probabilities based on Equation (4.1) where we set the weekday effects to 0. There, we see that starting on December 24th, the reporting pattern exhibits strong irregular behavior (recall that the reported cases for December 24th correspond to December 23rd to December 20th for delays  $\tau = 1, \dots, 4$ ) for the model using all observations. Additionally, in January, we a large spike in  $p_{t,1}$ , which could correspond to a backlog of cases being reported all at once. Again, the model that has the Christmas period removed, proceeds much smoother.

#### 4.1.4 Discussion

As we can see from the exemplary results, the model allows to accurately model the evolution of reported cases over time, while taking care of unwelcome reporting artifacts such as the weekday effect, delays and changes in reporting pattern due to holidays. The estimated growth factors allow inferring about the speed at which the cases proliferate, and can thus be a valuable tool for decision makers. With our model, we can identify the almost constant exponential growth in the summer of 2020 (Figure 4.2), which is difficult to see by only looking at the number of reported cases, due to the reporting artifacts and low number of cases.

method	$\hat{\sigma}_\rho$	$\hat{\sigma}_W$	$\hat{\sigma}_q$	$\hat{\sigma}_M$	$\hat{\sigma}_{W_q}$
<b>all observations</b>					
manual	0.0150	0.024	0.12	0.140	0.81
initial	0.0126	0.032	0.37	0.110	0.91
MLE	0.0126	0.032	0.38	0.110	0.91
<b>Christmas removed</b>					
manual	0.0150	0.024	0.12	0.140	0.81
initial	0.0087	0.028	0.16	0.048	0.38
MLE	0.0087	0.028	0.16	0.048	0.38

Table 4.2: Estimated parameters for the model during the Christmas period, for all observations or with observations during the Christmas period (19th December 2020 until January 17th 2021) removed. The manual parameter is based on the estimate of the models' showcase, i.e. the MLE result from Section 4.1.3.

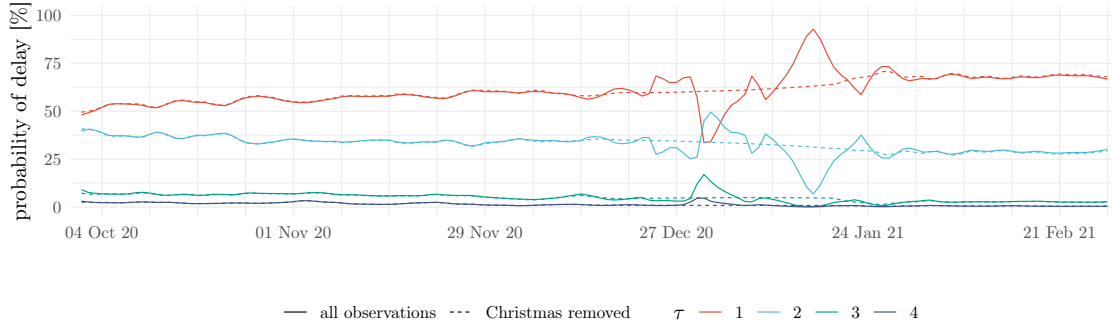


Figure 4.5: Importance sampling estimates of conditional expectation  $\mathbb{E}(p_{t,\tau}|Y)$  for the two sets of Christmas observations: the full lines correspond to all reported incidences and the dashed lines to observations between December 19th 2020 and January 17th 2021 removed. When using all observations, we clearly see the effect of the Christmas holidays, with a drop in next day reporting during the holidays, and an increase in next day reporting in the middle of January.

As our model explicitly models reporting delays, we can infer about them as well. From Figures 4.3 and 4.5 we can see that during the first year of the epidemic, reporting became faster: while only about 40% of cases were reported with delay of one day in April 2020, that fraction rose to more than 60% at the end of 2020, excluding the noisy Christmas period. By including reporting delays, our model is also capable of performing now- and forecasts of future reported cases. For forecasts, we evaluate the performance of seven day predictions from this model in the following section.

The SSM nature of our model has the additional advantage of being capable of naturally handling missing observations, either actual missing observations or synthetically missing observations, such as in the Christmas period. By removing available observations from the model, we are able to create a what-if scenario, letting the model automatically fill-in the faulty observations. Let us hasten to add that this should not be confused with redistributing the number of cases observed in the Christmas period to better fit the model, as we have not included any restrictions on the total number of cases being equal to the observed number of cases in this period. Technically, this is possible, by adding  $\log I_{t-1}, \dots, \log I_{t-D+1}$ , where  $D$  is the number of days removed, to the states and adding a single observation of  $\sum_{s=0}^{D-1} I_{t^*-s}$  at time  $t^*$ , the first day after the Christmas period is over. Removing the observations from December 19th 2020 to January 17th 2021 removes a total of 551 031 cases. In the Christmas model, the predictive distribution of cases for this time period has mean 618 000 (standard deviation 59,000) with a 95% prediction interval of (511 000, 743 000) (all numbers rounded to the next thousand to account for the Monte Carlo error). Thus our reconstruction of the total number of cases is compatible with the total number of cases removed, albeit slightly overestimating the total number of cases.

While we believe that our model already captures many of the relevant effects for modeling the daily evolution of cases, there are several worthwhile extensions conceivable. We here give an incomplete list of potential improvements:

- We have only used the reporting date in our model, but the data include also information (for some cases) on the symptom onset date. Including this would also allow to better remove the weekday effect, as infection dates, presumably, are less affected by weekdays than reported cases.
- In the same vein, including data on deaths would allow for estimates of the reporting dark figure, and its change over time, as long as immunization through infection and vaccination is low, i.e. at the beginning of the epidemic.

Most of these improvements require the use of additional data sources, which is straightforward to do with state space models: we just have to extend the states and dynamics accordingly.

ignores repeated Poisson noise -> need overdispersion

## 4.2 Regional growth factor model

### 4.2.1 Context

Modeling the epidemics spread on a regional level allows us to differentiate between localized and global outbreaks, such as the one in June 2020, highlighted in ???. Additionally, regional level prediction and growth factors are of interest on their own, because NPIs are enforced on the regional level. Moreover, having access to the spread on the regional level enables, e.g., regression of the growth rate against regional covariates, which in turn sheds light on which factors drive the epidemic.

Instead of modeling the number of cases per day and with delay as we did in Section 4.1, we will now model the total number of cases reported within one week for every county in Germany. Here we assume that a sufficient time period has passed, i.e. several days, see ??, such that the total number of cases is known sufficiently well. This weekly approach has several advantages: First, aggregating over the weekly data gets rid of the weekday effect, at the expense of a lower time resolution. Second, if we are interested in a retrospective analysis, it is sensible to assume all cases have been reported already, so we can avoid modeling the reporting delays.



However, modeling cases on the regional level comes with its own challenges, as we have to take care of accounting for the spatial spread, as well as an exchange between regions, cf. ??.

### 4.2.2 Model

Similar to the last section, we start by modeling the evolution of cases in time. We now have incidences  $I_{t,r}$  reported for reporting date  $t$  and region  $r$ , where there are a total of  $R$  regions.

comment on Gebietsreform in 2021 (?)

Again, we model the evolution of cases by

$$\log I_{t+1,r} \approx \log I_{t+1,r} + \log \rho_{t+1,r} \quad (4.3)$$

where  $\rho_{t+1,r}$  is the weekly growth factor in region  $r$ . Now we deviate from the previous model and model

$$\log \rho_{t,r} = \overline{\log \rho_t} + u_{t,r},$$

where  $\overline{\log \rho_t}$  is the average growth rate and  $u_{t,r}$  is the difference between the growth rate in region  $r$  and the country wide average. We will model  $u_{t,r}, r = 1, \dots, R$  to be jointly Gaussian, but correlated, which will enable us to model regional dependencies. To motivate our choice for the covariance structure, let us consider how cases are transferred between regions first.

As we are modeling cases on a regional level, we have to account for an exchange of cases as well. To illustrate our approach, suppose that we have for region  $r$   $S^r$  many secondary cases generated where the primary case belongs to region  $r$ , but the secondary case may belong to another region  $r'$ . Here „belonging to“ signifies that the case is reported in that region, which means that the infectee has registered their center of living to be in this region. Denote by  $p_{r,r'}$  the fraction of such cases and set  $p_{r,r} = 1 - \sum_{r' \neq r} p_{r,r'}$ .

Under these assumptions, the newly reported cases in region  $r$  are

$$\tilde{S}^r = \sum_{r'} p_{r',r} S^{r'} = (P^T S)_r$$

for  $P = (p_{r,r'})_{r,r'=1,\dots,R}$ . Assuming now that  $S^r, r = 1, \dots, R$  are random and i.i.d. with variance  $\sigma_S^2$ , we have

$$\text{Cov}(\tilde{S}) = \text{Cov}(P^T S, P^T S) = \sigma_S^2 P^T P.$$

However, modeling the correlation of newly reported cases turns out to be difficult: the cases will surely be modeled by a Poisson or Negative Binomial distribution, so we would have to decide on a copula to introduce this dependency structure. While this is feasible in principle, we opt for an easier way. Instead of modeling correlated incidences  $I_{t+1,r}$ , we model correlated growth rates  $\log \rho_{t+1,r}$ , by taking  $\text{Cov}(u_t)$  to be  $\sigma_S^2 P^T P$ . By Equation (4.3), conditional on  $I_{t,r}$ , this also captures regional correlation, without having to specify an involved joint distribution for the incidences.

As elaborated in ??, we want the regional effects  $u_{t,r}$  to be both flexible, but also, in some sense, stable over time. Thus, it makes sense to model  $u_t$  as a stationary process in time. The simplest, non-trivial, stationary process is a vector-autoregressive process

$$u_{t+1} = \alpha u_t + \varepsilon_{t+1,u}$$

where  $\alpha \in (-1, 1)$  and  $\varepsilon_{t+1,u} \sim \mathcal{N}(0, \Gamma)$ , where  $\Gamma$  is a positive definite matrix. By the above discussion, we set  $\Gamma = (1 - \alpha^2) \sigma_S^2 P^T P$  so that the stationary distribution of  $u_t, t = 0, \dots, n$  is  $\mathcal{N}(0, \sigma_S^2 P^T P)$ .

To setup our SSM, let  $X_t = (\overline{\log \rho_{t+1}}, u_{t,1}, \dots, u_{t,R})^T \in \mathbf{R}^{R+1}$ . For the observations, we let  $Y_t = (I_{t,1}, \dots, I_{t,R})^T$ , the number of cases observed in regions  $1, \dots, R$  in the  $t$ -th week.

We then model the number of cases at time  $t + 1$  in region  $r$ ,  $I_{t+1,r}$  to follow a negative binomial distribution, conditional on the states  $X_t$  to be

$$I_{t+1,r} | I_t, \overline{\log \rho_t}, u_{t,r} \sim \text{NegBinom}(\bar{\rho}_t \exp(u_{t,r}) P^T I_t, r),$$

conditionally independent. While the previous observations  $I_t$  are now conditioned on as well, recall from our discussion in the beginning of Chapter 3, that this is not problematic.

To fully specify the model, we have to provide the transfer probabilities  $p_{r,r'}$ . For these, we use official data by Germany's federal employment agency on commuters

ref

. From these data, we calculate  $q_{r,r'}$ , the fraction of socially insured employees that have their center of life in region  $r$ , but are registered to work in region  $r'$ . As this is only a crude approximation to the actual exchange between regions, we let

$$p_{r,r'} \propto \bar{q} + (1 - \bar{q}) \frac{q_{r,r'}}{\sum_{r'' \neq r'} q_{r,r''} + C q_{r,r}}$$

where we interpret  $\bar{q}$  as a constant socket of exchange between regions and  $C \geq 1$  as an additional proportion of stay at home inhabitants that are not captured by  $q_{r,r}$ , e.g. elderly or children.

Thus, our final model is parameterized by

$$\theta = \left( \log \sigma_S^2, \text{logit } \alpha, \log(C - 1), \text{logit } \bar{q}, \log \sigma_{\log \rho}^2, \log r \right),$$

where chose a parametrization that is unconstrained. The model has a linear signal

$$S_t = (\log \rho_t + u_{t,r})_{r=1, \dots, R},$$

which makes inference fast, as the approximating GLSSM in the EIS method only requires  $\mathcal{O}(nR)$  many parameters. Again, we use MLE to estimate  $\theta$ , using the methods from Section 3.6.

### 4.2.3 Results

- fit model by MLE + show inference for Toennies outbreak, interpret covariance matrix estimates
- predict incidences on regional level and show that we outperform simple Poisson / NB baseline that only uses a single region
- maybe: perform predictions for 1-4 weeks ahead, compare to regional FCH

### 4.2.4 Discussion

consider Armbruster2024Networkbased, Armillotta2023Inference

## 4.3 Nowcasting hospitalizations

compare SSM predictions to FCH submissions

### 4.3.1 Context

Judging the severity of the COVID-19 epidemic has been an ongoing challenge since its inception. As immunization against COVID-19 rose, strict enforcement of social distancing rules eased and testing regimes became less strict, case incidences became a less reliable and harder to interpret indicator of epidemic severity. Instead more direct indicators of morbidity, such as the number of deaths and ICU admissions and occupancy have come to the fore. But these indicators are late due to the substantial delays between infection and occurrence. An alternative indicator that captures



Figure 4.6: **TODO: redo figure with final model** Germany’s 7-day hospitalisation incidence changes due to various delays such as time to hospitalisation and delays in reporting. This figure shows the extent of these delays: incidences reported at the present date (red lines) severely underestimate the hospitalisation incidence (green solid lines) that is reported after 3 months. Our nowcasting model (blue dotted lines, 95% prediction intervals in shaded gray) deals with this problem by predicting the hospitalisation incidence based on past cases and their delays to hospitalisation.

the morbidity caused by COVID-19 but is earlier than the others is the number of hospitalisations of positive COVID-19 cases.

While hospitalisations occur earlier, they still come with substantial delay between the infection and subsequent admission to hospital. Additional difficulties arise due to delays in reporting, i.e. the time it takes until the hospital reports the new case to the national health authorities. The problem of accounting for delays in reporting for occurred, but not yet reported events has been termed **nowcasting**, i.e. forecasting of the indicator at time “now”. Predicting the number of hospitalisations is thus a mixture of both forecasting — which reported COVID-19 cases will end up in the hospital — and nowcasting — which cases have yet to be reported — and we will use the term nowcasting in this paper to mean this predictive mixture. In this section we focus on the situation in Germany where data on hospitalisations has been available since April 2021 provided by the German federal health care authority, the Robert Koch-Institut (RKI), via Github (Robert Koch-Institut, 2021).

Compared to other approaches in the COVID-19 NowcastHub, that tended to exclusively focus on modelling the delay distribution with parametric and non-parametric models, our model sidesteps this complex delay structure by decomposing delayed hospitalisations into weekly chunks (??) and incorporating case data. As cases and hospitalisations are explicitly linked by the case reporting date we forecast the number of hospitalisations in each chunk based on the current incidences and past fractions of hospitalisations in a comparable weekly chunk. We additionally quantify uncertainty by prediction intervals that are informed by the past performance of our model. This makes our model straightforward to understand, easy to implement and fast to run.

reformulate

The origin of nowcasting lie in accounting for incurred, but not reported claims in the actuarial sciences (Kaminsky, 1987), delays in reporting for AIDS (Lawless, 1994; Zeger, See, and Diggle, 1989) and other infectious diseases (Farrington et al., 1996). Popular statistical approaches include methods from survival analysis (Lawless, 1994) and generalized linear regression (Zeger, See, and Diggle, 1989). In the survial analysis setting one commonly models the reverse time discrete hazard parametrically and assumes multinomial sampling of the final number of cases, potentially accounting for overdispersion. This has been studied with frequentist (Midthune et al., 2005) and Bayesian (An Der Heiden and Hamouda, 2020; Höhle and An Der Heiden, 2014) methods. The generalized linear regression approach has origins in the chain ladder model from actuarial sciences (Renshaw and Verrall, 1998) and models the observed counts in the reporting triangle by a Poisson or negative binomial distribution. For both approaches, available covariates can be incorporated in a straightforward way. In the setting of real-time nowcasting, it is often beneficial to incorporate epidemic dynamics into the model, this can be achieved by splines (Höhle and An Der Heiden, 2014;

van de Kasstele et al., 2019) or by a latent process of infections (McGough et al., 2020).

Nowcasting methods have wide application in accounting for reporting delays (Midthune et al., 2005), early outbreak detection (Bastos et al., 2019; Salmon et al., 2015), and, in the recent COVID-19 epidemic, improving real-time monitoring of epidemic outbreaks (Akhmetzhanov, 2021; An Der Heiden and Hamouda, 2020; Günther et al., 2021; Schneble et al., 2021). Evaluating a forecasting model in a real-time public health setting is advantageous as it avoids hindsight bias (Desai et al., 2019), however nowcasting approach may have difficulties with bias and properly calibrated uncertainty if used in a real-time setting. This includes rapidly changing dynamics (Günther et al., 2021; van de Kasstele et al., 2019), both of the delay distribution and the underlying epidemic, retrospective changes in data (Midthune et al., 2005) and long delays with few observed cases (Noufaily et al., 2015).

To avoid the aforementioned hindsight bias one can make their predictions publicly available in real-time (Bracher et al., 2021; Ray et al., 2020). For the hospitalisations in Germany, we have participated in the German COVID-19 NowcastHub (*Nowcasts Der COVID-19 Hospitalisierungsinzidenz* 2022) since November 2021 where nowcasts are available in a public Github repository (*Hospitalization Nowcast Hub* 2022) with the “ILM-prop” model. The ideas, especially the model and the “double-weekday effect”, discussed in this section are based on this model. However, the “ILM-prop” model is based on simple point estimates for the proportion of hospitalisations per reported case, neglecting regularization over time. In this thesis we extend this model to the SSM setting of this thesis and investigate if the increased model complexity results in improved performance. In particular, we want to reduce computation time, as the previous model quantified uncertainty by past model performance, which requires running the model many times. If prediction uncertainty is based on predicting future observations in a SSM, we can reduce computation time drastically. However, this is only worthwhile, if the predictive performance is comparable to the computationally more intensive model.

To predict the number of hospitalisations we consider the reporting process of both reported COVID-19 cases and reported hospitalisations. Recall that the reporting date of a COVID-19 case is shared for both the case and its hospitalisation, i.e. the case and hospitalisation are linked through this date.

As hospitalisations are only available as 7-day rolling sums, we use 7-day rolling sums for daily reported incidences as well. To avoid dealing with the double weekday effect of both reporting date of the case and reporting date of the hospitalisation (see ??) we divide the future hospitalisations we wish to predict into chunks of one week, which gets rid of the weekday effect for the hospitalisations. This is depicted in ?. Our prediction of each of these weekly chunks then consists of the fraction of hospitalisations of reported cases in the past.

We use the publicly available data from the RKI discussed in ?? on daily reported COVID-19 cases (Robert Koch-Institut, 2024b) and weekly reported hospitalizations (Robert Koch-Institut, 2024a). Both datasets are updated daily.

Recall from ?? that COVID-19 cases are described by their date of reporting, and are subject to reporting delay and hospitalizations are reported by the *reporting date of the associated case*, and are subject to delay as well. As the date of symptom onset is not known for a substantial amount of incident cases, and is not reported for hospitalized cases, we focus our analysis on the date of reporting.

### 4.3.2 Model

In line with the structure of the data, we let  $H_{t,t+\tau}^a$  be the number of weekly hospitalizations in age group  $a$  with case reporting date  $t - 1, \dots, t - 7$  that are known on the day  $t + \tau$ , aggregated over all states. We suppress the dependence on age group in the following for ease of notation, but all modeling is to be performed for every age group separately.

As we focus on same-day nowcasting, our goal is to predict on day  $t$   $H_{t,t+D}$  the number of hospitalizations reported  $D$  days into the future, for simplicity assume that the maximal delay

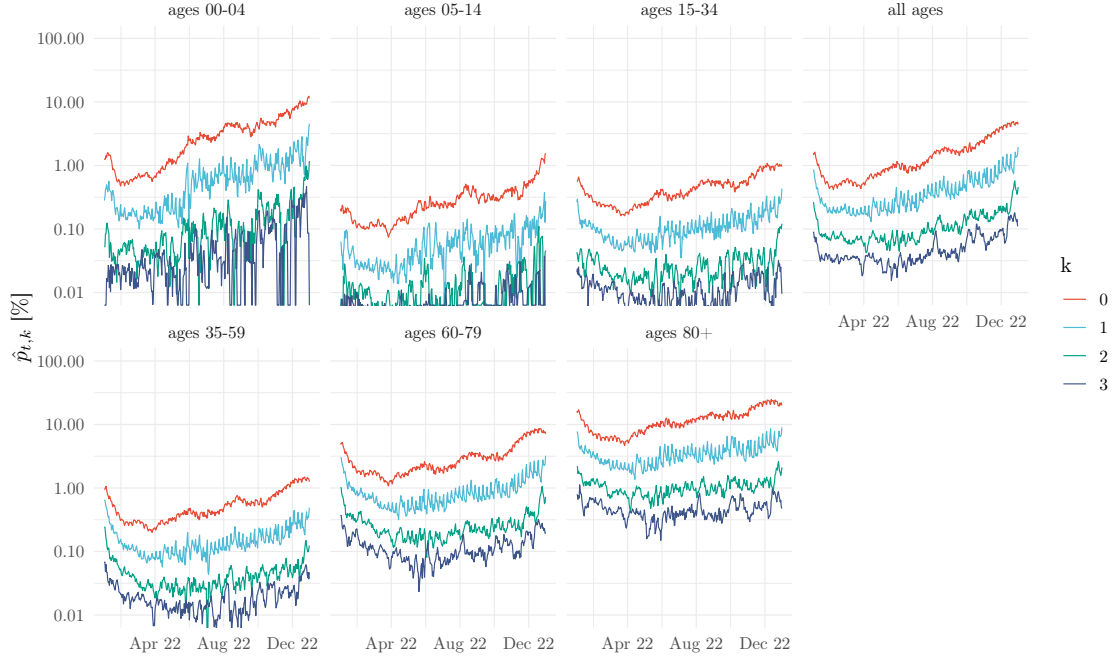


Figure 4.7

considered,  $D$ , is a multiple of 7. We decompose this target into a weekly telescoping sum

$$H_{t,t+D} = H_{t,t} + \sum_{k=1}^{D/7} (H_{t,t+7k} - H_{t,t+7(k-1)}),$$

$q_{t,0}$

where  $H_{t,t}$  is already known on day  $t$  and  $H_{t,t+7k} - H_{t,t+7(k-1)}$  is the increment in the hospitalization incidence from the  $(k-1)$ -st week to the  $k$ -th week. Recalling that any case attached to the hospitalization incidence on this date has case reporting date  $t$  we now crucially assume that the hospitalization reporting process consists of two independent events: hospitalization and its delayed reporting. More formally, let  $I_t^7$  be the seven day case incidence (again, modeled separately for every age group) on day  $t$ , defined in the same fashion as the hospitalization incidence. Thus

$$I_t^7 = \sum_{\tau=1}^7 I_{t-\tau,t}$$

where for  $\tau = 1, \dots, 7$   $I_{t-\tau,t}^7$  is the number of cases with reporting date  $t - \tau$  known on date  $t$ . Note that, similar to the hospitalization incidence,  $I_t^7$  does not contain cases with reporting date  $t$ , but rather cases with reporting dates  $t - 1, \dots, t - 7$ . While cases are also affected by reporting delays, these delays are on the order of days, rather than weeks, cf. ??, and averaging over the past week means that  $I_t^7$  is subject to only minor, negligible, reporting delays. We thus model

$$H_{t,t+7k} - H_{t,t+7(k-1)} | I_t^7, p_{t,k} \sim \text{Poisson}(\lambda_{t,k}) \quad \lambda_{t,k} = I_t^7 p_{t,k}, \quad (4.4)$$

conditionally independent for all  $t$  and  $k$ . Here  $p_{t,k}$  is the proportion of reported cases  $I_t^7$  that will become hospitalized after  $k$  weeks. For simplicity of notation, let  $H_{t,t-7} = 0$ , so that  $H_{t,t} - H_{t,t-7} = H_{t,t}$  has conditional rate  $\lambda_{t,0} = I_t^7 p_{t,0}$ .

Figure 4.7 displays the empirical delay probabilities  $\hat{p}_{t,k} = \frac{H_{t,t+7k} - H_{t,t+7(k-1)}}{I_t^7}$  during 2022 on the log scale for small  $k$ .

Ignoring noisy day-to-day variation, we see that within an age group, the delay probabilities evolve roughly in parallel. This encourages us split the delayed hospitalization probabilities  $p_{t,k}$  for all  $t$  and  $k$  into two parts

$$p_{t,k} = p_t q_{t,k}$$

where  $p_t$  is the time-varying proportion of hospitalization and  $q_{t,0}, \dots, q_{t, \frac{D}{7}}$  comprise the delay distribution. To make this identifiable we impose that  $\sum_{k=0}^{\frac{D}{7}} q_{t,k} = 1$ . Figure 4.7 implies that the delay probabilities evolve rather smoothly, so we let the log-probabilities  $\log p_t$  perform a second order random walk, i.e. we model

$$\begin{aligned} \log p_{t+1} &= \log p_t + v_t \\ v_{t+1} &= v_t + \varepsilon_{t+1,v}. \end{aligned}$$

For the delay distribution, we first reparameterize to consecutive conditional probabilities

$$q_{t,k}^c = \frac{q_{t,k}}{1 - \sum_{l=1}^{k-1} q_{t,l}},$$

i.e.  $q_{t,k}^c$  is the probability of a delay of exactly  $k$  weeks, conditional on having at least  $k$  weeks of delay. This reparameterization is a diffeomorphism from the open simplex  $\{p \in \mathbf{R}_{>0}^{D/7+1} \mid \|p\|_1 = 1\}$  to  $(0, 1)^{D/7} \times \{1\}$  which has the advantage that  $q_{t,k}$  only depends on  $q_{t,l}^c$  for  $l \leq k$  (rather than all of them, as was the case for the model in Section 4.1). We then model the logits of these reparameterized delay probabilities to perform independent random walks

$$\text{logit } q_{t+1,k}^c = \text{logit } q_{t,k}^c + \varepsilon_{t+1,q,k}$$

for  $\varepsilon_{t+1,q,k} \sim \mathcal{N}(0, \sigma_q^2)$ .

From Figure 4.7 we additionally observe a weekday effect, at least for small  $k$ . Thus, we additionally add two multiplicative weekday effects for  $q_0^c$  and  $q_1^c$ , i.e. we modify (4.4) to be

$$\lambda_{t,k} = I_t p_t q_{t,k} W_{t,k} \quad \text{for } k = 0, 1, \quad (4.5)$$

qs are now on logit scale

where  $W_{t,0}$  and  $W_{t,1}$  are two independent, multiplicative weekday effects as for the model in Section 4.1. The choice of having two weekday effects here is based on balancing the dimension of the model, and thus the computational resources required to run inferences and predictions, with its explainability and is based on numerical experiments. A more rigorous analysis, e.g. using information criteria, could be run as well, but is outside the scope of this thesis.

As always, we assume that the innovations

$$\begin{aligned} \varepsilon_{t+1} &= (\varepsilon_{t+1,v}, \varepsilon_{t+1,q,0}, \dots, \varepsilon_{t+1,q,D/7-1}, \varepsilon_{t+1,W,0}, \varepsilon_{t+1,W,1}) \\ \text{Cov}(\varepsilon_{t+1}) &= \Sigma_{t+1} = \text{diag}(\sigma_p^2, \sigma_q^2, \dots, \sigma_q^2, \sigma_W^2, \sigma_W^2) \end{aligned}$$

are centered, independent across all  $t$  and Gaussian.

These considerations lead to a PGSSM with linear signal

rethink this term once again

. Let the states and signals be given by

$$\begin{aligned} X_t &= (\log p_t, \text{logit } q_{t,0}, \dots, \text{logit } q_{t,D/7}, \log W_{t,0}, \dots, \log W_{t-5,0}, \log W_{t,1}, \dots, \log W_{t-5,1})^T \\ S_t &= (e_1 \quad \mathbf{0}_p \quad e_2 \quad \dots \quad e_p \quad e_2 \quad \mathbf{0}_{p \times 5} \quad e_3 \quad \mathbf{0}_{p \times 5}) X_t \\ &= \left( \log p_t, \text{logit } q_{t,0} + \log W_{t,0}, \text{logit } q_{t,1} + \log W_{t,1}, \text{logit } q_{t,2}, \dots, \text{logit } q_{t, \frac{D}{7}-1} \right)^T \end{aligned}$$

and let the observations be

$$Y_t = (H_{t,t}, H_{t,t+7} - H_{t,t}, \dots, H_{t,t+D} - H_{t,t+D-7})^T$$

with conditional distribution given by Equation (4.4) where  $\lambda_{t,k}$  is given by Equation (4.5).

Notice that we do not need to specify the evolution of cases over time, as we are interested only in nowcasting hospitalizations for the current day. On day  $t$  there are many missing observations, in particular, we only observe the first component of  $Y_{t-6}, \dots, Y_t$ , only the first two components of  $Y_{t-13}, \dots, Y_{t-7}$  and so on. These can be dealt with as in the Gaussian case, setting the corresponding rates  $\lambda_{t,k}$  manually to 0 and replacing the missing observations by 0. For the approximating GLSSMs, we use the missing data strategy discussed in ??

write this here

.

For the initial distribution of use  $X_0 \sim \mathcal{N}(\mathbb{E}X_0, \Sigma_0)$  where

$$\mathbb{E}X_0 = \begin{pmatrix} \log p_0 \\ 0 \\ \dots \\ 0 \end{pmatrix}$$

and  $\Sigma_0 = \sigma_0^2 I$  is a multiple of the identity matrix. We chose these initial conditions as they only introduce two further unknown parameters, making them amenable to maximum likelihood estimation. Of course specifying the same (large) variance  $\sigma_0^2$  for all states may simultaneously over- and under-estimate the initial variance in some components. As an alternative, one could implement the diffuse initialization of the Kalman filter, see (Ansley and Kohn, 1985; Koopman, 1997).

The model is parameterized by

$$\theta = (\log \sigma_p^2, \log \sigma_q^2, \log \sigma_W^2, \log \sigma_0^2, \log p_0)$$

which we estimate by MLE.

To fit the model for all age groups, we use at most 100 iterations for all occurrences of the LA, with a convergence threshold set to  $10^{-5}$  relative difference in  $z$  and  $\Omega$ . We use the same method for EIS, where we additionally use 1,000 samples to determine the optimal proposal, starting with the proposal given by th LA.

For MLE, we use 1,000 samples to determine the maximum likelihood estimate of the parameters, initializing at the initial guess given by Section 3.6.1.

To obtain prediction intervals of the states, signals and missing observations we use 10,000 samples. To estimate the ESS we use 10,000 samples.

### 4.3.3 Results

To demonstrate the capabilities of our model, we fit it the analysis period of the NowcastHub, i.e. to the period from 22nd November 2021 to 29th April 2022 (Wolfram et al., 2023). We do this for each of the seven age groups, including all age-groups together, the 00+ age group. For each of these age groups, we fit the model as described in last subsection.

In the original NowcastHub the truth against which nowcasts were to be evaluated was set to the data available 100 days after the last date of the study period, 8th August 2022, under the assumptions that there would be no late reporting after such a long period. However, it turns out that the hospitalization incidence is still subject to relevant data revisions long after the date of reporting of the COVID-19 case has passed see (Wolfram et al., 2023, Section 3.7). Thus, we choose to focus on a delay of 42 days (6 weeks), similar to the alternative time horizon of 40 days proposed as an alternative target in the same section.

Age group	EF [%]	weeks of delay
A00-04	61	5
A05-14	5	5
A15-34	74	7
A35-59	88	7
A60-79	93	8
A80+	98	8
A00+	97	8

Table 4.3: Efficiency factors (in %) and weeks of delay for the seven models (one per age group) presented in this section. For younger age groups, there are few long delays, which causes numerical instabilities due to the consecutive conditional probability parametrization chosen in this section. For each of the age groups, we chose the longest delays that still allowed for a reasonable fit, with a maximum delay of 8 weeks. While the efficiency factor for A05-14 is quite low, we use a large enough number of samples for the prediction of states and signals, so the ESS is still sufficiently large.

For the younger age groups, long delays are rare (see also Figure 4.7), which leads to numerical instabilities in the consecutive logit parametrization. If such numerical instabilities occur, we manually choose the maximum delay (in weeks) that still produces a reasonable model fit, which we define here as an EF above 5%. We show the resulting weeks of delay and EF in Table 4.3. The resulting posterior distributions of interest are displayed in Figure 4.8.

We see that, generally, hospitalization probabilities  $p_t$  grow larger as the age group under consideration becomes older; note the logarithmic  $y$ -axis. The exception here is the youngest age group A00-04. While infants are vulnerable to COVID-19 (Havers et al., 2024), this may also be explained by circumstantial testing in hospitals: children in age group A05-14 were largely subjected to mandatory testing at school, so we would expect the darkfigure of unreported cases in age group A00-04 to be large compared to the older age groups. As always, we stress that interpretations of our results are contingent on taking the considerations from ?? into account. Nevertheless, we see  $p_t$  drop in all age groups, except A00-A04, over the period considered. This is consistent with the rise of the Omicron variant of SARS-CoV-2 (Robert Koch-Institut, 2024c) which is associated with milder progression of disease.

We also observe a pronounced weekday effect  $W_{t,0}$  across all age groups, with a smaller proportion of  $I_t^7$  reported as hospitalized already on day  $t$  if  $t$  is a Sunday, as indicated by the vertical grid lines in Figure 4.8. To compensate,  $W_{t,1}$  is large when  $W_{t,0}$  is small. Again, A00-04 exhibits a more pronounced weekday effect, but the general pattern is consistent across all age groups. On the right-hand side of Figure 4.8 we see the delay probabilities  $q_{t,k}$  for  $k = 0, \dots, 3$ .

interpret, wait for mean delay

.

interpret table:hospita. ess after final results

To evaluate the predictive capabilities of our model, we use it to perform retrospective nowcasting of hospitalizations, emulating the setting of the German NowcastHub. We focus on same-day nowcasting, i.e. only nowcasting for the current day, with a maximum delay of 6 weeks, performing all predictions for every age group separately.

For every day,  $s$  say, in the period of 22nd November 2021 to 14th April 2022 we fit the model to the data of the past 50 days that were available on day  $s$ . We exclude the Easter period starting with Good Friday on 15th April 2022 to avoid having to deal with data artifacts with inconsistent reporting in that period.

Thus, the observations of the model consist of  $y_t$  for  $s - 100 < t \leq s$ , but as

$$y_t = (H_{t,t}, H_{t,t+7} - H_{t,t}, \dots, H_{t,t+D} - H_{t,t+D-7}),$$



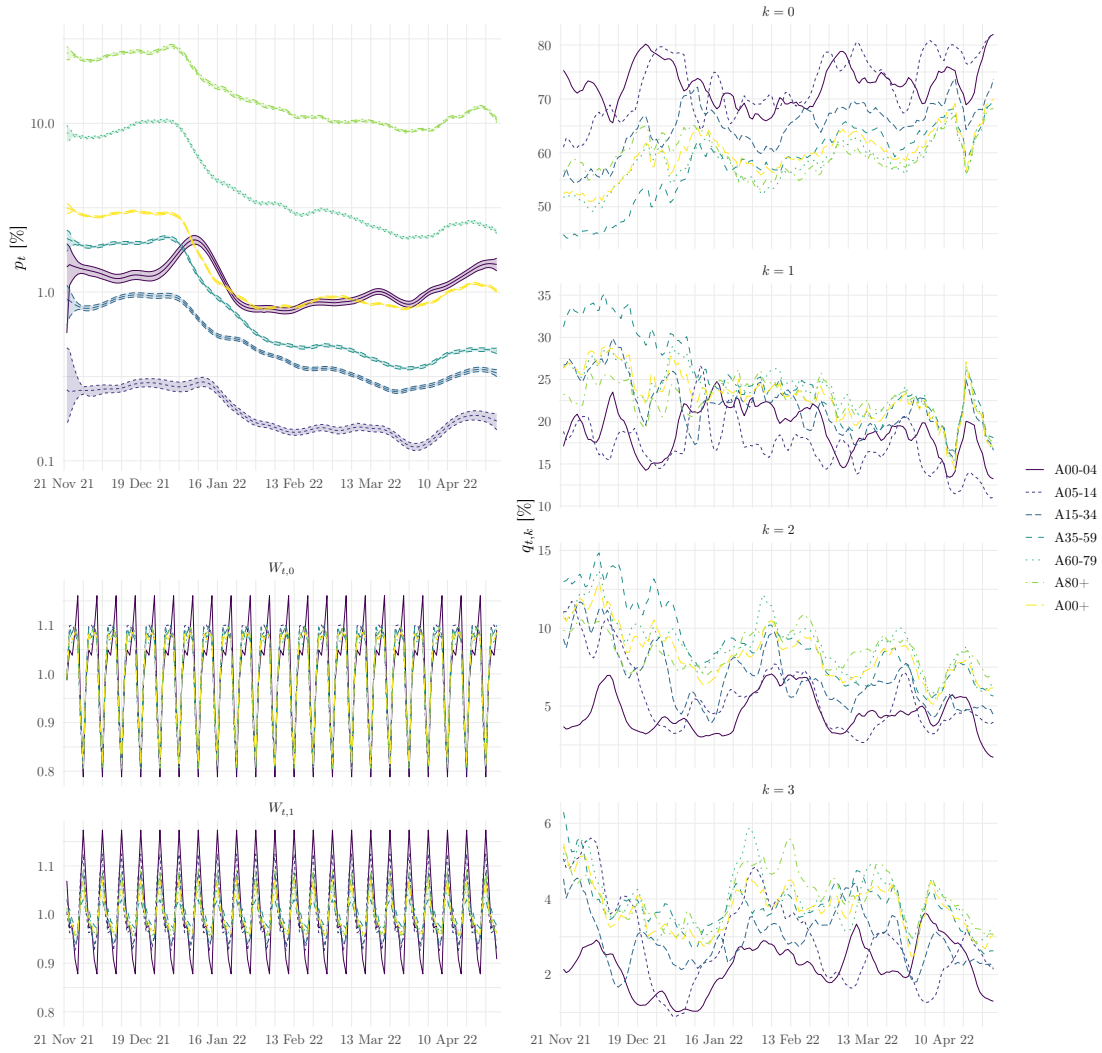


Figure 4.8: For each of the seven age groups (indicated by color and linetype), we show means of the smoothing distribution for the first four delay probabilities  $q_{t,k}$ ,  $k = 0, \dots, 3$ , the smoothed probabilities of hospitalization  $p_t$  and the two weekday effects  $W_{t,0}, W_{t,1}$ . Recall that we fit a separate model for each age group. For the smoothed delay probabilities, we additionally show 95% prediction intervals. Note the log-scale of the  $y$ -axis for the smoothed delay probabilities.

the  $k$ -th component of  $y_t$  is missing whenever  $t + 7(k - 1) > s$ . Taking the last day,  $s$ , as an example  $H_{s,s}$  is made available to the model, but  $H_{s,s+7k}$  for  $k > 0$  is not, and so  $y_s = (H_{s,s}, \mathbf{NA}, \dots, \mathbf{NA})$  where  $\mathbf{NA}$  indicates missing observations. Similarly, the last observation for which the second component is available is  $y_{s-7} = (H_{s-7,s-7}, H_{s-7,s} - H_{s-7,s-7}, \mathbf{NA}, \dots, \mathbf{NA})$  and so on.

Similar to the other models in this thesis, we can include these missing observations in a straightforward manner, by setting the corresponding rows of  $B_t$  to 0, setting the same entries of  $s_t$  to  $-\infty$ , such that  $\lambda_{t,k} = 0$ , and replacing missing observations by 0. For the approximating GLSSMs, we fix the rows and columns of  $\Omega_t$  and entries of  $z_t$  that correspond to missing signals  $s_t$  to 0. To make fitting the 158 resulting models computationally feasible, we omit the MLE step and fit the model using only the initial value from Section 3.6.1.

To nowcast the total number of hospitalizations, we use the method described in Section 3.6.2, i.e. using MC-integration to estimate quantiles of  $H_{s,D}$ . Accordingly, we draw  $N$  samples from the smoothing distribution  $S_s^i | Z = z$  with weights  $W^i$  and, conditional on these samples,  $\tilde{Y}_s^i | S_s^i \sim \text{Poisson}(\exp S_s^i)$ , independent of everything else, where we fix the first component be the known  $y_{s,1} = H_{s,s}$ . In total, we obtain  $N$  draws  $H_{s,D}^i = H_{s,s} + \sum_{k=1}^{D/7} \tilde{Y}_{s,k}^i$  with associated weights  $W^i$  from which we can estimate the desired quantiles. We use the same quantiles as in the NowcastHub, i.e. the 2.5%, 10%, 25%, 50%, 75%, 90%, 97.5% quantiles.

To assess the predictive performance of our model, we compare its predictions to the ILM-prop42 and the revised ensemble model, tailored to a maximum delay of 40 days from (Wolffram et al., 2023, Section 3.7).

- explore model fit & nowcasting for NCH period
- compare to same-day nowcasts provided by other models in the NCH (w/ WIS as performance metric)
- discuss usefulness of indicator vs. actual hospitalizations

#### 4.3.4 Discussion

# Chapter 5

## Discussion

Let us look back on the insights gained in this thesis and put the theoretical and applied results into a broader picture. To this end, let us give answers to the questions raised in Chapter 1 and see how far we have come. First, let us reiterate the ambitious goal of this thesis detailed in the introduction: provide mathematically sound statistical tools that allow practitioners to design and fit models suited for the epidemiological data arising in an ongoing epidemic. To this end we have contributed results to three areas, from more applied to more abstract:

- the results of concrete models applied to COVID-19 in Germany in Chapter 4,
- the general modeling strategy of PGSSMs for epidemiological models in the same chapters, as well as ??, and
- the rigorous mathematical study of the computational methods used to fit these models, especially the comparison of EIS and the CE-method in Chapter 3.

reformulate locality of epidemics more generally, effects particular to question at hand

In Chapter 2 we highlighted the need for explainable statistical models in infectious disease epidemiology. Throughout this thesis we use COVID-19 as a driving example, which is owed to the fact that this thesis is a product of this particular epidemic. However, the core questions surrounding past, current and future development of cases and derived indicators is of essence for all epidemics outbreaks, including seasonal influenza and other respiratory diseases

cite hubs, fraser influenza models

. As a solution to this challenging task, we presented SSMs, in particular PGSSMs, as a flexible framework for modeling the high-dimensional time-series data we are faced with. The wide range of applications in Chapter 4 demonstrates the flexibility of SSMs as a modeling tool.

We have seen that the widely available data surrounding the epidemic, in particular daily reports on infections and hospitalizations, can be leveraged to fit PGSSMs and the fitted models allow for straight-forward interpretation. While the fitting procedures for these models can be quite involved, the models themselves, especially the temporal dynamics, are not, as such our analyses can be disseminated to statisticians and practitioners alike. However, the interpretability of all models presented in this thesis is hampered by the quality of data available. Indeed, large constituents of our models are purely for dealing with weekday effects, and reporting delays.

The final contribution of this thesis is a mathematically rigorous analysis of the performance of the importance sampling methods used in the applications (Chapter 4). To the authors' knowledge, Theorems 3.6 and 3.9 gives the first joint analysis of the CE-method and EIS, and sheds insight on the poor performance of the CE-method in practice: we can expect its asymptotic covariance to be larger than that of EIS, as its meat matrix  $M$  is fixed, whereas that of EIS can be expected to be small when the optimal proposal is close to the target.



# Appendix A

## Reproducibility and code

All code used in to create figures and tables in this thesis is written in Python and R

cite

and available as open source software. Python is used for simulations, while R is used to create figures and tables of these results.

The code is split into two software packages:

- Importance Sampling for State Space Models (**isssm**)<sup>1</sup> is a Python package developed by the author. It implements frequentist inference for SSMs using the general methods described in this thesis, in particular the CE-method and EIS for PGSSMs.
- The SSMs for Epidemiology **ssm4epi** package contains Python and R code particular to this thesis, i.e. the code needed to reproduce all results and figures in this thesis.

The **ssm4epi** package is available as Jupyter Notebooks organized by chapters of this thesis. To reproduce the results of this thesis, follow the instructions in the associated documentation

ref to doc

. Simulations use a fixed seed that is set at the beginning of each notebook to ensure reproducibility.

The data produced by these Jupyter notebooks are available on zenodo

put them there

, and can be reproduced by running the notebooks. Figures and tables in this thesis that depend on simulation results can be reproduced similarly, using Jupyter notebooks with an R kernel. Dependent R packages can be found in the **setup.R** file in this thesis' GitHub repository.

---

<sup>1</sup><https://stefanheyder.github.io/isssm>



## Appendix B

# Additional calculations

### Equation (3.10)

To show Equation (3.10) we calculate the second moment of  $w(X)X$ ,

$$\begin{aligned}\mathbb{E}(w(X)X)^2 &= \int w(x)^2 x^2 g(x) \, dx \\ &= \int \sigma^2 \exp\left(-x^2 \left(1 - \frac{1}{\sigma^2}\right)\right) x^2 \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{x^2}{2\sigma^2}\right) \, dx \\ &= \int \sigma x^2 \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2} \left(2 - \frac{2}{\sigma^2} + \frac{1}{\sigma^2}\right)\right) \, dx \\ &= \int \sigma x^2 \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2} \frac{2\sigma^2 - 1}{\sigma^2}\right) \, dx \\ &= \tau \sigma \int x^2 \frac{1}{\sqrt{2\pi\tau^2}} \exp\left(-\frac{x^2}{2\tau^2}\right) \, dx \\ &= \tau^3 \sigma = \frac{\sigma^4}{(2\sigma^2 - 1)^{\frac{3}{2}}}\end{aligned}$$

where  $\tau^2 = \frac{\sigma^2}{2\sigma^2 - 1}$ .

### Example 3.5

We have  $I(\psi_{\text{CE}})^{-1} = \frac{(2\mu^2 + 2\tau^2)^2}{2}$  by standard properties of the single parameter Gaussian exponential family, so  $B_{\text{CE}} = \frac{2}{(2\mu^2 + 2\tau^2)^2}$ .

Additionally,

$$\begin{aligned}M_{\text{CE}} &= \text{Cov}_{\mathbf{P}}(T) = \mathbf{P} \left( ((\text{id} - \mu)^2 - \tau^2 - \mu^2)^2 \right) \\ &= \mathbf{P} \left( (\text{id}^2 - \tau^2 + 2\mu \text{id})^2 \right) \\ &= \mathbf{P} \left( \text{id}^4 - 2\text{id}^2 \tau^2 + 4\mu \text{id}^3 + \tau^4 + 4\mu^2 \text{id}^2 \right) \\ &= \nu + 4\mu^2 \tau^2 - \tau^4,\end{aligned}$$

as  $\mathbf{P} \text{id}^3 = 0$  and  $\mathbf{P} \text{id}^2 = \tau^2$ . In total

$$V_{\text{CE}} = B_{\text{CE}} M_{\text{CE}} B_{\text{CE}} = \frac{4(\nu + 4\mu^2 \tau^2 - \tau^4)}{(2\mu^2 + 2\tau^2)^4} = \frac{\nu + 4\mu^2 \tau^2 - \tau^4}{4(\mu^2 + \tau^2)^4}.$$

For EIS, we have  $\log p(x) = -\frac{1}{2} \frac{x^2}{\tau^2} - \frac{1}{2\sqrt{\pi\tau^2}}$ , so the log-weights are, up to an additive constant, given by

$$-\frac{1}{2} \frac{\text{id}^2}{\tau^2} - T\psi_{\text{EIS}}.$$

$M_{\text{EIS}}$  is then given by

$$\mathbf{P} \left( (\text{id} - \mu)^4 (\log w - \mathbf{P} \log w)^2 \right)$$

which is the expectation of a sixth order polynomial with respect to the standard normal distribution  $\mathbf{P}$ , so its value is analytically tractable and turns out to be<sup>1</sup>

$$M_{\text{EIS}} = \frac{\mu^2 (2\mu^6 + 45\mu^4\tau^2 + 15\tau^6)}{(2\mu^2 + \tau^2)^2}.$$

---

<sup>1</sup>See `03_08_comparison.ipynb` for calculations.



# Bibliography

- Agapiou, S. et al. (Jan. 14, 2017). *Importance Sampling: Intrinsic Dimension and Computational Cost*. DOI: [10.48550/arXiv.1511.06196](https://doi.org/10.48550/arXiv.1511.06196). arXiv: [1511.06196](https://arxiv.org/abs/1511.06196) [stat]. Pre-published.
- Akhmetzhanov, A. R. (2021). “Estimation of Delay-Adjusted All-Cause Excess Mortality in the USA: March-December 2020.” In: *Epidemiology and Infection*. DOI: [10.1017/s0950268821001527](https://doi.org/10.1017/s0950268821001527). pmid: [34210370](https://pubmed.ncbi.nlm.nih.gov/34210370/).
- An Der Heiden, M. et al. (Apr. 22, 2020). “Schätzung Der Aktuellen Entwicklung Der SARS-CoV-2-Epidemie in Deutschland – Nowcasting.” In: *Epidemiologisches Bulletin*. DOI: [10.25646/6692.4](https://doi.org/10.25646/6692.4).
- Ansley, C. F. et al. (1985). “Estimation, Filtering, and Smoothing in State Space Models with Incompletely Specified Initial Conditions.” In: *The Annals of Statistics*, pp. 1286–1316. JSTOR: [2241356](https://www.jstor.org/stable/2241356). URL: [https://www.jstor.org/stable/2241356?casa\\_token=gV7gKqDaiGEAAAAA:4tRfNK4F6dxh-q6oMa9nZ3i72JTLy-urE2vOVFKSyONYmltb4kbUT-hE7LfLFcTLaSCK20Z7v4jLK1HEaHzzvYDuvhFc2ZdQ7](https://www.jstor.org/stable/2241356?casa_token=gV7gKqDaiGEAAAAA:4tRfNK4F6dxh-q6oMa9nZ3i72JTLy-urE2vOVFKSyONYmltb4kbUT-hE7LfLFcTLaSCK20Z7v4jLK1HEaHzzvYDuvhFc2ZdQ7) (visited on 10/15/2024).
- Arroyo-Marioli, F. et al. (Jan. 13, 2021). “Tracking R of COVID-19: A New Real-Time Estimation Using the Kalman Filter.” In: *PLOS ONE* 16.1. DOI: [10.1371/journal.pone.0244474](https://doi.org/10.1371/journal.pone.0244474). pmid: [33439880](https://pubmed.ncbi.nlm.nih.gov/33439880/).
- Assimakis, N. et al. (2012). “Information Filter and Kalman Filter Comparison: Selection of the Faster Filter.” In: *Information Engineering*. Vol. 2. 1, pp. 1–5. URL: [http://madam.users.uth.gr/papers/3%20IJIE\\_2012.pdf](http://madam.users.uth.gr/papers/3%20IJIE_2012.pdf) (visited on 06/24/2024).
- Bastos, L. S. et al. (2019). “A Modelling Approach for Correcting Reporting Delays in Disease Surveillance Data.” In: *Statistics in Medicine*. DOI: [10.1002/sim.8303](https://doi.org/10.1002/sim.8303).
- Bazaraa, M. S. et al. (2006). *Nonlinear Programming: Theory and Algorithms*. 3. ed. Hoboken, NJ: Wiley-Interscience. 853 pp. ISBN: 978-0-471-48600-8.
- Bengtsson, T. et al. (2008). “Curse-of-Dimensionality Revisited: Collapse of the Particle Filter in Very Large Scale Systems.” In: *Probability and statistics: Essays in honor of David A. Freedman* 2, pp. 316–334. DOI: [10.1214/193940307000000518](https://doi.org/10.1214/193940307000000518).
- Billingsley, P. (1995). *Probability and Measure*. 3rd ed. Wiley Series in Probability and Mathematical Statistics. New York: Wiley. 593 pp. ISBN: 978-0-471-00710-4.
- Biswas, M. et al. (Dec. 9, 2020). “Association of Sex, Age, and Comorbidities with Mortality in COVID-19 Patients: A Systematic Review and Meta-Analysis.” In: *Intervirology* 64.1, pp. 36–47. ISSN: 0300-5526. DOI: [10.1159/000512592](https://doi.org/10.1159/000512592).
- Bracher, J. et al. (Aug. 27, 2021). “A Pre-Registered Short-Term Forecasting Study of COVID-19 in Germany and Poland during the Second Wave.” In: *Nat Commun* 12.1 (1), p. 5173. ISSN: 2041-1723. DOI: [10.1038/s41467-021-25207-0](https://doi.org/10.1038/s41467-021-25207-0).
- Branch, M. A. et al. (Jan. 1999). “A Subspace, Interior, and Conjugate Gradient Method for Large-Scale Bound-Constrained Minimization Problems.” In: *SIAM J. Sci. Comput.* 21.1, pp. 1–23. ISSN: 1064-8275, 1095-7197. DOI: [10.1137/S1064827595289108](https://doi.org/10.1137/S1064827595289108).
- Britton, T. et al. (2019). *Stochastic Epidemic Models with Inference*. Ed. by T. Britton et al. Springer.
- Brooks, S. et al., eds. (2011). *Handbook for Markov Chain Monte Carlo*. Boca Raton: Taylor & Francis. 592 pp. ISBN: 978-1-4200-7941-8.
- Brown, L. D. (1986). *Fundamentals of Statistical Exponential Families: With Applications in Statistical Decision Theory*. Lecture Notes-Monograph Series v. 9. Hayward, Calif: Institute of Mathematical Statistics. 283 pp. ISBN: 978-0-940600-10-2.

- Byambasuren, O. et al. (Dec. 31, 2020). “Estimating the Extent of Asymptomatic COVID-19 and Its Potential for Community Transmission: Systematic Review and Meta-Analysis.” In: *Journal of the Association of Medical Microbiology and Infectious Disease Canada* 5.4, pp. 223–234. DOI: [10.3138/jammi-2020-0030](https://doi.org/10.3138/jammi-2020-0030).
- Chan, J. C. C. et al. (Sept. 1, 2012). “Improved Cross-Entropy Method for Estimation.” In: *Stat Comput* 22.5, pp. 1031–1040. ISSN: 1573-1375. DOI: [10.1007/s11222-011-9275-7](https://doi.org/10.1007/s11222-011-9275-7).
- Chan, J. C. C. et al. (May 1, 2012). *Marginal Likelihood Estimation with the Cross-Entropy Method*. DOI: [10.2139/ssrn.2055042](https://doi.org/10.2139/ssrn.2055042). Pre-published.
- Chatterjee, S. et al. (Apr. 1, 2018). “The Sample Size Required in Importance Sampling.” In: *Ann. Appl. Probab.* 28.2. ISSN: 1050-5164. DOI: [10.1214/17-AAP1326](https://doi.org/10.1214/17-AAP1326).
- Chopin, N. et al. (Feb. 1, 2017). “Leave Pima Indians Alone: Binary Regression as a Benchmark for Bayesian Computation.” In: *Statist. Sci.* 32.1. ISSN: 0883-4237. DOI: [10.1214/16-ST581](https://doi.org/10.1214/16-ST581).
- Chopin, N. et al. (2020). *An Introduction to Sequential Monte Carlo*. Springer Series in Statistics. Cham, Switzerland: Springer. 378 pp. ISBN: 978-3-030-47844-5.
- Cover, T. M. et al. (2006). *Elements of Information Theory*. 2nd ed. Hoboken, N.J: Wiley-Interscience. 748 pp. ISBN: 978-0-471-24195-9.
- Danielsson, J. et al. (1993). “Accelerated Gaussian Importance Sampler with Application to Dynamic Latent Variable Models.” In: *Journal of Applied Econometrics* 8.S1, S153–S173. ISSN: 1099-1255. DOI: [10.1002/jae.3950080510](https://doi.org/10.1002/jae.3950080510).
- Desai, A. N. et al. (Aug. 2019). “Real-Time Epidemic Forecasting: Challenges and Opportunities.” In: *Health Security* 17.4, pp. 268–275. ISSN: 2326-5094. DOI: [10.1089/hs.2019.0022](https://doi.org/10.1089/hs.2019.0022).
- Diekmann, O. et al. (2013). *Mathematical Tools for Understanding Infectious Diseases Dynamics*. Princeton Series in Theoretical and Computational Biology. Princeton: Princeton University Press. 502 pp. ISBN: 978-0-691-15539-5.
- Du, Z. et al. (July 1, 2022). “Reproduction Numbers of Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2) Variants: A Systematic Review and Meta-analysis.” In: *Clinical Infectious Diseases* 75.1, e293–e295. ISSN: 1058-4838. DOI: [10.1093/cid/ciac137](https://doi.org/10.1093/cid/ciac137).
- Durbin, J. et al. (2012). *Time Series Analysis by State Space Methods*. 2nd ed. Oxford Statistical Science Series 38. Oxford: Oxford University Press. 346 pp. ISBN: 978-0-19-964117-8.
- Durbin, J. et al. (Sept. 1, 1997). “Monte Carlo Maximum Likelihood Estimation for Non-Gaussian State Space Models.” In: *Biometrika* 84.3, pp. 669–684. ISSN: 0006-3444. DOI: [10.1093/biomet/84.3.669](https://doi.org/10.1093/biomet/84.3.669).
- (2002). “A Simple and Efficient Simulation Smoother for State Space Time Series Analysis.” In: *Biometrika* 89.3, pp. 603–616. DOI: [10.1093/biomet/89.3.603](https://doi.org/10.1093/biomet/89.3.603).
- Ehre, M. et al. (Mar. 31, 2023). “Certified Dimension Reduction for Bayesian Updating with the Cross-Entropy Method.” In: *SIAM/ASA J. Uncertainty Quantification* 11.1, pp. 358–388. DOI: [10.1137/22M1484031](https://doi.org/10.1137/22M1484031).
- Engbert, R. et al. (Dec. 8, 2020). “Sequential Data Assimilation of the Stochastic SEIR Epidemic Model for Regional COVID-19 Dynamics.” In: *Bull Math Biol* 83.1, p. 1. ISSN: 1522-9602. DOI: [10.1007/s11538-020-00834-8](https://doi.org/10.1007/s11538-020-00834-8).
- Engel, M. et al. (Jan. 15, 2023). “Bayesian Updating and Marginal Likelihood Estimation by Cross Entropy Based Importance Sampling.” In: *Journal of Computational Physics* 473, p. 111746. ISSN: 0021-9991. DOI: [10.1016/j.jcp.2022.111746](https://doi.org/10.1016/j.jcp.2022.111746).
- Evensen, G. (1994). “Sequential Data Assimilation with a Nonlinear Quasi-Geostrophic Model Using Monte Carlo Methods to Forecast Error Statistics.” In: *Journal of Geophysical Research: Oceans* 99.C5, pp. 10143–10162. ISSN: 2156-2202. DOI: [10.1029/94JC00572](https://doi.org/10.1029/94JC00572).
- Farrington, C. P. et al. (1996). “A Statistical Algorithm for the Early Detection of Outbreaks of Infectious Disease.” In: *Journal of The Royal Statistical Society Series A-statistics in Society*. DOI: [10.2307/2983331](https://doi.org/10.2307/2983331).
- Fraser, D. et al. (Aug. 1969). “The Optimum Linear Smoother as a Combination of Two Optimum Linear Filters.” In: *IEEE Trans. Automat. Contr.* 14.4, pp. 387–390. ISSN: 0018-9286. DOI: [10.1109/TAC.1969.1099196](https://doi.org/10.1109/TAC.1969.1099196).
- Frühwirth-Schnatter, S. (1994). “Data Augmentation and Dynamic Linear Models.” In: *Journal of Time Series Analysis* 15.2, pp. 183–202. ISSN: 1467-9892. DOI: [10.1111/j.1467-9892.1994.tb00184.x](https://doi.org/10.1111/j.1467-9892.1994.tb00184.x).

- “Discrete Spatial Variation” (2010). In: *Handbook of Spatial Statistics*. Ed. by A. E. Gelfand et al. CRC Press. ISBN: 978-0-429-13650-4.
- Günther, F. et al. (2021). “Nowcasting the COVID-19 Pandemic in Bavaria.” In: *Biometrical Journal* 63.3, pp. 490–502. ISSN: 1521-4036. DOI: [10.1002/bimj.202000112](https://doi.org/10.1002/bimj.202000112).
- Haberman, S. J. (1989). “Concavity and Estimation.” In: *The Annals of Statistics* 17.4, pp. 1631–1661. ISSN: 0090-5364. DOI: [10.1214/aos/1176347385](https://doi.org/10.1214/aos/1176347385). JSTOR: [2241655](https://www.jstor.org/stable/2241655).
- Havers, F. P. et al. (Sept. 26, 2024). “COVID-19–Associated Hospitalizations and Maternal Vaccination Among Infants Aged <6 Months — COVID-NET, 12 States, October 2022–April 2024.” In: *MMWR Morb. Mortal. Wkly. Rep.* 73.38, pp. 830–836. ISSN: 0149-2195, 1545-861X. DOI: [10.15585/mmwr.mm7338a1](https://doi.org/10.15585/mmwr.mm7338a1).
- Heyder, S. et al. (Oct. 4, 2023). “Measures of COVID-19 Spread.” In: *Covid-19 pandisziplinär und international: Gesundheitswissenschaftliche, gesellschaftspolitische und philosophische Hintergründe*. Ed. by A. Kraemer et al. Medizin, Kultur, Gesellschaft. Wiesbaden: Springer Fachmedien, pp. 51–66. ISBN: 978-3-658-40525-0. DOI: [10.1007/978-3-658-40525-0\\_3](https://doi.org/10.1007/978-3-658-40525-0_3).
- Höhle, M. et al. (2014). “Bayesian Nowcasting during the STEC O104:H4 Outbreak in Germany, 2011.” In: *Biometrics* 70.4, pp. 993–1002. ISSN: 1541-0420. DOI: [10.1111/biom.12194](https://doi.org/10.1111/biom.12194).
- Homem-de-Mello, T. (July 20, 2007). “A Study on the Cross-Entropy Method for Rare-Event Probability Estimation.” In: *INFORMS Journal on Computing*. DOI: [10.1287/ijoc.1060.0176](https://doi.org/10.1287/ijoc.1060.0176).
- Hospitalization Nowcast Hub* (Oct. 31, 2022). KITmetricslab. URL: <https://github.com/KITmetricslab/hospitalization-nowcast-hub> (visited on 11/09/2022).
- Hughes, T. D. et al. (Mar. 18, 2023). “The Effect of SARS-CoV-2 Variant on Respiratory Features and Mortality.” In: *Sci Rep* 13.1, p. 4503. ISSN: 2045-2322. DOI: [10.1038/s41598-023-31761-y](https://doi.org/10.1038/s41598-023-31761-y).
- Jazwinski, A. H. (1970). *Stochastic Processes and Filtering Theory*. New York: Academic Press. 376 pp.
- Johnson, N. L. et al. (1994). *Continuous Univariate Distributions*. 2nd ed. Wiley Series in Probability and Mathematical Statistics. New York: Wiley. 2 pp. ISBN: 978-0-471-58495-7 978-0-471-58494-0.
- Julier, S. J. et al. (July 28, 1997). “New Extension of the Kalman Filter to Nonlinear Systems.” In: *Signal Processing, Sensor Fusion, and Target Recognition VI*. Vol. 3068. International Society for Optics and Photonics, pp. 182–194. DOI: [10.1117/12.280797](https://doi.org/10.1117/12.280797).
- Jungbacker, B. et al. (Dec. 1, 2007). “Monte Carlo Estimation for Nonlinear Non-Gaussian State Space Models.” In: *Biometrika* 94.4, pp. 827–839. ISSN: 0006-3444. DOI: [10.1093/biomet/asm074](https://doi.org/10.1093/biomet/asm074).
- Kaminsky, K. S. (Apr. 1, 1987). “Prediction of IBNR Claim Counts by Modelling the Distribution of Report Lags.” In: *Insurance Mathematics & Economics* 6.2, pp. 151–159. DOI: [10.1016/0167-6687\(87\)90024-2](https://doi.org/10.1016/0167-6687(87)90024-2).
- Kappen, H. J. et al. (Mar. 1, 2016). “Adaptive Importance Sampling for Control and Inference.” In: *J Stat Phys* 162.5, pp. 1244–1266. ISSN: 1572-9613. DOI: [10.1007/s10955-016-1446-7](https://doi.org/10.1007/s10955-016-1446-7).
- Katzfuss, M. et al. (Oct. 1, 2016). “Understanding the Ensemble Kalman Filter.” In: *The American Statistician* 70.4, pp. 350–357. ISSN: 0003-1305, 1537-2731. DOI: [10.1080/00031305.2016.1141709](https://doi.org/10.1080/00031305.2016.1141709).
- Kermack, W. O. et al. (Aug. 1927). “A Contribution to the Mathematical Theory of Epidemics.” In: *Proc. R. Soc. Lond. A* 115.772, pp. 700–721. ISSN: 0950-1207, 2053-9150. DOI: [10.1098/rspa.1927.0118](https://doi.org/10.1098/rspa.1927.0118).
- Kong, A. (1992). “A Note on Importance Sampling Using Standardized Weights.” In: *University of Chicago, Dept. of Statistics, Tech. Rep* 348, p. 14.
- Kong, A. et al. (Mar. 1994). “Sequential Imputations and Bayesian Missing Data Problems.” In: *Journal of the American Statistical Association* 89.425, pp. 278–288. ISSN: 0162-1459. DOI: [10.1080/01621459.1994.10476469](https://doi.org/10.1080/01621459.1994.10476469).
- Koopman, S. J. (Dec. 1997). “Exact Initial Kalman Filtering and Smoothing for Nonstationary Time Series Models.” In: *Journal of the American Statistical Association* 92.440, pp. 1630–1638. ISSN: 0162-1459, 1537-274X. DOI: [10.1080/01621459.1997.10473685](https://doi.org/10.1080/01621459.1997.10473685).
- Koopman, S. J. et al. (1992). “Exact Score for Time Series Models in State Space Form.” In: *Biometrika* 79.4, pp. 823–826. ISSN: 0006-3444. DOI: [10.2307/2337237](https://doi.org/10.2307/2337237). JSTOR: [2337237](https://www.jstor.org/stable/2337237).
- Koopman, S. J. et al. (Jan. 2, 2015). “Numerically Accelerated Importance Sampling for Nonlinear Non-Gaussian State-Space Models.” In: *Journal of Business & Economic Statistics* 33.1, pp. 114–127. ISSN: 0735-0015, 1537-2707. DOI: [10.1080/07350015.2014.925807](https://doi.org/10.1080/07350015.2014.925807).

- Koopman, S. J. et al. (2019). "Modified Efficient Importance Sampling for Partially Non-Gaussian State Space Models." In: *Statistica Neerlandica* 73.1, pp. 44–62. ISSN: 1467-9574. DOI: [10.1111/stan.12128](https://doi.org/10.1111/stan.12128).
- Laplace, P. S. (Aug. 1986). "Memoir on the Probability of the Causes of Events." In: *Statistical Science* 1.3, pp. 364–378. ISSN: 0883-4237, 2168-8745. DOI: [10.1214/ss/1177013621](https://doi.org/10.1214/ss/1177013621).
- Lauritzen, S. L. (1996). *Graphical Models*. Oxford Statistical Science Series 17. Oxford : New York: Clarendon Press ; Oxford University Press. 298 pp. ISBN: 978-0-19-852219-5.
- Lawless, J. F. (1994). "Adjustments for Reporting Delays and the Prediction of Occurred but Not Reported Events." In: *Canadian Journal of Statistics* 22.1, pp. 15–31. ISSN: 1708-945X. DOI: [10.2307/3315826.n1](https://doi.org/10.2307/3315826.n1).
- Liang, K.-Y. et al. (May 1995). "Inference Based on Estimating Functions in the Presence of Nuisance Parameters." In: *Statistical Science* 10.2, pp. 158–173. ISSN: 0883-4237, 2168-8745. DOI: [10.1214/ss/1177010028](https://doi.org/10.1214/ss/1177010028).
- Liesenfeld, R. et al. (Sept. 1, 2003). "Univariate and Multivariate Stochastic Volatility Models: Estimation and Diagnostics." In: *Journal of Empirical Finance* 10.4, pp. 505–531. ISSN: 0927-5398. DOI: [10.1016/S0927-5398\(02\)00072-5](https://doi.org/10.1016/S0927-5398(02)00072-5).
- Lloyd-Smith, J. O. et al. (Nov. 2005). "Superspreading and the Effect of Individual Variation on Disease Emergence." In: *Nature* 438.7066, pp. 355–359. ISSN: 1476-4687. DOI: [10.1038/nature04153](https://doi.org/10.1038/nature04153).
- McGough, S. F. et al. (2020). "Nowcasting by Bayesian Smoothing: A Flexible, Generalizable Model for Real-Time Epidemic Tracking." In: *PLOS Computational Biology*. DOI: [10.1371/journal.pcbi.1007735](https://doi.org/10.1371/journal.pcbi.1007735). pmid: [32251464](https://pubmed.ncbi.nlm.nih.gov/32251464/).
- Midthune, D. N. et al. (Mar. 1, 2005). "Modeling Reporting Delays and Reporting Corrections in Cancer Registry Data." In: *Journal of the American Statistical Association* 100.469, pp. 61–70. ISSN: 0162-1459. DOI: [10.1198/016214504000001899](https://doi.org/10.1198/016214504000001899).
- Morf, M. et al. (Aug. 1975). "Square-Root Algorithms for Least-Squares Estimation." In: *IEEE Transactions on Automatic Control* 20.4, pp. 487–497. ISSN: 1558-2523. DOI: [10.1109/TAC.1975.1100994](https://doi.org/10.1109/TAC.1975.1100994).
- Mossong, J. et al. (Mar. 25, 2008). "Social Contacts and Mixing Patterns Relevant to the Spread of Infectious Diseases." In: *PLOS Medicine* 5.3, e74. ISSN: 1549-1676. DOI: [10.1371/journal.pmed.0050074](https://doi.org/10.1371/journal.pmed.0050074).
- Neiswanger, W. et al. (July 23, 2014). "Asymptotically Exact, Embarrassingly Parallel MCMC." In: *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence*. UAI'14. Arlington, Virginia, USA: AUAI Press, pp. 623–632. ISBN: 978-0-9749039-1-0.
- Nocedal, J. et al. (2006). *Numerical Optimization*. Second edition. Springer Series in Operation Research and Financial Engineering. New York, NY: Springer. 664 pp. ISBN: 978-0-387-30303-1 978-1-4939-3711-0.
- Noufaily, A. et al. (2015). "Modelling Reporting Delays for Outbreak Detection in Infectious Disease Data." In: *Journal of The Royal Statistical Society Series A-statistics in Society*. DOI: [10.1111/rssa.12055](https://doi.org/10.1111/rssa.12055).
- Nowcasts Der COVID-19 Hospitalisierungsinzidenz* (2022). URL: <https://covid19nowcasthub.de/> (visited on 11/09/2022).
- Rao, C. R. (2002). *Linear Statistical Inference and Its Applications*. 2. ed., Paperback ed. Wiley Series in Probability and Statistics. New York: Wiley. 625 pp. ISBN: 978-0-471-21875-3.
- Ray, E. L. et al. (Aug. 22, 2020). "Ensemble Forecasts of Coronavirus Disease 2019 (COVID-19) in the U.S." In: *medRxiv*, p. 2020.08.19.20177493. DOI: [10.1101/2020.08.19.20177493](https://doi.org/10.1101/2020.08.19.20177493).
- Renshaw, A. E. et al. (1998). "A Stochastic Model Underlying the Chain-Ladder Technique." In: *British Actuarial Journal* 4.4, pp. 903–923. DOI: [10.1017/S1357321700000222](https://doi.org/10.1017/S1357321700000222).
- Richard, J.-F. et al. (Dec. 1, 2007). "Efficient High-Dimensional Importance Sampling." In: *Journal of Econometrics* 141.2, pp. 1385–1411. ISSN: 0304-4076. DOI: [10.1016/j.jeconom.2007.02.007](https://doi.org/10.1016/j.jeconom.2007.02.007).
- Ripley, B. D. (2009). *Stochastic Simulation*. Vol. 316. John Wiley & Sons.
- Robert Koch-Institut (Oct. 1, 2021). *COVID-19-Hospitalisierungen in Deutschland*. Version 2021-10-01. Zenodo. DOI: [10.5281/ZENODO.5519056](https://doi.org/10.5281/ZENODO.5519056).
- (Feb. 7, 2022). *SARS-CoV-2 Infektionen in Deutschland*. Version 2022-02-07. Zenodo. DOI: [10.5281/ZENODO.4681153](https://doi.org/10.5281/ZENODO.4681153).



- (Aug. 19, 2024a). *COVID-19-Hospitalisierungen in Deutschland*. Version 2024-08-19. Zenodo. DOI: [10.5281/ZENODO.5519056](https://doi.org/10.5281/ZENODO.5519056).
- (Aug. 18, 2024b). *SARS-CoV-2 Infektionen in Deutschland*. Version 2024-08-18. Zenodo. DOI: [10.5281/ZENODO.4681153](https://doi.org/10.5281/ZENODO.4681153).
- (Oct. 23, 2024c). *SARS-CoV-2 Sequenzdaten aus Deutschland*. Version 2024-10-22. Zenodo. DOI: [10.5281/ZENODO.5139363](https://doi.org/10.5281/ZENODO.5139363).
- Rubinstein, R. Y. (Sept. 1, 1999). “The Cross-Entropy Method for Combinatorial and Continuous Optimization.” In: *Methodology and Computing in Applied Probability* 1.2, pp. 127–190. ISSN: 1573-7713. DOI: [10.1023/A:1010091220143](https://doi.org/10.1023/A:1010091220143).
- Rubinstein, R. Y. et al. (2004). *The Cross-Entropy Method: A Unified Approach to Combinatorial Optimization, Monte-Carlo Simulation and Machine Learning*. New York, NY: Springer New York. ISBN: 978-1-4757-4321-0.
- Rubinstein, R. Y. et al. (Nov. 6, 2009). “How to Deal with the Curse of Dimensionality of Likelihood Ratios in Monte Carlo Simulation.” In: *Stochastic Models* 25.4, pp. 547–568. ISSN: 1532-6349. DOI: [10.1080/15326340903291248](https://doi.org/10.1080/15326340903291248).
- Salmon, M. et al. (2015). “Bayesian Outbreak Detection in the Presence of Reporting Delays.” In: *Biometrical Journal*. DOI: [10.1002/bimj.201400159](https://doi.org/10.1002/bimj.201400159). pmid: [26250543](https://pubmed.ncbi.nlm.nih.gov/26250543/).
- Salzberger, B. et al. (Apr. 1, 2021). “Epidemiology of SARS-CoV-2.” In: *Infection* 49.2, pp. 233–239. ISSN: 1439-0973. DOI: [10.1007/s15010-020-01531-3](https://doi.org/10.1007/s15010-020-01531-3).
- Schäfer, F. et al. (Jan. 2021). “Sparse Cholesky Factorization by Kullback–Leibler Minimization.” In: *SIAM J. Sci. Comput.* 43.3, A2019–A2046. ISSN: 1064-8275. DOI: [10.1137/20M1336254](https://doi.org/10.1137/20M1336254).
- Schneble, M. et al. (Mar. 2021). “Nowcasting Fatal COVID-19 Infections on a Regional Level in Germany.” In: *Biometrical Journal* 63.3, pp. 471–489. ISSN: 0323-3847, 1521-4036. DOI: [10.1002/bimj.202000143](https://doi.org/10.1002/bimj.202000143).
- Schneider, W. (1986). *Der Kalmanfilter Als Instrument Zur Diagnose Und Schätzung Variabler Parameter in Ökonometrischen Modellen*. Arbeiten Zur Angewandten Statistik Bd. 27. Heidelberg: Physica-Verlag. 490 pp. ISBN: 978-3-7908-0359-4.
- Shephard, N. (1994). “Partial Non-Gaussian State Space.” In: *Biometrika* 81.1, pp. 115–131. DOI: [10.1093/biomet/81.1.115](https://doi.org/10.1093/biomet/81.1.115).
- Shephard, N. et al. (Sept. 1, 1997). “Likelihood Analysis of Non-Gaussian Measurement Time Series.” In: *Biometrika* 84.3, pp. 653–667. ISSN: 0006-3444. DOI: [10.1093/biomet/84.3.653](https://doi.org/10.1093/biomet/84.3.653).
- Song, J. et al. (May 1, 2021). “Maximum Likelihood-Based Extended Kalman Filter for COVID-19 Prediction.” In: *Chaos, Solitons & Fractals* 146, p. 110922. ISSN: 0960-0779. DOI: [10.1016/j.chaos.2021.110922](https://doi.org/10.1016/j.chaos.2021.110922).
- Stark, P. B. et al. (1995). “Bounded-Variable Least-Squares: An Algorithm and Applications.” In: *Computational Statistics* 10, pp. 129–129. URL: <https://digitalassets.lib.berkeley.edu/sdtr/ucb/text/394.pdf> (visited on 06/05/2024).
- Tierney, L. et al. (Mar. 1986). “Accurate Approximations for Posterior Moments and Marginal Densities.” In: *Journal of the American Statistical Association* 81.393, pp. 82–86. ISSN: 0162-1459, 1537-274X. DOI: [10.1080/01621459.1986.10478240](https://doi.org/10.1080/01621459.1986.10478240).
- Tierney, L. et al. (Sept. 1989). “Fully Exponential Laplace Approximations to Expectations and Variances of Nonpositive Functions.” In: *Journal of the American Statistical Association* 84.407, pp. 710–716. ISSN: 0162-1459, 1537-274X. DOI: [10.1080/01621459.1989.10478824](https://doi.org/10.1080/01621459.1989.10478824).
- Tomori, D. V. et al. (Mar. 26, 2021). “Individual Social Contact Data Reflected SARS-CoV-2 Transmission Dynamics during the First Wave in Germany Better than Population Mobility Data – an Analysis Based on the COVIMOD Study.” In: p. 2021.03.24.21254194. DOI: [10.1101/2021.03.24.21254194](https://doi.org/10.1101/2021.03.24.21254194).
- Van de Kasstele, J. et al. (2019). “Nowcasting the Number of New Symptomatic Cases during Infectious Disease Outbreaks Using Constrained P-Spline Smoothing.” In: *Epidemiology*. DOI: [10.1097/ede.0000000000001050](https://doi.org/10.1097/ede.0000000000001050). pmid: [31205290](https://pubmed.ncbi.nlm.nih.gov/31205290/).
- Van der Vaart, A. W. (2000). *Asymptotic Statistics*. Cambridge: Cambridge University Press.
- White, H. (1982). “Maximum Likelihood Estimation of Misspecified Models.” In: *Econometrica* 50.1, pp. 1–25. ISSN: 0012-9682. DOI: [10.2307/1912526](https://doi.org/10.2307/1912526). JSTOR: [1912526](https://www.jstor.org/stable/1912526).
- Whitt, W. (Nov. 1976). “Bivariate Distributions with Given Marginals.” In: *The Annals of Statistics* 4.6, pp. 1280–1289. ISSN: 0090-5364, 2168-8966. DOI: [10.1214/aos/1176343660](https://doi.org/10.1214/aos/1176343660).

- Wolffram, D. et al. (Aug. 11, 2023). “Collaborative Nowcasting of COVID-19 Hospitalization Incidences in Germany.” In: *PLOS Computational Biology* 19.8, e1011394. ISSN: 1553-7358. DOI: [10.1371/journal.pcbi.1011394](https://doi.org/10.1371/journal.pcbi.1011394).
- Wu, N. et al. (May 1, 2023). “Long-Term Effectiveness of COVID-19 Vaccines against Infections, Hospitalisations, and Mortality in Adults: Findings from a Rapid Living Systematic Evidence Synthesis and Meta-Analysis up to December, 2022.” In: *The Lancet Respiratory Medicine* 11.5, pp. 439–452. ISSN: 2213-2600, 2213-2619. DOI: [10.1016/S2213-2600\(23\)00015-2](https://doi.org/10.1016/S2213-2600(23)00015-2). pmid: [36780914](https://pubmed.ncbi.nlm.nih.gov/36780914/).
- Zeger, S. L. et al. (1989). “Statistical Methods for Monitoring the AIDS Epidemic.” In: *Statistics in Medicine* 8.1, pp. 3–21. DOI: [10.1002/sim.4780080104](https://doi.org/10.1002/sim.4780080104).
- Zhang, W. et al. (Jan. 2014). “Applications of the Cross-Entropy Method to Importance Sampling and Optimal Control of Diffusions.” In: *SIAM J. Sci. Comput.* 36.6, A2654–A2672. ISSN: 1064-8275. DOI: [10.1137/14096493X](https://doi.org/10.1137/14096493X).
- Zhu, X. et al. (Oct. 1, 2021). “Extended Kalman Filter Based on Stochastic Epidemiological Model for COVID-19 Modelling.” In: *Computers in Biology and Medicine* 137, p. 104810. ISSN: 0010-4825. DOI: [10.1016/j.compbiomed.2021.104810](https://doi.org/10.1016/j.compbiomed.2021.104810).







# Symbols

$R_c$	case reproduction number
$R_t$	time-varying reproduction number
$X_{:t}$	$(X_0, \dots, X_t)$ for $t \in \mathbf{N}_0$
$X_{s:t}$	$(X_s, \dots, X_t)$ for $s, t \in \mathbf{N}_0$ , $s < t$
$\mathbf{R}_{>0}$	positive real numbers
$\Theta$	set of all parameters
Pois	Poisson distribution
$\rho$	growth factor
$\theta$	a parameter
$d$	doubling time
$p$	generic density
$r$	exponential growth rate
$test$	doubling
$w$	generation time distribution
$\mathbf{N}_0$	Natural numbers, including 0
$\mathbf{N}$	Natural numbers, excluding 0
$\mathbf{R}^{p \times m}$	matrices with $p$ rows and $m$ columns of real entries
$\mathbb{E}$	expected value



# Abbreviations

**BLUP** best linear unbiased predictor. 37

**CDF** cumulative distribution function. 63

**CE-method** Cross-Entropy method. ix, 2, 5, 8, 23, 28–34, 36, 37, 39, 41–44, 46–49, 51–54, 63–73, 75, 97, 99

**CLT** central limit theorem. 34, 35, 39, 42

**COD** curse of dimensionality. 31

**COVID-19** Coronavirus disease 2019. 1, 2, 4, 14, 77, 93, 94, 97

**CRN** common random number. 31, 53, 54, 57, 61, 62

**EF** efficiency factor. 24–26, 94

**EGSSM** Exponential Family Partially Gaussian state space model. 17, 18, 44–46, 53, 55, 64, 65

**EIS** Efficient Importance Sampling. ix, 2, 5, 8, 15, 23, 28, 36–39, 41–46, 48, 53, 58, 61–73, 75, 80, 82, 88, 93, 97, 99, 102

**EKF** Extended Kalman filter. 12, 14

**EnKF** Ensemble Kalman filter. 12, 14

**ESS** effective sample size. 23–25, 28, 43, 66, 93, 94

**FFBS** Forwards Filter, Backwards Sampling. 13, 14, 59, 62

**GLSSM** Gaussian linear state space model. ix, 5, 8, 9, 11–15, 43–46, 48, 53, 54, 56–58, 60, 62, 72, 80, 88, 93, 96

**IS** Importance Sampling. 19

**KL-divergence** Kullback Leibler divergence. 2, 8, 21–23, 28, 29, 36, 51

**LA** Laplace approximation. ix, 18, 27, 28, 43–46, 53, 58, 60, 61, 63–66, 71, 72, 75, 80, 82, 93

**MC-integration** Monte-Carlo integration. 18, 45

**MCMC** Markov chain Monte Carlo. 8, 18, 19, 31, 43, 57

**MLE** maximum likelihood estimator. 8, 12, 57, 62, 80, 81, 85, 88, 93, 96

**MSE** mean-squared error. 21, 24

**NPI** non-pharmaceutical intervention. 78, 86

**PGSSM** Partially Gaussian state space model. 8, 15, 17, 19, 20, 45, 46, 56, 58, 61, 62, 80, 92, 97, 99

**PSD** positive semi-definite. 12

**RKI** Robert Koch-Institut. 1, 78, 89, 90

**SARS-CoV-2** Severe acute respiratory syndrome coronavirus 2. 4, 94

**SMC** sequential Monte Carlo. 8, 18, 19, 57

**SPD** symmetric positive definite. 14, 24, 65, 71

**SSM** state space model. 2, 6–9, 14, 15, 18, 27, 28, 41, 43, 46–48, 56, 57, 71, 78, 86, 87, 90, 97, 99

**UKF** unscented Kalman filter. 12

**VM-method** Variance-Minimization method. 42, 43

# Declaration

Put your declaration here.

*Ilmenau, October 2023*