# A Stochastic Variational Inference Approach for Semiparametric Distributional Regression

## Final Kolloquium MSc Applied Statistic

Sebastian Lorek

First supervisor: Prof. Dr. Thomas Kneib
Second supervisor: Gianmarco Callegher

2023-11-28

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA
SEIT 1737

# Table of Contents I

# Introduction

# Statistical Infrence

Frequentist inference:

- ML estimation (Newton-Rhapson/Fisher-Scoring)
- Hyperparameter estimation
- REML for latent parameters (random-effects)
- SGD, Automatic differentiation in machine/deep learning

Bayesian inference:

- Conjugate models, full conditional conjugate models, non-conjugate models
- Gibbs sampling
- Rejection/Importance sampling
- MCMC (IWLS, HMC, NUTS)
- MAP estimation in combination with the Laplace approximation

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

# Bayesian Inference

- Focal point of interest is the posterior distribution
- General posterior specification for Bayesian regression

$$p(\boldsymbol{\theta}|\mathbf{y}, \mathcal{D}) = \frac{p(\mathbf{y}, \boldsymbol{\theta}| \mathcal{D})}{p(\mathbf{y}|\mathcal{D})} \tag{1}$$

$$= \frac{p(\mathbf{y}| \boldsymbol{\theta}, \mathcal{D})p(\boldsymbol{\theta})}{p(\mathbf{y}|\mathcal{D})} \tag{2}$$

$$= \frac{p(\mathbf{y}| \boldsymbol{\theta}, \mathcal{D})p(\boldsymbol{\theta})}{\int p(\mathbf{y}| \boldsymbol{\theta}, \mathcal{D})p(\boldsymbol{\theta}) \, d\boldsymbol{\theta}} \tag{3}$$

$$\propto p(\mathbf{y}| \boldsymbol{\theta}, \mathcal{D})p(\boldsymbol{\theta}) \tag{4}$$

- Primary challenge is the calculation of the normalizing constant/evidence
- We need approximate methods that bypass the direct calculation of the normalizing constant

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

- MCMC methods have been **so far** the work horse for Bayesian inference in classical statistics
- Set up Markov chain (adhere to detailed balance and ergodicity) and sample . . .
- Enjoyes nice properties
  - Assurance of convergence to the true posterior
- But unfortunately does not scale well for modern applications
  - High dimensional parameter spaces (tausands of parameters)
  - Large dataset
  - Many latent and hyper-parameters
- Trade some **accuracy** for **scalability**
- Make use of SGD and automatic differentiation which works well for inference in machine/deep learning (frequentist inference)

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

# Theory

# Variational Inference

- VI is a method from machine learning to **approximate** probability densities (Jordan et al. 1999)
- Approximate posterior with a variational distribution $q(\boldsymbol{\theta})$ from a predefined (parametric) variational family $\mathcal{Q}$

$$q(\boldsymbol{\theta}) \in \mathcal{Q}$$

- Use optimization (SGD) to find a member that is as closely as possible to the true posterior
- What means close in terms of distributions ?

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

# Kullback-Leibler divergence

- Divergence measure (Kullback and Leibler 1951) that quantifies the proximity between two probability distributions

$$D_{KL} = \int q(\boldsymbol{\theta}) \ln \left( \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{y}, \mathcal{D})} \right) d\boldsymbol{\theta}$$
$$= E_{q(\boldsymbol{\theta})} \left[ \ln \left( \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{y}, \mathcal{D})} \right) \right]$$

- In short $D_{KL}(q \| p)$
- It holds that $D_{KL}(q \| p) \geq 0$
- Has some nice properties but also drawbacks (not a distance/metric)

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA

# Optimization objective

- Use optimization to find a variational distribution that is as close as possible in terms of the divergence measure to the true posterior

$$
\begin{aligned}
q^*(\boldsymbol{\theta}) &= \underset{q(\boldsymbol{\theta}) \in \mathcal{Q}}{\arg\min} \, \mathsf{D}_{\mathsf{KL}}\left(q(\boldsymbol{\theta}) \| \, p(\boldsymbol{\theta}|\mathbf{y}, \mathcal{D})\right) \\
&= \underset{q(\boldsymbol{\theta}) \in \mathcal{Q}}{\arg\min} \int q(\boldsymbol{\theta}) \ln\left(\frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{y}, \mathcal{D})}\right) \, d\boldsymbol{\theta} \\
&= \underset{q(\boldsymbol{\theta}) \in \mathcal{Q}}{\arg\min} \, \mathsf{E}_{q(\boldsymbol{\theta})}\left[\ln\left(\frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{y}, \mathcal{D})}\right)\right].
\end{aligned}
$$

- Flexibility of $\mathcal{Q}$ significantly influences the optimization
  - Complex $\mathcal{Q} \rightarrow$ better approximation, but increased complexity during optimization
  - Simple $\mathcal{Q} \rightarrow$ worse approximation, but simpler optimization
- Objective offers theoretical insights but remains **infeasible** to compute, due to containing the posterior (evidence)

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

# Evidence Lower Bound

- Think about a way to introduce a mathematical **equivalent** objective that does not depend on the evidence
- Lets start with the log-evidence

$$\ln(p(\mathbf{y}|\mathcal{D})) = \int q(\theta|\phi) \ln(p(\mathbf{y}|\mathcal{D})) \, d\theta \tag{5}$$

$$= \int q(\theta|\phi) \ln\left( \frac{p(\mathbf{y}|\mathcal{D}) p(\theta|\mathbf{y}, \mathcal{D})}{p(\theta|\mathbf{y}, \mathcal{D})} \right) d\theta \tag{6}$$

$$= \int q(\theta|\phi) \ln\left( \frac{p(\mathbf{y}, \theta|\mathcal{D})}{p(\theta|\mathbf{y}, \mathcal{D})} \right) d\theta \tag{7}$$

$$= \int q(\theta|\phi) \ln\left( \frac{p(\mathbf{y}, \theta|\mathcal{D})}{q(\theta|\phi)} \frac{q(\theta|\phi)}{p(\theta|\mathbf{y}, \mathcal{D})} \right) d\theta \tag{8}$$

$$= \int q(\theta|\phi) \ln\left( \frac{p(\mathbf{y}, \theta|\mathcal{D})}{q(\theta|\phi)} \right) d\theta + \int q(\theta|\phi) \ln\left( \frac{q(\theta|\phi)}{p(\theta|\mathbf{y}, \mathcal{D})} \right) d\theta \tag{9}$$

$$= \mathsf{E}_{q(\theta|\phi)}\left[ \ln\left( \frac{p(\mathbf{y}, \theta|\mathbf{X})}{q(\theta|\phi)} \right) \right] + \mathsf{D}_{\mathsf{KL}}(q(\theta|\phi)||p(\theta|\mathbf{y}, \mathcal{D})) \tag{10}$$

- Parameterize $q$ with $\phi$

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN

# Optimization objective revisited

- Rewrite 10 as

$$\ln(p(\mathbf{y}|\mathbf{X})) - D_{KL}(q||p) = E_{q(\theta|\phi)}\left[\ln\left(\frac{p(\mathbf{y},\theta|\mathcal{D})}{q(\theta|\phi)}\right)\right] \tag{11}$$

- And take arg max w.r.t. $\phi$

$$\underset{\phi}{\arg\max}\ \ln(p(\mathbf{y}|\mathcal{D})) - D_{KL}(q||p) = \underset{\phi}{\arg\max}\ E_{q(\theta|\phi)}\left[\ln\left(\frac{p(\mathbf{y},\theta|\mathcal{D})}{q(\theta|\phi)}\right)\right]$$

$$\underset{\phi}{\arg\max}\ -D_{KL}(q||p) = \underset{\phi}{\arg\max}\ E_{q(\theta|\phi)}\left[\ln\left(\frac{p(\mathbf{y},\theta|\mathcal{D})}{q(\theta|\phi)}\right)\right]$$

$$\underset{\phi}{\arg\min}\ D_{KL}(q||p) = \underset{\phi}{\arg\max}\ E_{q(\theta|\phi)}\left[\ln\left(\frac{p(\mathbf{y},\theta|\mathcal{D})}{q(\theta|\phi)}\right)\right]$$

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

- New optimization objective

$$\hat{\phi} = \underset{\phi}{\arg\max}\, \mathsf{E}_{q(\boldsymbol{\theta}|\phi)} \left[ \ln \left( \frac{p(\mathbf{y}, \boldsymbol{\theta}|\mathcal{D})}{q(\boldsymbol{\theta}|\phi)} \right) \right]$$

$$= \underset{\phi}{\arg\max}\, \mathsf{ELBO}(\phi)$$

- Take a breath 😩

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

- What does the ELBO ?

$$\text{ELBO}(\phi) = E_{q(\theta|\phi)} \left[ \ln \left( \frac{p(\mathbf{y}, \theta | \mathcal{D})}{q(\theta|\phi)} \right) \right] \tag{12}$$

$$= E_{q(\theta|\phi)} \left[ \ln \left( \frac{p(\mathbf{y}|\mathcal{D}, \theta) p(\theta)}{q(\theta|\phi)} \right) \right] \tag{13}$$

$$= E_{q(\theta|\phi)} \left[ \ln(p(\mathbf{y}|\mathcal{D}, \theta)) \right] + E_{q(\theta|\phi)} \left[ \ln(p(\theta)) \right] - E_{q(\theta|\phi)} \left[ \ln(q(\theta|\phi)) \right] \tag{14}$$

$$= E_{q(\theta|\phi)} \left[ \ln(p(\mathbf{y}|\mathcal{D}, \theta)) \right] - D_{\text{KL}} \left( q(\theta|\phi) || p(\theta) \right) \tag{15}$$

- Actually something similar to MAP/ML 🤔

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA
SEIT 1737

# Variational family

- What is the $\mathcal{Q}$ and thus $q(\boldsymbol{\theta}|\phi)$ ?
- Start simple

$$q(\boldsymbol{\theta}|\phi) = \prod_{j=1}^{J} q_j(\theta_j|\ \phi_j)$$

- Known as mean-field variational family
- $q_j$ factors in the variational distribution
  - Some parameteric distribution that respects the parameter space of $\boldsymbol{\theta}_j$
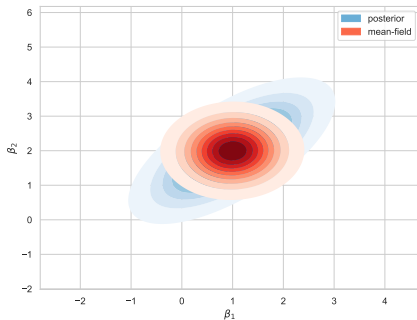  - All model parameters are independent

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

Figure 1: Mean-field approximation to a 2-D multivariate normal distribution, based on Blei et al. (ibid., p. 9, figure 1).

- Augment the variational distribution to blocks of parameters
- Structured mean-field variational inference (Wainwright and Jordan 2007)

$$q(\boldsymbol{\theta}|\boldsymbol{\phi}) = \prod_{j=1}^{J} q_j(\boldsymbol{\theta}_j|\ \boldsymbol{\phi}_j)$$

- Captures interdependencies for blocks of parameters

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

# Coordinate ascent variational inference (side node)

- Traditional way of solving opt. objective is CAVI (Blei et al. 2017)
  - However CAVI does not scale well
  - Closely connected to Gibbs sampling
  - Only works for conditional conjugate models
  - If you search for VI CAVI is still all over the place, see f.e. wikipedia
- Of course we want to be able to also conduct inference in non-conjugate models ✗

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

# Stochastic variational inference

- SGD to optimize the ELBO (Hoffman et al. 2012)
- 2 sources of stochasticity
  - We use a subset $\mathcal{I}$ of the data in each iteration (ELBO remains unbiased)
  - We need to evaluate the integral in the ELBO

$$\mathsf{ELBO}(\phi)_{\mathcal{I}} = \int q(\boldsymbol{\theta}|\phi) \ln \left( \frac{p(\mathbf{y}_{\mathcal{I}}, \boldsymbol{\theta}|\mathcal{D}_{\mathcal{I}})}{q(\boldsymbol{\theta}|\phi)} \right) d\boldsymbol{\theta} \tag{16}$$

$$\nabla_{\phi} \mathsf{ELBO}(\phi) = \nabla_{\phi} \int q(\boldsymbol{\theta}|\phi) \ln \left( \frac{p(\mathbf{y}_{\mathcal{I}}, \boldsymbol{\theta}|\mathcal{D}_{\mathcal{I}})}{q(\boldsymbol{\theta}|\phi)} \right) d\boldsymbol{\theta} \tag{17}$$

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

- Common method to solve this problem is Monte Carlo integration

$$\nabla_\phi \text{ELBO}(\phi)_\mathcal{I} = \nabla_\phi \int q(\boldsymbol{\theta}|\phi) \ln\left(\frac{p(\mathbf{y}_\mathcal{I}, \boldsymbol{\theta}|\mathcal{D}_\mathcal{I})}{q(\boldsymbol{\theta}|\phi)}\right) d\boldsymbol{\theta} \tag{18}$$

$$\approx \nabla_\phi \frac{1}{S} \sum_{s=1}^{S} \ln\left(\frac{p(\mathbf{y}_\mathcal{I}, \boldsymbol{\theta}^s|\mathcal{D}_\mathcal{I})}{q(\boldsymbol{\theta}^s|\phi)}\right), \ \boldsymbol{\theta}^s \sim q(\boldsymbol{\theta}|\phi). \tag{19}$$

- But this does not work 😢
- If we change $\phi$ even infinitessimal the samples $\boldsymbol{\theta}^s$ are invalid, which we used to calculate $\nabla_\phi \text{ELBO}(\phi)_\mathcal{I}$ in the first place

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA
SEIT 1737

Figure 2: Plate notation of the dependence structure in a Bayesian regression model for VI, when sampling from the variational distribution.

# Reparameterization gradient estimator

- Reparapemeterize $\boldsymbol{\theta} = \mathbf{g}_{\phi}(\boldsymbol{\epsilon})$, with a bijective function $\mathbf{g}_{\phi}$ such that we can use SGD with Monte Carlo integration (Kingma and Welling 2013; Kucukelbir et al. 2016; Rezende et al. 2014)
- Amounts to using the (multivariate) change of variable theorem for probability density functions

$$\mathbf{g}_{\phi} : \mathbb{R}^d \rightarrow \mathbb{R}^d; \ \boldsymbol{\theta}, \boldsymbol{\epsilon} \in \mathbb{R}^d$$

$$\boldsymbol{\theta} = \mathbf{g}_{\phi}(\boldsymbol{\epsilon}), \ \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$q(\boldsymbol{\theta}|\phi) = \begin{cases} p_{\boldsymbol{\epsilon}}(\mathbf{g}_{\phi}^{-1}(\boldsymbol{\theta})) \left| \det(\mathbf{J}_{\mathbf{g}_{\phi}^{-1}}) \right|, & \text{if } \boldsymbol{\theta} \text{ is in the codomain of } \mathbf{g}_{\phi} \\ 0, & \text{else.} \end{cases}$$

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN

Figure 3: Plate notation of the dependence structure in a Bayesian regression model in VI, using the "reparameterization-trick".

- This "trick" allows us to pull the gradient operator inside of the monte carlo integral and use the chain rule

$$\nabla_\phi \mathsf{ELBO}(\phi)_\mathcal{I} \approx \frac{1}{S} \sum_{s=1}^{S} \nabla_\phi \ln \left( \frac{p(\mathbf{y}_\mathcal{I}, \boldsymbol{\theta} = \mathbf{g}_\phi(\boldsymbol{\epsilon}) | \mathcal{D}_\mathcal{I})}{q(\boldsymbol{\theta} = \mathbf{g}_\phi(\boldsymbol{\epsilon}) | \phi)} \right) \Bigg|_{\boldsymbol{\theta} = \boldsymbol{\theta}^s}$$

$$\approx \frac{1}{S} \sum_{s=1}^{S} \nabla_{\boldsymbol{\theta}} \ln \left( \frac{p(\mathbf{y}_\mathcal{I}, \boldsymbol{\theta} = \mathbf{g}_\phi(\boldsymbol{\epsilon}) | \mathcal{D}_\mathcal{I})}{q(\boldsymbol{\theta} = \mathbf{g}_\phi(\boldsymbol{\epsilon}) | \phi)} \right) \Bigg|_{\boldsymbol{\theta} = \boldsymbol{\theta}^s} \nabla_\phi \mathbf{g}_\phi(\boldsymbol{\epsilon}^s),$$

with $\boldsymbol{\theta}_\phi^s = \mathbf{g}_\phi(\boldsymbol{\epsilon}^s),\ \boldsymbol{\epsilon}^s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),\ s = 1, \ldots, S$

- Opens the door for backpropagation and thus automatic diff. 💡

# "Black-box" variational inference

- Using SVI with the reparameterization gradient estimator
- Researcher only formulates a probabilistic model and provides a dataset (Kucukelbir et al. 2016), inference algo. is model "agnostic"
- What is $\mathbf{g}_\phi$?
  - Linear and non-linear choices
  - We consider linear choice

$$\boldsymbol{\theta}_j = \mathbf{L}_j \boldsymbol{\epsilon}_j + \boldsymbol{\mu}_j, \boldsymbol{\epsilon}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\boldsymbol{\theta}_j \sim \mathcal{N}(\boldsymbol{\mu}_j, \mathbf{L}_j \mathbf{L}_j^{\mathrm{T}})$$

- For positive restricted parameters we need to chain another transformation layer via exp transformation
- Chaining variable transformations 💬
  - Normalizing flows (Rezende et al. 2014)
  - Allow for expressive $\mathcal{Q}$s but are more difficult to optimize

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

# Optimization

- We use Adam (Kingma and Ba 2014)

$$\hat{\phi}^{t} = \hat{\phi}^{t-1} + \rho_t \nabla_\phi \mathsf{ELBO}(\phi)_{\mathcal{I}^t}\Big|_{\phi = \hat{\phi}^{t-1}}$$

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA

# Full Algorithm

**Data:** $\mathcal{D}_{\text{train}}$; $\mathcal{D}_{\text{val}}$

**Require:** Learning rate $\alpha$; stopping threshold $\varepsilon$; mini-batch size $M$; share train $w$; num. var. samples $S$; num. epochs: $E$

Initialize $\hat{\phi}^0$; set $t = 1$

**for** $e = 1$ **to** $E$ **do**
$\quad n_{\text{train}} = |\mathcal{D}_{\text{train}}| * w$
$\quad$ create the mini-batches, $\mathcal{B} = \{\ldots, \mathcal{I}^k, \ldots\}$, $k = 1, \ldots, n_{\text{train}} // M \ (+1)$
$\quad$ **for** $k = 1$ **to** $n_{train} // M \ (+1)$ **do**
$\quad\quad$ sample noise, $\epsilon_j^s \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $s = 1, \ldots, S$, $\forall j$
$\quad\quad$ calculate approx. gradient, $\nabla_{\boldsymbol{\phi}}\text{ELBO}(\boldsymbol{\phi})_{\mathcal{I}^k}\big|_{\boldsymbol{\phi} = \hat{\phi}^{t-1}}$
$\quad\quad$ update variational parameters, $\hat{\phi}^t = \hat{\phi}^{t-1} + \rho_t \nabla_{\boldsymbol{\phi}}\text{ELBO}(\boldsymbol{\phi})_{\mathcal{I}^k}\big|_{\boldsymbol{\phi} = \hat{\phi}^{t-1}}$
$\quad\quad$ calculate approx. ELBO, $\text{ELBO}(\hat{\phi}^t)_{\mathcal{D}_{\text{val}}}$
$\quad\quad t = t + 1$
$\quad$ **end**
$\quad$ **if** $t > 200$ **then**
$\quad\quad \Delta\text{ELBO} = |\text{ELBO}(\hat{\phi}^t)_{\mathcal{D}_{\text{val}}} - \text{ELBO}(\hat{\phi}^{t-200})_{\mathcal{D}_{\text{val}}}|$
$\quad$ **else**
$\quad\quad \Delta\text{ELBO} = \infty$
$\quad$ **end**
$\quad$ **if** $\Delta ELBO < \varepsilon$ **then**
$\quad\quad$ **break**
$\quad$ **end**
**end**

**Result:** $\hat{\phi}$; $\text{ELBO}(\hat{\phi})_{\mathcal{D}_{\text{val}}}$

**Algorithm 1:** BBVI algorithm.

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN

# Impact Analysis



Figure 4: ELBO traces for 3 different SGD runs, using different initializations but the same seed. We use a batch size of 128, 64 samples from the variational distribution and a learning rate of 1e-2.
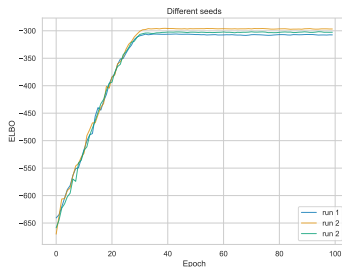


Figure 5: ELBO traces for 3 different SGD runs, using different seeds but the same initialization. Otherwise same configuration as in Figure 4.
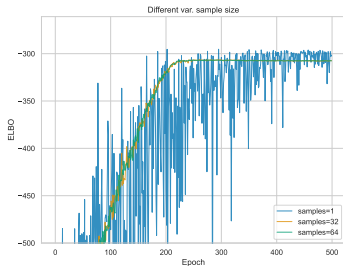
GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA

Figure 6: ELBO traces for 3 different SGD runs, using different variational sample sizes. We use batch VI with a learning rate of 1e-2 and the same seed as in Figure 4.



Figure 7: ELBO traces for 3 different SGD runs, using different batch sizes. We use a variational sample size of 64 with a learning rate of 1e-2 and the same seed as in Figure 4.

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN

# Semiparametric distributional regression

- Not only normally distributed responses
- Linear predictors with inverse link function
- Structured addtive linear predictors so fixed and smooth effects
    - B-spline basis functions
- Augmented with priors
    - Bayesian P-splines

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

# Application

# Implementation

- Developed a small `python` package `tigerpy`, which consists of two libararies
- A model building library `tigerpy.model`
  - Construct the model
  - Uses the idea of probabilistic graphical models
- An inference library `tigerpy.bbvi`
  - Runs the inference algorithm
- Aligned with concepts found in `liesel` (Riebl et al. 2022)

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA
SEIT 1737

# Technology

When walking about the countryside of Italy, the people will not hesitate to tell you that JAX has "una anima di pura programmazione funzionale". (`JAX` docs, Bradbury et al. (2018))

- There are great things about `JAX` 😍
  - Uses a `numpy` flawored API
  - Closely follows the math (`jax.grad`)
  - Is fast (if you follow `JAX`s principles)
- There are things that cause headaches 🥴
  - Pure functions
  - Tracing
  - Efficiency considerations in JIT-compiled code

- `tigerpy.model` constructs under the hood a DAG
- Employs the `networkx` package for constructing, traversing and visualizing the DAG



Figure 8: The DAG visualization for location-scale regression from the method `.visualize_graph()`.

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737

# Simulation Studies

# Open Problems

📄 Blei, David M., Alp Kucukelbir, and Jon D. McAuliffe (Apr. 2017). "Variational Inference: A Review for Statisticians". In: *Journal of the American Statistical Association* 112.518, pp. 859–877.

📄 Bradbury, James, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang (2018). *JAX: composable transformations of Python+NumPy programs*. Version 0.3.13. URL: http://github.com/google/jax.

📄 Hoffman, Matt, David M. Blei, Chong Wang, and John Paisley (2012). *Stochastic Variational Inference*.

📄 Jordan, Michael I., Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul (1999). "An Introduction to Variational Methods for Graphical Models". In: *Machine Learning* 37.2, pp. 183–233.

📄 Kingma, Diederik P and Max Welling (2013). *Auto-Encoding Variational Bayes*.

📄 Kingma, Diederik P. and Jimmy Ba (2014). *Adam: A Method for Stochastic Optimization*.

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA
SEIT 1737

📄 Kucukelbir, Alp, Dustin Tran, Rajesh Ranganath, Andrew Gelman, and David M. Blei (2016). *Automatic Differentiation Variational Inference*.

📄 Kullback, S. and R. A. Leibler (Mar. 1951). "On Information and Sufficiency". In: *The Annals of Mathematical Statistics* 22.1, pp. 79–86.

📄 Rezende, Danilo Jimenez, Shakir Mohamed, and Daan Wierstra (2014). *Stochastic Backpropagation and Approximate Inference in Deep Generative Models*.

📄 Riebl, Hannes, Paul F. V. Wiemann, and Thomas Kneib (2022). *Liesel: A Probabilistic Programming Framework for Developing Semi-Parametric Regression Models and Custom Bayesian Inference Algorithms*.

📄 Wainwright, Martin J. and Michael I. Jordan (2007). "Graphical Models, Exponential Families, and Variational Inference". In: *Foundations and Trends® in Machine Learning* 1.1–2, pp. 1–305.

GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN IN PUBLICA COMMODA SEIT 1737