



CZECH TECHNICAL UNIVERSITY IN
PRAGUE
Faculty of Nuclear Sciences and Physical
Engineering



The Role of a priori Distributions for Sparse Parametrizations of Models

Bc. Lukáš Kulička

Student's Scientific Conference 2021
Supported from CTU Grant SVK 29/21/F4

June 24, 2021

Contents

1. Motivation
2. Bayesian Approach & Shrinkage priors
3. Practical Application
4. References

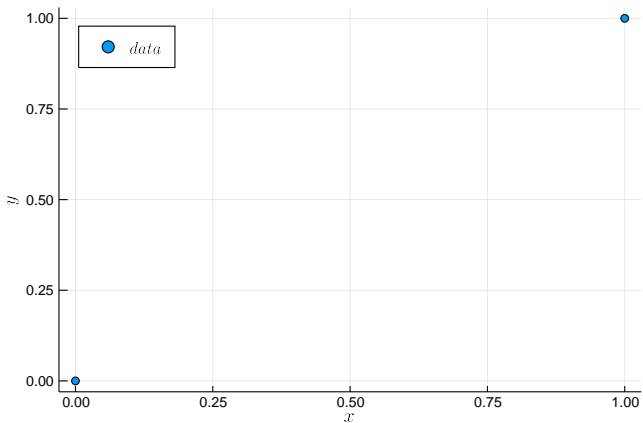


Figure 1: Example of sparse data space.

- Least squares problem:

$$\begin{aligned}y_1 &= w_1 x_{11} + w_2 x_{12}^2 \\y_2 &= w_1 x_{21} + w_2 x_{22}^2\end{aligned} \Rightarrow \mathbb{X} = \begin{pmatrix} 0 & 0 \\ 1 & 1 \end{pmatrix} \quad (1)$$

- $h(\mathbb{X}) = 1$.
- How to find $\hat{\theta}$?

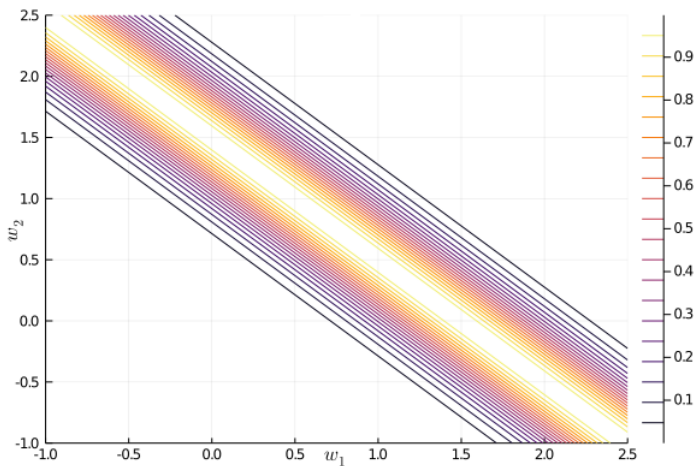


Figure 2: Contour plot of parameters likelihood w_1, w_2 from Eq. (1).

Bayesian Approach

- Bayes' rule:

$$\text{Aposteriori distribution} \propto \text{Likelihood} \times \text{Apriori distribution.}$$

- L_1 norm:

$$w_j \sim \text{Laplace}(0, \lambda) \propto \exp(-\lambda \|\boldsymbol{\theta}\|_1), \quad (2)$$

we choose λ arbitrarily firmly.

- ARD:

$$p(\boldsymbol{\theta}|\boldsymbol{\alpha}) = \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{\alpha}^{-1}\mathbb{I}) \wedge p(\alpha_j) = \text{St}(0, \sigma^2, \nu), \quad \forall j \in \hat{d} \quad (3)$$

- Spike & Slab:

$$w_j \sim \lambda_j \mathcal{N}(0, c^2) + (1 - \lambda_j) \mathcal{N}(0, \epsilon^2), \quad (4)$$

where $\epsilon \ll c$ and $\lambda_j \in \{0, 1\} \sim \text{Bernoulli}$

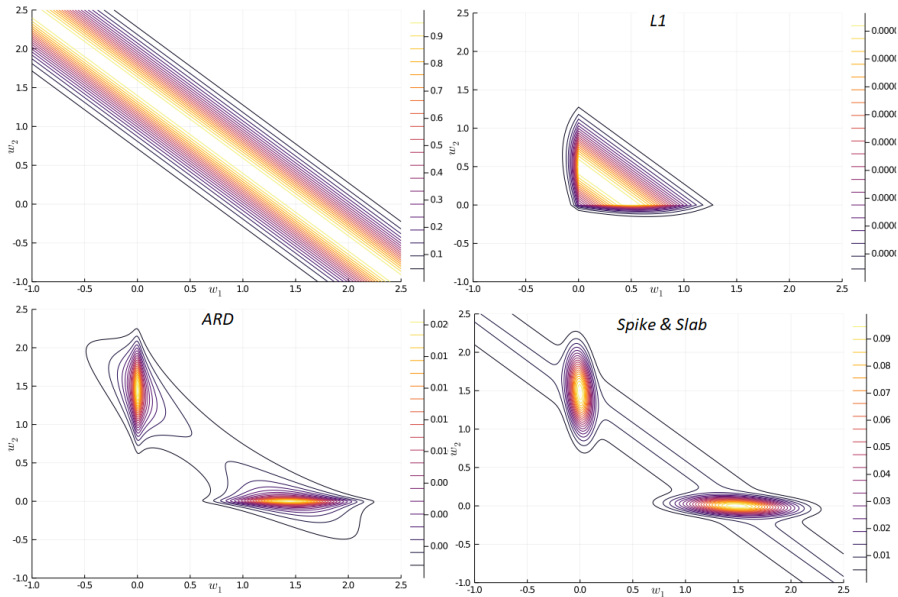


Figure 3: Contour plots of likelihood with priors (2), (3), (4).

Mathematical Tools – ELBO

- Evidence Lower Bound:

$$\begin{aligned}\ln p(\mathbf{x}) &= \int q(\mathbf{z}_\theta) \left(\ln \left\{ \frac{p(\mathbf{x}, \mathbf{z}_\theta)}{q(\mathbf{z}_\theta)} \right\} \right) d\mathbf{z}_\theta - \int q(\mathbf{z}_\theta) \left(\ln \left\{ \frac{p(\mathbf{z}_\theta|\mathbf{x})}{q(\mathbf{z}_\theta)} \right\} \right) d\mathbf{z}_\theta \\ &= \mathcal{L}(q(\mathbf{z}_\theta)) + KL(q(\mathbf{z}_\theta) \| p(\mathbf{z}_\theta|\mathbf{x}))\end{aligned}\tag{5}$$

- $p(\mathbf{z}_\theta|\mathbf{x})$ – unknown.
- $q(\mathbf{z}_\theta|\mathbf{x})$ – known, chosen.

$$\hat{\theta} = \operatorname{argmax}_{\theta} \mathcal{L}\tag{6}$$

Mathematical Tools – Reparameterization Trick

- Optimization $\mathbb{E}_{q_{\theta}(\mathbf{z})} [f(\mathbf{z})]$.

$$\mathcal{S}_{\theta}(\mathbf{z}) = \varepsilon \sim q(\varepsilon) \quad \mathbf{z} = \mathcal{S}_{\theta}^{-1}(\varepsilon) \quad (7)$$

$$\mathbb{E}_{q_{\theta}(\mathbf{z})} [f(\mathbf{z})] = \mathbb{E}_{q(\varepsilon)} [f(\mathcal{S}_{\theta}^{-1}(\varepsilon))] \quad (8)$$

- $\mathcal{N}(\mu, \sigma^2) \Rightarrow \mathcal{S}_{\mu, \sigma^2}(\mathbf{z}) = \frac{\mathbf{z} - \mu}{\sigma} \sim \mathcal{N}(0, 1)$.

$$\mathbb{E}_{\mathcal{N}(\mu, \sigma^2)} [f(\mathbf{z})] = \mathbb{E}_{\mathcal{N}(0, 1)} [f(\mu + \varepsilon\sigma)] \quad (9)$$

- Non-linear tasks.

Classic Regression

- Probability model

$$p(\mathbf{y}, \theta \mid \mathbb{X}, \alpha, \omega) = \mathcal{N}(\mathbb{X}\theta, \omega\mathbb{I})\mathcal{N}(0, \text{diag}(\alpha)) \prod_i \Gamma(\gamma_0, \delta_0) \quad (10)$$

- ETEX data
- ARD with $\alpha = 0.1$, $\omega = 5$ and $\Gamma(0.1, 0.1)$ prior.

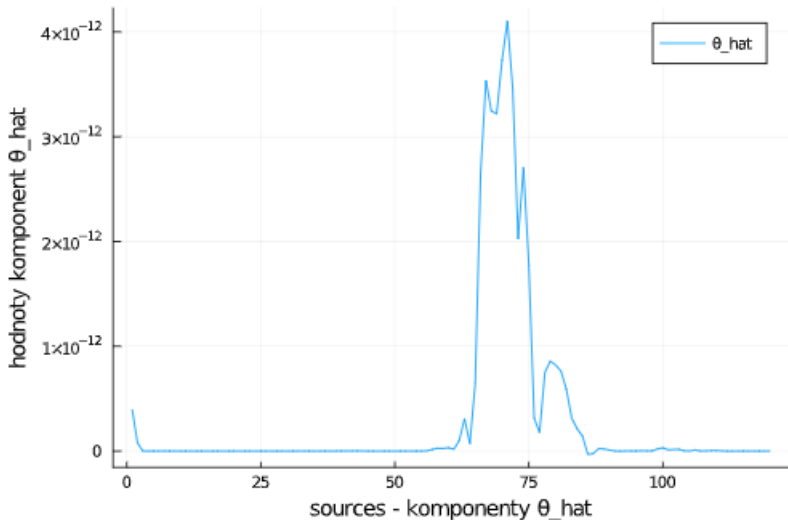


Figure 4: Finding sparse solution for ETEX data with ARD.

Logistic Regression

- 1000 generated observations from $\mathcal{N}_2(\mathbf{0}, \mathbb{I})$. Let $\boldsymbol{\theta}^{(\text{true})} = (0, 10)$.
- Generating of $\mathbf{y}^{(\text{true})}$:

$$\mathbf{y}^{(\text{true})} = \boldsymbol{\sigma} \left(\mathbb{X} \cdot \boldsymbol{\theta}^{(\text{true})T} \right) \stackrel{\text{Bernoulli}}{\Rightarrow} \text{binary-class vector.} \quad (11)$$

- Model:

$$\text{model}(\boldsymbol{\theta}) = \boldsymbol{\sigma} \left(\mathbb{X} \cdot \boldsymbol{\theta}^T \right). \quad (12)$$

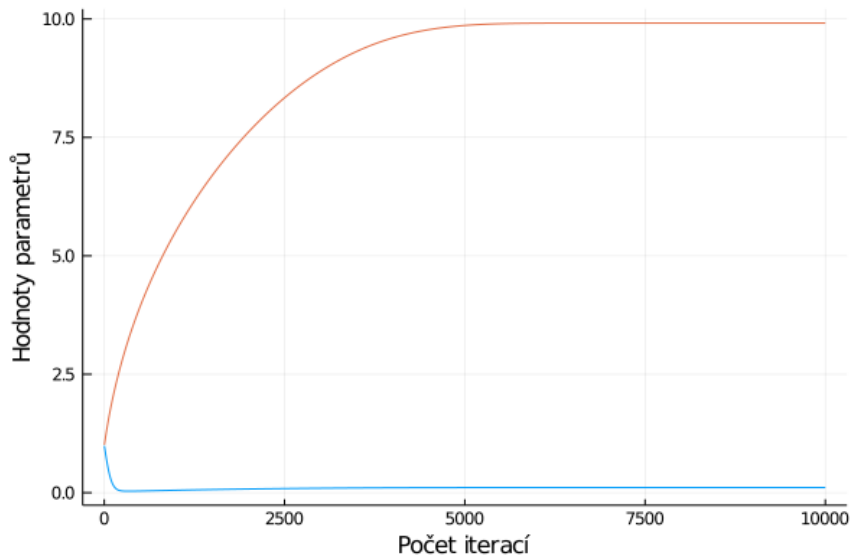


Figure 5: Learning process of θ with L_1 norm.

Tree Structures

- Model:

$$\mathbf{y} = \sigma(w_1 \mathbf{x}_1 + \max(\mathbb{W} \cdot \mathbf{Z})) \quad (13)$$

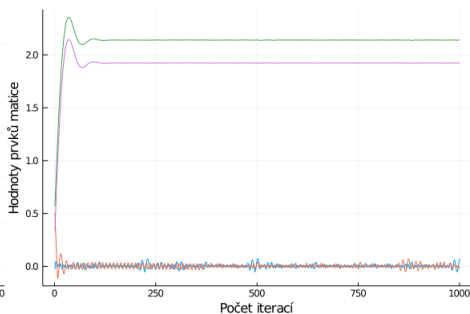
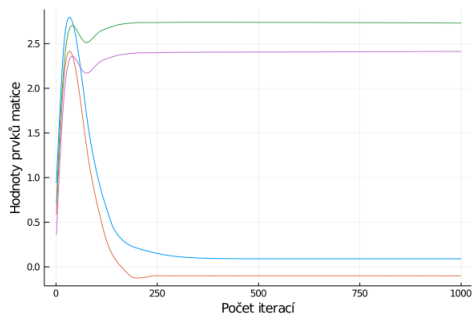


Figure 6: $\mathbb{W}^{(\text{true})} = \begin{pmatrix} 2 & 0 \\ 0 & 4 \end{pmatrix}$, right with penalty L_1 norm ($\lambda = 0.01$).

Thank you for your attention.

References

- KULIČKA, Lukáš. *Klasifikace dat popsaných stromovou strukturou*. Praha, 2020. Bakalářská práce. České vysoké učení technické v Praze, Fakulta jaderná a fyzikálně inženýrská, Katedra matematiky. Vedoucí práce doc. Ing. Václav Šmídl, Ph.D.

Author: Bc. Lukáš Kulička

Title: The Role of a priori Distributions for Sparse
Parametrizations of Models

Programme: Applied Mathematical Stochastic Methods

Student's Scientific Conference 2021

Supported from CTU Grant SVK 29/21/F4

Lukas.Kulicka@cvut.cz

