

Learning Approximately Objective Priors

Galina Boeva

MIPT, 2023

October 3, 2023

- 1 Motivation
- 2 Definition
- 3 Methods
- 4 Empirical results

Motivation

Main idea

It's difficult to have informative priors, so they turn to uninformative priors, but for many models it's impossible to deduce them. Therefore, the authors propose to approximate these prior.

Definition

Likelihood

$$p(\mathcal{D}|\theta) = \prod_{i=1}^N p(x_i|\theta), \quad (1)$$

where θ are the model parameters and \mathcal{D} is the dataset, which is comprised of N i.i.d. observations $x_i \in \mathcal{X}$. $p(\theta)$ denotes the prior, $p(\theta|\mathcal{D})$ the posterior, and $p(\mathcal{D}) = \int_{\theta} p(\mathcal{D}|\theta)p(\theta)d\theta$ the marginal likelihood.

Reference prior

Reference prior is the distribution that maximizes the mutual information between the parameters θ and the data \mathcal{D} :

$$p^*(\theta) = \arg \max_{p(\theta)} \mathbb{I}(\theta, \mathcal{D}) = \arg \max_{p(\theta)} \mathbb{H}[\theta] - \mathbb{H}[\theta|\mathcal{D}] \quad (2)$$

Definition

Mutual information

Mutual information in terms of a Kullback-Leibler divergence:

$$I(\theta, \mathcal{D}) = \int_{\mathcal{D}} p(\mathcal{D}) \text{KLD}[p(\theta|\mathcal{D})||p(\theta)] d\mathcal{D}. \quad (3)$$

$$I(\theta, \mathcal{D}) = - \int_{\theta} p(\theta) \log \frac{p(\theta)}{\bar{f}(\theta)} d\theta, \quad f(\theta) = \exp \left\{ \int_{\mathcal{D}} p(\mathcal{D}|\theta) \log p(\theta|\mathcal{D}) d\mathcal{D} \right\}. \quad (4)$$

Jeffreys Prior

$$\pi(\theta) \propto \sqrt{\det \mathcal{F}[\theta]}. \quad (5)$$

Bernstein Von Mises theorem

$$p(\theta|\mathcal{D}) \approx N(\theta_{MLE}, \mathcal{F}^{-1}(\theta)) \quad (6)$$

Information lower bound

Argmax λ

$$\begin{aligned}\lambda_* &= \arg \max_{\lambda} \mathbb{I}(\theta, \mathcal{D}) = \arg \max_{\lambda} \int_{\theta} p_{\lambda}(\theta) \int_{\mathcal{D}} p(\mathcal{D}|\theta) \log \frac{p(\mathcal{D}, \theta)}{p_{\lambda}(\theta)p(\mathcal{D})} d\mathcal{D} d\theta \\ &= \arg \max_{\lambda} \int_{\theta} p_{\lambda}(\theta) \int_{\mathcal{D}} p(\mathcal{D}|\theta) \log \frac{p(\mathcal{D}, \theta)}{p(\mathcal{D})} d\mathcal{D} d\theta = \\ &= \arg \max_{\lambda} \mathbb{E}_{\theta_{\lambda}} [-\mathbb{H}_{\mathcal{D}|\theta}[\mathcal{D}] - \mathbb{E}_{\mathcal{D}|\theta}[\log p(\mathcal{D})]]\end{aligned}\quad (7)$$

Bounding $\log p(\mathcal{D})$

Using variational Renyi bound:

$$\log p(\mathcal{D}) \leq \frac{1}{1-\alpha} \log \mathbb{E}_{\theta} (p(\mathcal{D}|\theta)^{(1-\alpha)}), \alpha \leq 0. \quad (8)$$

Information lower bound

Lower bound for λ

$$\mathbb{I}(\theta, \mathcal{D}) \geq \mathbb{E}_{\theta_\lambda} \left[-\mathbb{H}_{\mathcal{D}|\theta}[\mathcal{D}] - \frac{1}{1-\alpha} \log \mathbb{E}_\theta [p(\mathcal{D}|\theta)^{1-\alpha}] \right]. \quad (9)$$

$$\mathbb{I}(\theta, \mathcal{D}) \geq \mathbb{J}_{\mathcal{RP}}(\lambda) = \mathbb{E}_{\theta_\lambda} \left[-\mathbb{H}_{\mathcal{D}|\theta}[\mathcal{D}] - \mathbb{E}_{\mathcal{D}|\theta} [\max_s \log p(\mathcal{D}|\hat{\theta}_s)] \right]. \quad (10)$$

$$\mathbb{J}_{\mathcal{RP}}(\lambda) = \mathbb{E}_{\theta_\lambda} \mathbb{E}_{\mathcal{D}|\theta} \left[\log p(\mathcal{D}|\theta) - \max_s \log p(\mathcal{D}|\hat{\theta}_s) \right]. \quad (11)$$

$$\begin{aligned} \hat{\mathbb{J}}_{\mathcal{RP}}(\lambda) &= \frac{1}{S} \sum_{s=1}^S \mathbb{H}[p(\mathcal{D}|\hat{\theta}_s) || p(\mathcal{D}|\hat{\theta}_s)] - \mathbb{H}_{\mathcal{D}|\hat{\theta}_s}[\mathcal{D}] = \\ &= \frac{1}{S} \sum_{s=1}^S \text{KLD}[p(\mathcal{D}|\hat{\theta}_s) || p(\mathcal{D}|\hat{\theta}_s)] \end{aligned} \quad (12)$$

Particle descent

Stein Variational Gradient Descent

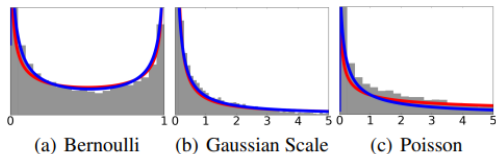
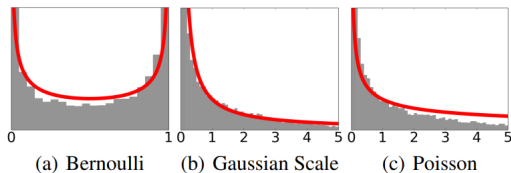
$$\bar{\theta}_j^{t+1} = \bar{\theta}_j^t + \eta \phi[\theta]$$

$$\phi[\theta] = \frac{1}{K} \sum_{k=1}^K \kappa(\bar{\theta}_k^t, \bar{\theta}_j^t) \nabla_{\bar{\theta}_k} \log p(\bar{\theta}_k^t) + \nabla_{\bar{\theta}_k} \kappa(\bar{\theta}_k^t, \bar{\theta}_j^t).$$

A-SVGD for RP Approximations

$$\begin{aligned} \nabla_{\bar{\theta}} \log f(\bar{\theta}) &= \nabla_{\bar{\theta}} \int_{\mathcal{D}} p(\mathcal{D}|\bar{\theta}) \log p(\bar{\theta}|\mathcal{D}) d\mathcal{D} = \\ &= \int_{\mathcal{D}} p(\mathcal{D}|\bar{\theta}) \log \frac{p(\mathcal{D}|\bar{\theta}) 1_{\Theta}}{p(\mathcal{D})} d\mathcal{D} = -\nabla_{\bar{\theta}} \mathbb{H}_{\mathcal{D}|\bar{\theta}}[\mathcal{D}] - \nabla_{\bar{\theta}} \mathbb{E}_{\mathcal{D}|\bar{\theta}}[\log p(\mathcal{D})] \\ &= \nabla_{\bar{\theta}} \frac{1}{S} \sum_s \text{KLD}[p(\mathcal{D}|\bar{\theta}) || p(\mathcal{D}|\hat{\theta}_s)]. \end{aligned}$$

Empirical results: Approximation



Empirical results: Approximation Quality

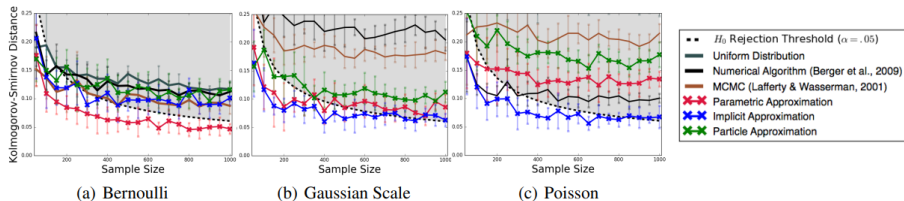


Figure: Quantifying the Approximation Quality.

Empirical results: VAE

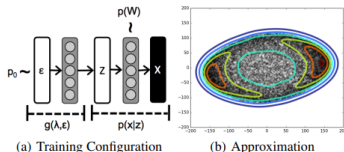


Figure: Learning the Variational Autoencoder's Reference Prior. (a) computational pipeline from the implicit prior through the VAE decoder; (b) RP approximation (contours are generated via kernel density estimation on 10, 000 samples).

- 1 **Main article** Learning Approximately Objective Priors.