

Bayesian multimodeling: Variational inference-2

MIPT

2022

Model selection: coherent Bayesian inference

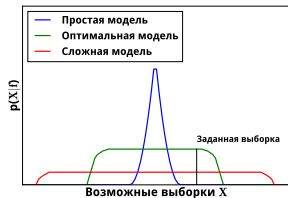
First level: find optimal parameters:

$$w = \arg \max \frac{p(\mathcal{D}|w)p(w|h)}{p(\mathcal{D}|h)},$$

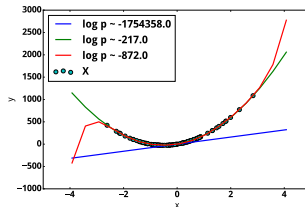
Second level: find optimal model:

Evidence:

$$p(\mathcal{D}|h) = \int_w p(\mathcal{D}|w)p(w|h)dw.$$



Model selection scheme



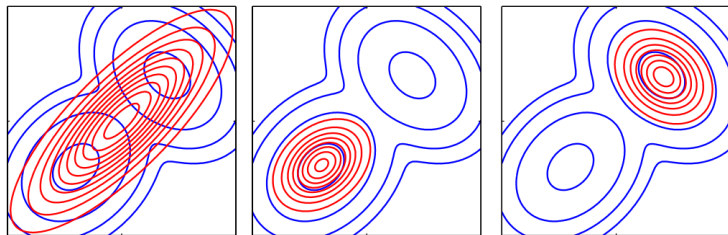
Polynomial regression example

Evidence lower bound, ELBO

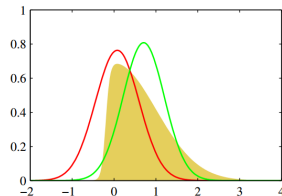
Evidence lower bound is a method of approximation of intractable distribution $p(w|\mathcal{D}, h)$ with a distribution $q(w) \in \mathcal{Q}$.

Evidence lower bound estimation often reduces to optimization problem

$$\log p(\mathcal{D}|h) \geq \text{KL}(q(w)||p(w|\mathcal{D})) = - \int_w q(w) \log \frac{p(w|\mathcal{D})}{q(w)} dw = E_w \log p(\mathcal{D}|w) - \text{KL}(q(w)||p(w|h))$$



Variational inference vs. expectation propagation (Bishop)



Laplace Approximation vs
Variational inference

ELBO estimation

ELBO maximization

$$\int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{y}, \mathbf{w} | \mathbf{X}, \mathbf{h})}{q(\mathbf{w})} d\mathbf{w}$$

is equivalent to KL-divergence minimization between $q(\mathbf{w}) \in \mathfrak{Q}$ and posteriod distribution $p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathbf{h})$:

$$\hat{q} = \arg \max_{q \in \mathfrak{Q}} \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{y}, \mathbf{w} | \mathbf{X}, \mathbf{h})}{q(\mathbf{w})} d\mathbf{w} \Leftrightarrow$$

$$\hat{q} = \arg \min_{q \in \mathfrak{Q}} D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathbf{h})),$$

$$D_{\text{KL}}(q(\mathbf{w}) || p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathbf{h})) = \int_{\mathbf{w}} q(\mathbf{w}) \log \left(\frac{q(\mathbf{w})}{p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathbf{h})} \right) d\mathbf{w}.$$

MCMC and variational inference

MCMC idea: Sample from the simple distribution and accept them, if the ratio is greater than some threshold:

$$\min \left(1, \frac{p(w^\tau | y, X, h)}{p(w^{\tau-1} | y, X, h)} \right),$$

where w^τ is set based on the previous sample:

$$w^\tau = T(w^{\tau-1}).$$

Salimans et al., 2014: let's interpret the sequence of some operator T application as a variational optimization:

$$T^1 \circ \dots \circ T^\eta(w) \rightarrow p(w^\tau | y, X, h).$$

Maclaurin et. al, 2015: use gradient descent as such operator. Do not reject samples at all.

Optimization operator, Maclaurin et. al, 2015

Definition

Let T be an algorithm of changing model parameters w' using previous parameter values w :

$$w' = T(w).$$

Definition

Let L be a continuous loss function.

Define a gradient descent operator in the following way:

$$T(w) = w - \beta \nabla L(w, y, \mathcal{D}).$$

Gradient descent for evidence estimation

Consider posterior probability maximization:

$$L = -\log p(\mathfrak{D}, w|h) = - \sum_{\mathfrak{D} \in \mathfrak{D}} \log p(\mathfrak{D}|w, h)p(w|h)$$

Optimize neural network in a multi-start regime with r initial parameter values w_1, \dots, w_r using (stochastic) gradient descent:

$$w' = T(w).$$

The parameter vectors w_1, \dots, w_r are from some latent distribution $q(w)$.

Entropy

We can rewrite variational inference using differential entropy term:

$$\log p(\mathfrak{D}|\mathbf{f}) \geq \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathfrak{D}, \mathbf{w}|\mathbf{h})}{q(\mathbf{w})} d\mathbf{w} = \\ \mathbb{E}_{q(\mathbf{w})}[\log p(\mathfrak{D}, \mathbf{w}|\mathbf{h})] + S(q(\mathbf{w})),$$

where $S(q(\mathbf{w}))$ is a differential entropy:

$$S(q(\mathbf{w})) = - \int_{\mathbf{w}} q(\mathbf{w}) \log q(\mathbf{w}) d\mathbf{w}.$$

Gradient descent for evidence estimation

Statement

Let L be a Lipschitz function, and optimization operator be a bijection. Then entropy difference for two steps is:

$$S(q'(w)) - S(q(w)) \simeq \frac{1}{r} \sum_{g=1}^r (-\beta \text{Tr}[H(w'^g)] - \beta^2 \text{Tr}[H(w'^g)H(w'^g)]).$$

Final estimation for the τ optimization step:

$$\begin{aligned} \log \hat{p}(Y|\mathcal{D}, h) &\sim \frac{1}{r} \sum_{g=1}^r L(w_{\tau}^g, \mathcal{D}, Y) + S(q^0(w)) + \\ &+ \frac{1}{r} \sum_{b=1}^{\tau} \sum_{g=1}^r (-\beta \text{Tr}[H(w_b^g)] - \beta^2 \text{Tr}[H(w_b^g)H(w_b^g)]), \end{aligned}$$

w_b^g is a parameter vector for optimization g on the step b , $S(q^0(w))$ is an initial entropy.

How to calculate Hessian trace?

Problem

$$\text{Tr}[H(w_b^g)]$$

Statement

Let U be a symmetric matrix and v be the random vector with the following properties:

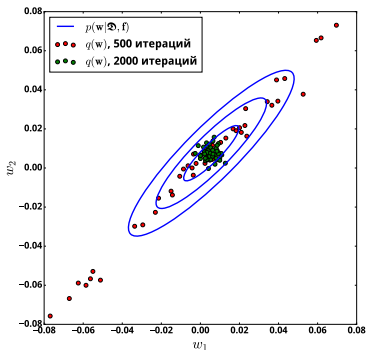
- ① $E v_i = 0$;
- ② $\text{Var}(v_i) = 1$.

Then

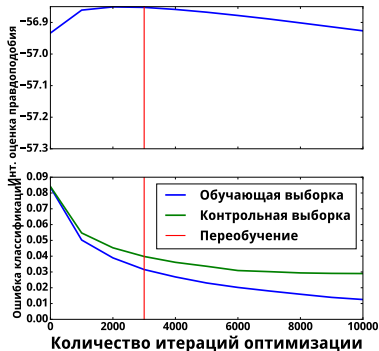
$$E v^T U v = \text{Tr}[U].$$

Overfitting, Maclaurin et. al, 2015

Gradient descent does not optimize KL-divergence $KL(q(w)||p(w|\mathcal{D}, h))$. Evidence estimation gets worse while optimization tends to the optimal parameter values. This can be considered as a overfitting start.



Convergence



Overfitting start

Stochastic gradient Langevin dynamics

A modification of SGD:

$$T = w - \beta \nabla L + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \frac{\beta}{2})$$

where β changes with a number of iterations:

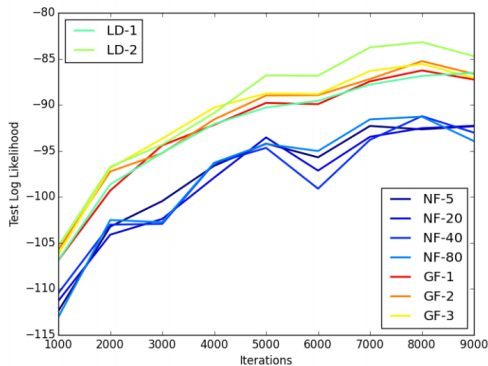
$$\sum_{\tau=1}^{\infty} \beta_{\tau} = \infty, \quad \sum_{\tau=1}^{\infty} \beta_{\tau}^2 < \infty.$$

Statement [Welling, 2011]. Distribution $q^{\tau}(w)$ converges to posterior distribution $p(w|X, f)$.
Entropy adjustment:

$$\hat{S}(q^{\tau}(w)) \geq \frac{1}{2} |w| \log \left(\exp \left(\frac{2S(q^{\tau}(w))}{|w|} \right) + \exp \left(\frac{2S(\epsilon)}{|w|} \right) \right).$$

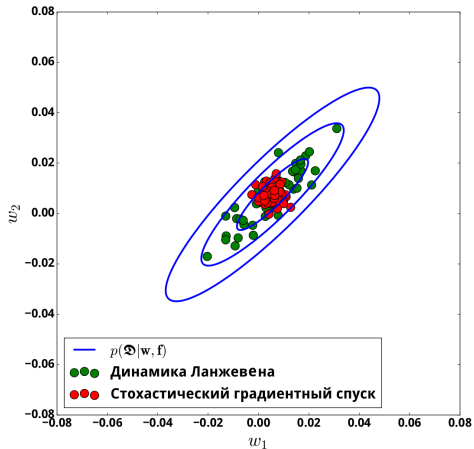
Stochastic gradient Langevin dynamics for generative models

Altieri et al., 2015: sample latent variable z and use SGLD as a normalizing flow.



SGLD vs SGD

Parameter distribution after 2000 iterations:



Reparametrization trick: problems

Reparamterization idea:

$$\varepsilon = S_{\theta}(w), \quad w = S_{\theta}^{-1}(\varepsilon).$$

Then:

$$\nabla_{\theta} E_q f(w) = E_q \nabla_{\theta} f(S_{\theta}^{-1}(\varepsilon)) = E_q \nabla_w f(S_{\theta}^{-1}(\varepsilon)) \nabla_{\theta} S^{-1}(\varepsilon).$$

Example:

$$w \sim \mathcal{N}(\mu, \sigma^2) \rightarrow S(w) = \frac{w - \mu}{\sigma} \sim \mathcal{N}(0, 1).$$

Challenge: calculation of S^{-1} is an expensive operation.

Implicit reparametrization trick

$$\nabla_{\theta} E_q f(w) = E_q \nabla_w f(w) \nabla_{\theta} w.$$

Use a total gradient formula for $\varepsilon = S_{\theta}(w)$:

$$\nabla_w S_{\theta}(w) \nabla_{\theta} w + \nabla_{\theta} S_{\theta}(w) = 0 \rightarrow$$

$$\rightarrow \nabla_{\theta} w = -(\nabla_w S_{\theta}(w))^{-1} \nabla_{\theta} S_{\theta}.$$

Obtain an expression without inverse function for S .

For 1d samples we can use, for example:

$$S(w) = F(w|\theta) \sim \mathcal{U}(0, 1).$$

Table 4: Test negative log-likelihood (lower is better) for VAE on MNIST. Mean \pm standard deviation over 5 runs. The von Mises-Fisher results are from [9].

Prior	Variational posterior	$D = 2$	$D = 5$	$D = 10$	$D = 20$	$D = 40$
$\mathcal{N}(0, 1)$	$\mathcal{N}(\mu, \sigma^2)$	131.1 ± 0.6	107.9 ± 0.4	92.5 ± 0.2	88.1 ± 0.2	88.1 ± 0.0
Gamma(0.3, 0.3)	Gamma(α, β)	132.4 ± 0.3	108.0 ± 0.3	94.0 ± 0.3	90.3 ± 0.2	90.6 ± 0.2
Gamma(10, 10)	Gamma(α, β)	135.0 ± 0.2	107.0 ± 0.2	92.3 ± 0.2	88.3 ± 0.2	88.3 ± 0.1
Uniform(0, 1)	Beta(α, β)	128.3 ± 0.2	107.4 ± 0.2	94.1 ± 0.1	88.9 ± 0.1	88.6 ± 0.1
Beta(10, 10)	Beta(α, β)	131.1 ± 0.4	106.7 ± 0.1	92.1 ± 0.2	87.8 ± 0.1	87.7 ± 0.1
Uniform($-\pi, \pi$)	vonMises(μ, κ)	127.6 ± 0.4	107.5 ± 0.4	94.4 ± 0.5	90.9 ± 0.1	91.5 ± 0.4
vonMises(0, 10)	vonMises(μ, κ)	130.7 ± 0.8	107.5 ± 0.5	92.3 ± 0.2	87.8 ± 0.2	87.9 ± 0.3
Uniform(S^D)	vonMisesFisher($\boldsymbol{\mu}, \kappa$)	132.5 ± 0.7	108.4 ± 0.1	93.2 ± 0.1	89.0 ± 0.3	90.9 ± 0.3

Model selection: coherent Bayesian inference

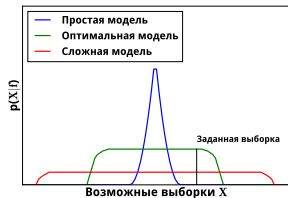
First level: find optimal parameters:

$$w = \arg \max \frac{p(\mathcal{D}|w)p(w|h)}{p(\mathcal{D}|h)},$$

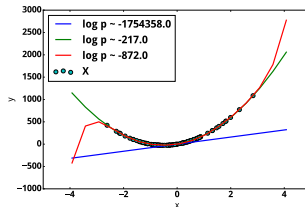
Second level: find optimal model:

Evidence:

$$p(\mathcal{D}|h) = \int_w p(\mathcal{D}|w)p(w|h)dw.$$

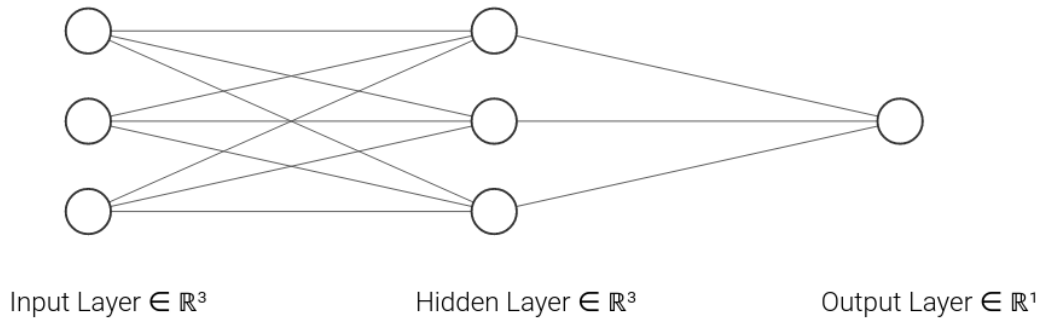


Model selection scheme

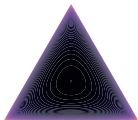


Polynomial regression example

Discrete variational optimization



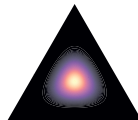
Discrete distribution: relaxation



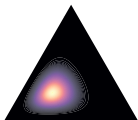
$$\bar{\alpha} = [1, 1, 1], t = 0.9$$



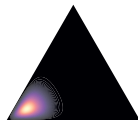
$$t = 1.0$$



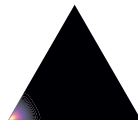
$$t = 10.0$$



$$\bar{\alpha} = [0.5, 0.25, 0.25], t = 30.0$$



$$[0.75, 0.125, 0.125]$$



$$[0.9, 0.05, 0.05]$$

Discrete distribution in variational inference

Relaxation:

- Dirichlet distribution (+ Implicit reparametrization trick)
- Gumbel-softmax:

$$p(w) = \Gamma(k)\tau^{k-1} \left(\sum_{i=1}^k \alpha_i / w_i \right)^{-k} \prod_{i=1}^k (\alpha_i / w_i \tau + 1)$$

- ▶ Reparameterization works well
 - ▶ KL divergence is intractable
- Invertible Gaussian reparametrization:

$$p(w) = \text{softmax}(\alpha), \quad \alpha \sim \mathcal{N},$$

(there should be a ϵ in the denominator for invertibility of the function)

- ▶ Reparameterization works well
 - ▶ $KL(w_1|w_2) = KL(\alpha_1|\alpha_2)$
 - ▶ Poor interpretation

Local reparametrization

Let $y = \text{ReLU}(XW)$ and parameter matrix W be distributed normally: $w_{i,j} \sim \mathcal{N}(\mu_{i,j}, \sigma_{i,j}^2)$.
Then XW is a Gaussian matrix:

$$G = XW, \quad G_{i,j} \sim \mathcal{N}\left(\sum_k x_{i,k} \mu_{k,j}, \sum_k x_{i,k}^2 \sigma_{k,j}^2\right).$$

Instead of sampling parameters, sample elements from G (units after activation).

Using CLT:

$$\sum_k x_{i,k} w_{k,j} \sim \mathcal{N}(\cdot, \cdot).$$

Conclusion: we can use the local reparametrization for discrete parameters too.

Rényi divergence

$$D_{\alpha}(p(w)||q(w)) = \frac{1}{\alpha - 1} \log \int p(w)^{\alpha} q(w)^{1-\alpha} dw.$$

Table 1: Special cases in the Rényi divergence family.

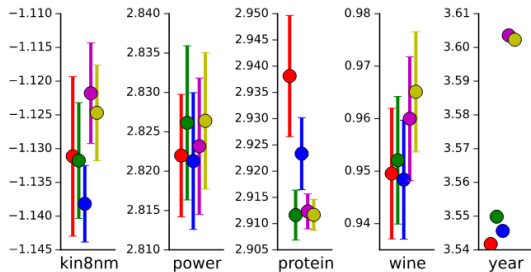
α	Definition	Notes
$\alpha \rightarrow 1$	$\int p(\boldsymbol{\theta}) \log \frac{p(\boldsymbol{\theta})}{q(\boldsymbol{\theta})} d\boldsymbol{\theta}$	<i>Kullback-Leibler (KL) divergence</i> , used in VI ($\text{KL}[q p]$) and EP ($\text{KL}[p q]$)
$\alpha = 0.5$	$-2 \log(1 - \text{Hel}^2[p q])$	function of the square <i>Hellinger distance</i>
$\alpha \rightarrow 0$	$-\log \int_{p(\boldsymbol{\theta}) > 0} q(\boldsymbol{\theta}) d\boldsymbol{\theta}$	zero when $\text{supp}(q) \subseteq \text{supp}(p)$ (not a divergence)
$\alpha = 2$	$-\log(1 - \chi^2[p q])$	proportional to the χ^2 -divergence
$\alpha \rightarrow +\infty$	$\log \max_{\boldsymbol{\theta} \in \Theta} \frac{p(\boldsymbol{\theta})}{q(\boldsymbol{\theta})}$	<i>worst-case regret</i> in <i>minimum description length principle</i> [24]

mass-covering

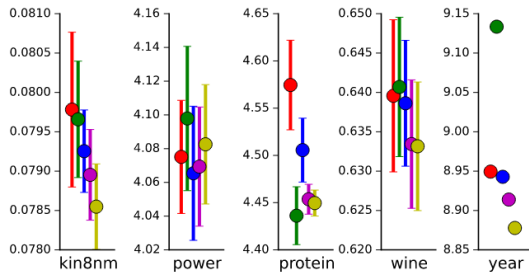
zero-forcing

● $\alpha \rightarrow -\infty$ (max) ● $\alpha=0.0$ ● $\alpha=0.5$ ● $\alpha=1.0$ (VI) ● $\alpha \rightarrow +\infty$

average negative test LL/nats



average test RMSE



References

- Bishop C. M. Pattern recognition //Machine learning. – 2006. – T. 128. – №. 9.
- MacKay D. J. C., Mac Kay D. J. C. Information theory, inference and learning algorithms. – Cambridge university press, 2003.
- Salimans, Tim, Diederik Kingma, and Max Welling, 2015. Markov chain monte carlo and variational inference: Bridging the gap
- Altieri: <http://approximateinference.org/accepted/AltieriDuvenaud2015.pdf>
- Stephan Mandt, Matthew D. Hoffman, David M. Blei, 2017. Stochastic Gradient Descent as Approximate Bayesian Inference
- Бахтеев О. Ю., Стрижов В. В. Выбор моделей глубокого обучения субоптимальной сложности //Автоматика и телемеханика. – 2018. – №. 8. – С. 129-147.
- Figurnov M., Mohamed S., Mnih A. Implicit reparameterization gradients //arXiv preprint arXiv:1805.08498. – 2018.
- Jang E., Gu S., Poole B. Categorical reparameterization with gumbel-softmax //arXiv preprint arXiv:1611.01144. – 2016.
- Potapczynski A., Loaiza-Ganem G., Cunningham J. P. Invertible gaussian reparameterization: Revisiting the gumbel-softmax //arXiv preprint arXiv:1912.09588. – 2019.
- Maddison C. J., Mnih A., Teh Y. W. The concrete distribution: A continuous relaxation of discrete random variables //arXiv preprint arXiv:1611.00712. – 2016.
- Shayer O., Levi D., Fetaya E. Learning discrete weights using the local reparameterization trick //arXiv preprint arXiv:1710.07739. – 2017.
- Li Y., Turner R. E. Rényi Divergence Variational Inference //arXiv preprint arXiv:1602.02311. – 2016.
- Hutchinson, M. F. (1990). A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines. Communications in Statistics - Simulation and Computation, 19(2), 433–450.