

Bayesian multimodeling: variational inference

MIPT

2021

Variational calculus

Variational calculus problem is to find maxima and minima of functionals: mappings from a set of functions to the real numbers.

Example

Find a PDF p that gives maximum of entropy $H = - \int_w \log p(w)p(w)dw$.

- p — function to find
- H — functional

If a function is set from a predefined set of functions, we can consider the variational calculus problem as an approximation problem.

Model selection: coherent Bayesian inference

First level: find optimal parameters:

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

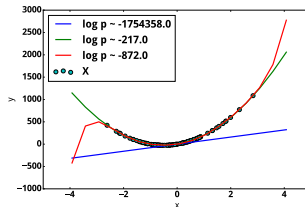
Second level: find optimal model:

Evidence:

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



Model selection scheme



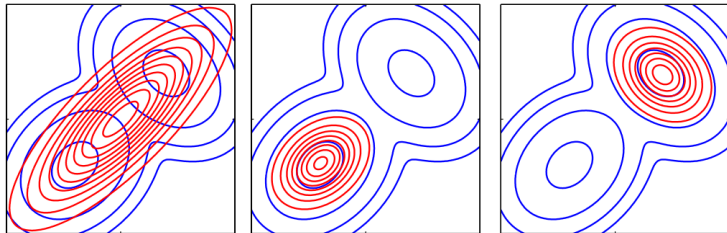
Polynomial regression example

Evidence lower bound, ELBO

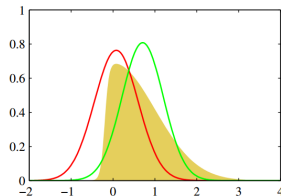
Evidence lower bound is a method of approximation of intractable distribution $p(w|\mathcal{D}, h)$ with a distribution $q(w) \in \mathcal{Q}$.

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

Minimum description length principle

$$\text{MDL}(f, \mathcal{D}) = L(f) + L(\mathcal{D}|f),$$

where f is a model, \mathcal{D} is a dataset, L is a description length in bits.

$$\text{MDL}(f, \mathcal{D}) \sim L(f) + L(\mathbf{w}^*|f) + L(\mathcal{D}|\mathbf{w}^*, f),$$

\mathbf{w}^* — optimal parameters.

\mathbf{f}_1	$L(\mathbf{f}_1)$	$L(\mathbf{w}_1^* \mathbf{f}_1)$	$L(\mathcal{D} \mathbf{w}_1^*, \mathbf{f}_1)$
\mathbf{f}_2	$L(\mathbf{f}_2)$	$L(\mathbf{w}_2^* \mathbf{f}_2)$	$L(\mathcal{D} \mathbf{w}_2^*, \mathbf{f}_2)$
\mathbf{f}_3	$L(\mathbf{f}_3)$	$L(\mathbf{w}_3^* \mathbf{f}_3)$	$L(\mathcal{D} \mathbf{w}_3^*, \mathbf{f}_3)$

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 $p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \mathbf{h})$:

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

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

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

ELBO and sample size

Statement

Let $m \gg 0$, $\lambda > 0$, $\frac{m}{\lambda} \in \mathbb{N}$, $\frac{m}{\lambda} \gg 0$. Then optimization

$$\mathbb{E}_q \log p(y|X, w) - \lambda D_{\text{KL}}(q(w) || p(w|y, X, h))$$

is equivalent to optimization of ELBO for a random subsample \hat{y}, \hat{X} with size $\frac{m}{\lambda}$.

See also, [β -VAE, Fixing Broken ELBO].

ELBO usage

ELBO: when to use?

- Evidence estimation;
- Latent distribution estimation (topic modeling, dimension reduction).

Why ELBO?

- reduces the problem of ELBO estimation to optimization;
- scales easily (compare with Laplace approximation);
- easy to use in comparison to MC-based methods.

ELBO can give a very biased evidence estimation.

ELBO: normal distribution

Let $q \sim \mathcal{N}(\boldsymbol{\mu}_q, \mathbf{A}_q)$.

Then ELBO equals to:

$$\int_{\mathbf{w}} q(\mathbf{w}) \log p(\mathbf{Y}|\mathbf{X}, \mathbf{w}, \mathbf{h}) d\mathbf{w} - D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})) \simeq$$
$$\sum_{i=1}^m \log p(y_i|x_i, \hat{\mathbf{w}}) - D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})) \rightarrow \max_{\mathbf{A}_q, \boldsymbol{\mu}_q}, \quad \hat{\mathbf{w}} \sim q.$$

If prior $p(\mathbf{w}|\mathbf{h})$ is normal:

$$p(\mathbf{w}|\mathbf{h}) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{A}),$$

KL-divergence $D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h}))$ is computed analytically:

$$D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})) = \frac{1}{2} (\text{tr}(\mathbf{A}^{-1}\mathbf{A}_q) + (\boldsymbol{\mu} - \boldsymbol{\mu}_q)^{\text{T}} \mathbf{A}^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_q) - n + \ln |\mathbf{A}| - \ln |\mathbf{A}_q|).$$

Graves, 2011

Prior: $p(w|\sigma) \sim \mathcal{N}(\mu, \sigma I)$.

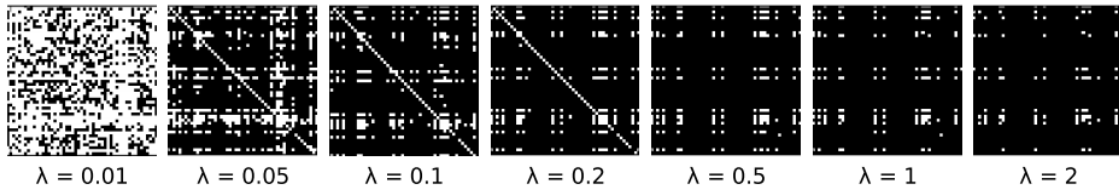
Variational distribution: $q(w) \sim \mathcal{N}(\mu_q, \sigma_q I)$.

Greedy hyperparameter optimization:

$$\mu = \hat{E}w, \quad \sigma = \hat{D}w.$$

Parameter pruning w_i using relative PDF:

$$\lambda = \frac{q(0)}{q(\mu_{i,q})} = \exp\left(-\frac{\mu_i^2}{2\sigma_i^2}\right).$$



ELBO: normal distribution

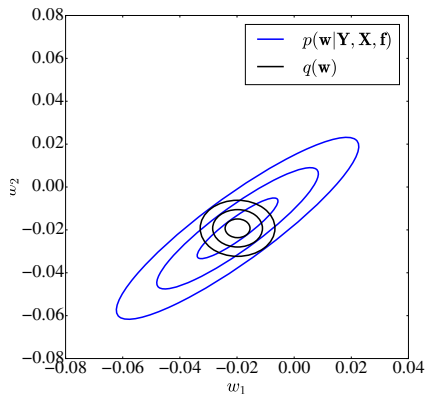
“Common” loss function:

$$L = \sum_{x,y \in \mathcal{D}} -\log p(y|x, \mathbf{w}) + \lambda \|\mathbf{w}\|_2^2.$$

Variational inference with
($p(\mathbf{w}|\mathbf{h}) \sim \mathcal{N}(0, 1)$):

$$L = \sum_{x,y} \log p(y|x, \hat{\mathbf{w}}) + \\ + \frac{1}{2} (\text{tr}(\mathbf{A}_q) + \boldsymbol{\mu}_q^T \mathbf{A}^{-1} \boldsymbol{\mu}_q - \ln |\mathbf{A}_q|).$$

Poor approximation example q



Local reparametrization

How to calculate $E_q \log p(y|X, w)$?

- Graves, 2011: 1 sample per iteration. Use the following properties:

$$w \sim \mathcal{N}(\mu, \sigma^2) \rightarrow w \sim \varepsilon \sigma + \mu, \quad \varepsilon \sim \mathcal{N}(0, 1).$$

- ▶ Poor expectation approximation
- Naive solution: sample 1 iteration per element in batch
 - ▶ BackProp will be very slow

Local reparametrization, Kingma et al., 2015

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:

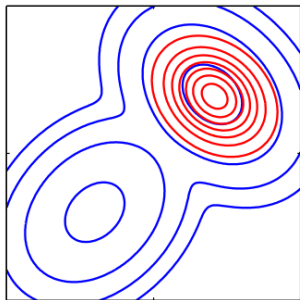
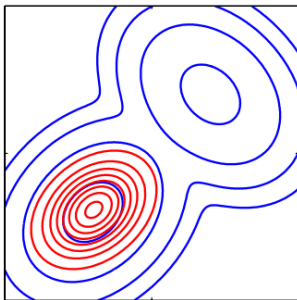
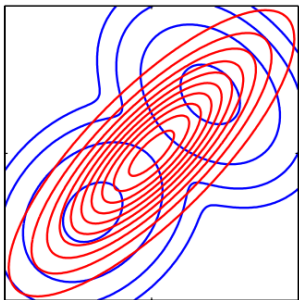
$$G = XW, \quad G_{i,j} \sim \mathcal{N}\left(\sum_k x_{i,k} \mu_{k,j}, \sum_k x_{i,k}^2 \sigma_{k,j}^2\right).$$

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

Example

Batch size = 64, matrix W dim is 64×64 .

- Graves: one sample, $64 \times 64 = 4096$ elements. Poor approximation.
- Naive solution: sample parameters 64 times, $64 \times 64 \times 64 = 262144$ elements. Better approximation (in theory).
- Local reparametrization: sample G , $64 \times 64 = 4096$ elements. Better approximation.



Expectation propagation

Minka, 2001: represent prior and approximation distribution via multiplication of factors:

$$p(w|\mathcal{D}) = \prod_i f_i, \quad q(w) = \prod_i \tilde{f}_i.$$

Main idea — minimize $KL(p(w|\mathcal{D})||q(w))$.

- Select factor \tilde{f}_i to approximate, «removing» it from consideration, changing into real factor value:

$$q^i \propto f_i \prod_{j \neq i} \tilde{f}_j$$

- Set moments of q^i equal to moments to distribution to approximate (correct, if q is from exponential distribution)
- Repeat until convergence

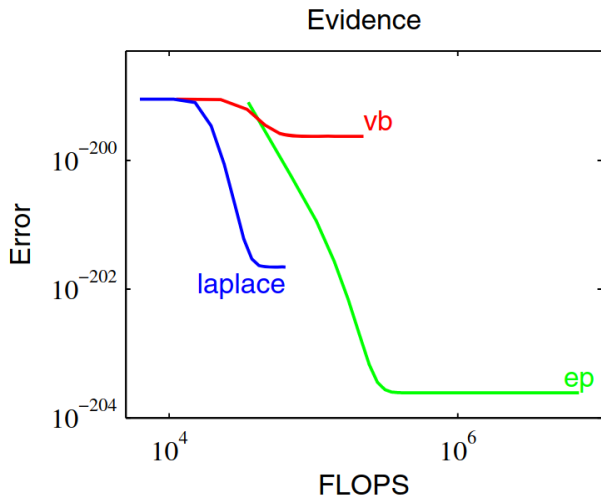
Expectation propagation: pros and cons

Cons:

- Assumption about posterior distribution (rather slight)
- Original version works only for q from exponential distribution
- No convergence guarantee

Pros:

- Minimizes KL, not it's lower bound



Plot for the 2-component Gaussian mixture.

Probabilistic backpropagation

Combination of Expectation propagation and backpropagation.

Backward pass:

Update parameters using Bayes rule:

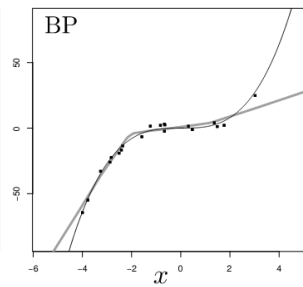
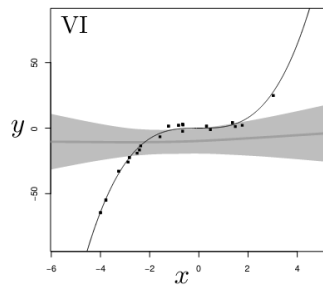
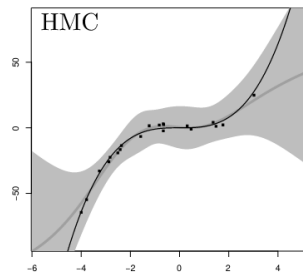
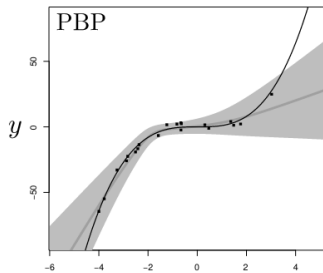
$$p(w_i|\mathcal{D}) = Z^{-1}p(\mathcal{D}|w_i, w^i)p(w) \rightarrow p(w_i|\mathcal{D}) = Z^{-1}p(\mathcal{D}|w_i, w^i)\mathcal{N}(w|\mu, \sigma^2).$$

Problem is to calculate Z .

Forward pass:

Compute Z approximately with ab assumption $f(x, w) \sim \mathcal{N}(m, v)$.

For m, v with ReLU activation there exists an iterative algorithm.



Reference

- Bishop C. M. Pattern recognition //Machine learning. – 2006. – Т. 128. – №. 9.
- MacKay D. J. C., Mac Kay D. J. C. Information theory, inference and learning algorithms. – Cambridge university press, 2003.
- Бахтеев О. Ю., Стрижов В. В. Выбор моделей глубокого обучения субоптимальной сложности //Автоматика и телемеханика. – 2018. – №. 8. – С. 129-147.
- Graves A. Practical variational inference for neural networks //Advances in neural information processing systems. – 2011. – Т. 24.
- Louizos C., Ullrich K., Welling M. Bayesian compression for deep learning //arXiv preprint arXiv:1705.08665. – 2017.
- Kingma D. P., Salimans T., Welling M. Variational dropout and the local reparameterization trick //Advances in neural information processing systems. – 2015. – Т. 28. – С. 2575-2583.
- Higgins I. et al. beta-vae: Learning basic visual concepts with a constrained variational framework. – 2016.
- Alemi A. et al. Fixing a broken ELBO //International Conference on Machine Learning. – PMLR, 2018. – С. 159-168.
- Minka T. P. Expectation propagation for approximate Bayesian inference //arXiv preprint arXiv:1301.2294. – 2013.
- Hernández-Lobato J. M., Adams R. Probabilistic backpropagation for scalable learning of bayesian neural networks //International conference on machine learning. – PMLR, 2015. – С. 1861-1869.