

Bayesian multimodeling: variational inference

MIPT

2023

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:

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

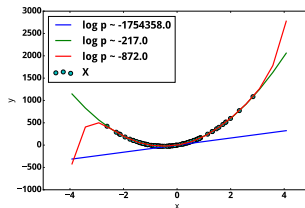
Second level: find optimal model:

Evidence:

$$p(\mathcal{D}|\mathbf{h}) = \int_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathbf{h})d\mathbf{w}.$$



Model selection scheme



Polynomial regression example

Local variational optimization, idea

Consider a problem of approximation of $f(x) = \exp(-x)$ by linear function.

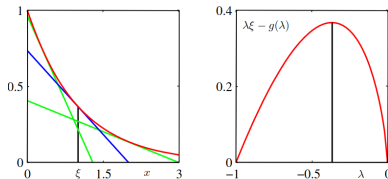
Any linear function will be a lower bound on $f(x)$ if it corresponds to a tangent. Use Taylor series:

$$y(x) = f(x_0) + f'(x_0)(x - x_0)$$

or, using $\lambda = -f(x_0)$

$$y(x) = \lambda x - \lambda + \lambda \log(-\lambda).$$

Where is the variational optimization here?



Local variational optimization, idea

Consider a problem of approximation of $f(x) = \exp(-x)$ by linear function.

Any linear function will be a lower bound on $f(x)$ if it corresponds to a tangent. Use Taylor series:

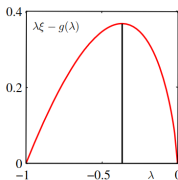
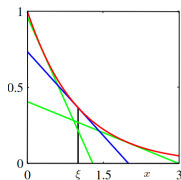
$$y(x) = f(x_0) + f'(x_0)(x - x_0)$$

or, using $\lambda = -f(x_0)$

$$y(x) = \lambda x - \lambda + \lambda \log(-\lambda).$$

We must find the tightest bound:

$$\max_{\lambda} \lambda x - \lambda + \lambda \log(-\lambda)$$



Local variational optimization and Evidence

Using similar approach we can approximate more interesting functions, for example sigmoid :

$$\log \sigma(x) = -\frac{x}{2} - \log(e^{\frac{x}{2}} + e^{\frac{-x}{2}}).$$

Note that $f(x) = -\log(e^{\frac{x}{2}} + e^{\frac{-x}{2}})$ is convex. Its approximation is:

$$\max_{x^2} \left(\lambda x^2 - f\left(\sqrt{x^2}\right) \right).$$

Optimal value gives:

$$\sigma(x) \geq \sigma(x_0) \exp \left((x - x_0) / 2 - \lambda(x_0) (x_0^2 - x^2) \right).$$

The evidence integral becomes quadratic \Rightarrow we can use an approximation by Gaussian, similar to Laplace approximation.

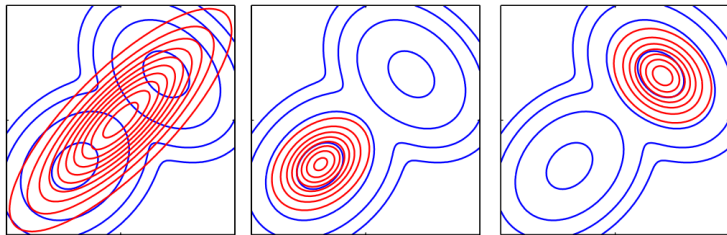
Evidence lower bound, ELBO

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

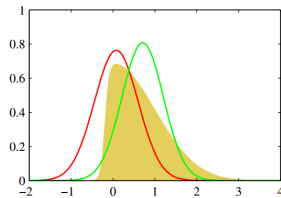
Evidence lower bound estimation often reduces to optimization problem

$$\log p(\mathcal{D}|\mathbf{h}) \geq$$

$$\geq \text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathcal{D})) = - \int_{\mathbf{w}} q(\mathbf{w}) \log \frac{p(\mathbf{w}|\mathcal{D})}{q(\mathbf{w})} d\mathbf{w} = \mathbb{E}_{\mathbf{w}} \log p(\mathcal{D}|\mathbf{w}) - \text{KL}(q(\mathbf{w})||p(\mathbf{w}|\mathbf{h})).$$



Variational inference vs. expectation propagation (Bishop)



Laplace Approximation vs

Variational inference

Minimum description length principle

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

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

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

\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 posterior 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

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

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

See also, [β -VAE, Fixing Broken ELBO, Talk from Olga Grebenkova, 2021].

What's better: ELBO or Laplace approximation?

(or MCMC)?

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:

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

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:

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

Graves, 2011

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

Variational distribution: $q(\mathbf{w}) \sim \mathcal{N}(\boldsymbol{\mu}_q, \sigma_q \mathbf{I})$.

Greedy hyperparameter optimization:

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

Parameter pruning w_i using relative PDF:

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



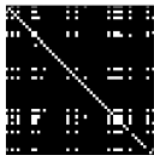
$\lambda = 0.01$



$\lambda = 0.05$



$\lambda = 0.1$



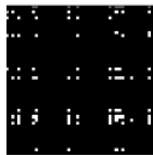
$\lambda = 0.2$



$\lambda = 0.5$



$\lambda = 1$



$\lambda = 2$

ELBO: normal distribution

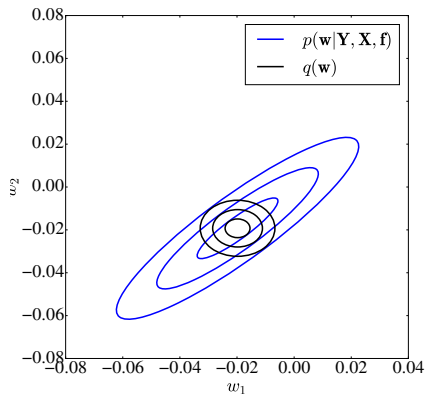
“Common” loss function:

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

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

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

Poor approximation example q



Local reparametrization

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

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

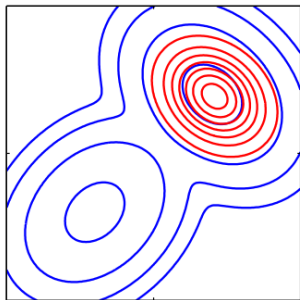
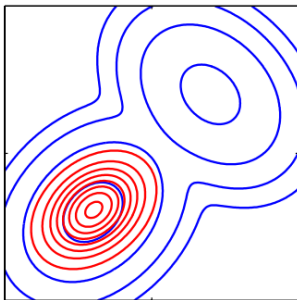
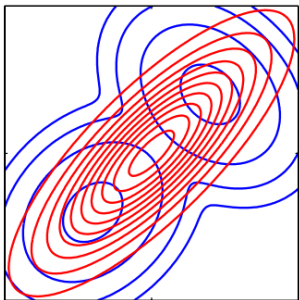
$$\mathbf{G} = \mathbf{XW}, \quad G_{i,j} \sim \mathcal{N}\left(\sum_k x_{i,k} \mu_{k,j}, \sum_k x_{i,k}^2 \sigma_{k,j}^2\right).$$

Instead of sampling parameters, sample elements from \mathbf{G} (units after activation).

Example

Batch size = 64, matrix \mathbf{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 \mathbf{G} , $64 \times 64 = 4096$ elements. Better approximation.



Expectation propagation

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

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

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

- 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 of distribution to approximate (correct, if q is from exponential distribution)
- Repeat until convergence

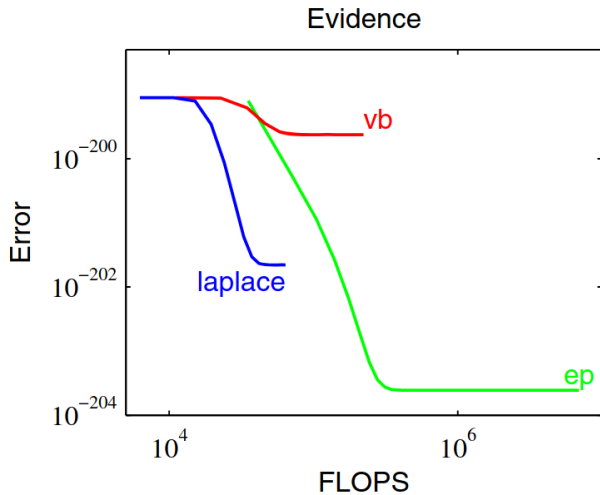
Expectation propagation: pros and cons

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.¹

Probabilistic model

Given data $\mathcal{D} = \{x_n, y_n\}_{n=1}^N$, made up of D -dimensional feature vectors and corresponding scalar target variables, we assume that $y_n = f(x_n; W) + \varepsilon_n$, where $f(; \mathbf{W})$ is the output of a multi-layer neural network with weights given by \mathbf{W} and $\varepsilon_n \sim \mathcal{N}(0, \gamma^{-1})$.

Prior distributions:

$$p(W|\lambda) = \prod_{\mathbf{w} \in \mathbf{W}} \mathcal{N}(w|0, \lambda^{-1}),$$

$$p(\lambda) = \Gamma(\lambda|\alpha_0^\lambda, \beta_0^\lambda),$$

$$p(\gamma) = \Gamma(\gamma|\alpha_0^\gamma, \beta_0^\gamma).$$

¹See talk of Polina Barabanshchikova, 2022

Theory

Likelihood for the weights \mathbf{W} and the noise precision γ is

$$p(\mathbf{y}|\mathbf{W}, \mathbf{X}, \gamma) = \prod_{n=1}^N \mathcal{N}(y_n|f(\mathbf{x}_n; \mathbf{W}), \gamma^{-1}).$$

The posterior distribution for \mathbf{W} , γ , λ

$$p(\mathbf{W}, \gamma, \lambda|\mathcal{D}) = \frac{p(\mathbf{y}|\mathbf{W}, \mathbf{X}, \gamma)p(\mathbf{W}|\lambda)p(\lambda)p(\gamma)}{p(\mathbf{y}|\mathbf{X})}.$$

Probabilistic backpropagation (PBP) approximates the exact posterior with a factored distribution given by

$$q(\mathbf{W}, \gamma, \lambda) = \prod_{\mathbf{w} \in \mathbf{W}} \mathcal{N}(\mathbf{w}|m_{\mathbf{w}}, v_{\mathbf{w}}) \times \Gamma(\gamma|\alpha^{\gamma}, \beta^{\gamma})\Gamma(\lambda|\alpha^{\lambda}, \beta^{\lambda}).$$

Method description

Stages of PBP

1. In the first phase, the input data is propagated forward through the network. PBP sequentially approximates the marginal posterior distributions of each weight with a collection of one-dimensional Gaussians that match their marginal means and variances. At the end of this phase, PBP computes the logarithm of the marginal probability of the target variable.

Method description

Stages of PBP

1. In the first phase, the input data is propagated forward through the network. PBP sequentially approximates the marginal posterior distributions of each weight with a collection of one-dimensional Gaussians that match their marginal means and variances. At the end of this phase, PBP computes the logarithm of the marginal probability of the target variable.
2. In the second phase, the gradients of this quantity with respect to the means and variances of the approximate Gaussian posterior are propagated back. These derivatives are used to update the means and variances of the posterior approximation.

Method description

Update rule

Let $f(w)$ encode an arbitrary likelihood function for the single weight w and let our current beliefs regarding the scalar w be captured by a distribution $q(w)$. After seeing the data, our beliefs about w are updated according to Bayes' rule:

$$s(w) = Z^{-1}f(w)q(w),$$

where Z is the normalization constant.

We approximate this posterior with a distribution q^{new} that has the same form as q . The parameters of q^{new} are chosen to minimize the KL divergence between s and q^{new} .

Method description

Update rule (Example)

Assume that $q(w) = \mathcal{N}(w|m, v)$. In this case, the parameters of the new Gaussian beliefs $q^{new}(w) = \mathcal{N}(w|m^{new}, v^{new})$ that minimize the KL divergence between s and q^{new} can be obtained by

$$m^{new} = m + v \frac{\partial \log Z}{\partial m},$$

$$v^{new} = v - v^2 \left[\left(\frac{\partial \log Z}{\partial m} \right)^2 - 2 \frac{\partial \log Z}{\partial v} \right].$$

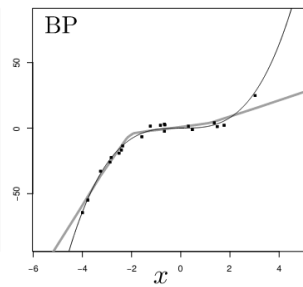
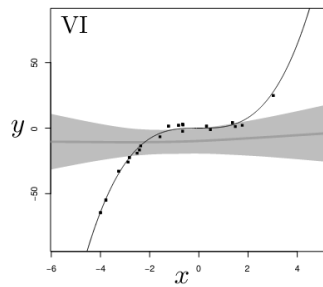
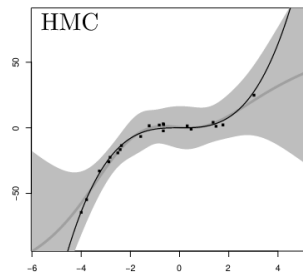
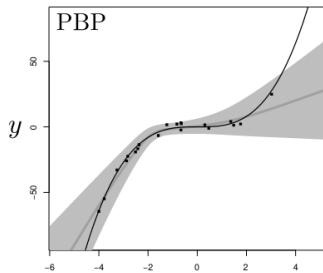
Method description

Update rule (Example)

Assume that $q(w) = \mathcal{N}(w|m, v)$. In this case, the parameters of the new Gaussian beliefs $q^{new}(w) = \mathcal{N}(w|m^{new}, v^{new})$ that minimize the KL divergence between s and q^{new} can be obtained by

$$m^{new} = m + v \frac{\partial \log Z}{\partial m},$$
$$v^{new} = v - v^2 \left[\left(\frac{\partial \log Z}{\partial m} \right)^2 - 2 \frac{\partial \log Z}{\partial v} \right].$$

Remark: Z is approximated during forward pass. Then its derivative is used to update the parameters of marginal distributions.



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.