

Ekaterina Lobacheva



Approximate Bayesian inference

2021



Yandex



EPFL



Slides are partially based on lectures of Dmitry Vetrov, Dmitry Kropotov and Kirill Struminsky, deepbayes.ru/2018

Bayesian ML models

Training stage:

$$p(\theta \mid X_{tr}, Y_{tr}) = \frac{p(Y_{tr} \mid X_{tr}, \theta) p(\theta)}{\int p(Y_{tr} \mid X_{tr}, \theta) p(\theta) d\theta}$$

Testing stage:

May be intractable

$$p(y \mid x, X_{tr}, Y_{tr}) = \int p(y \mid x, \theta) p(\theta \mid X_{tr}, Y_{tr}) d\theta$$

What can we do if they are intractable?

Approximate inference

Probabilistic model: $p(x, \theta) = p(x | \theta)p(\theta)$

Variational Inference

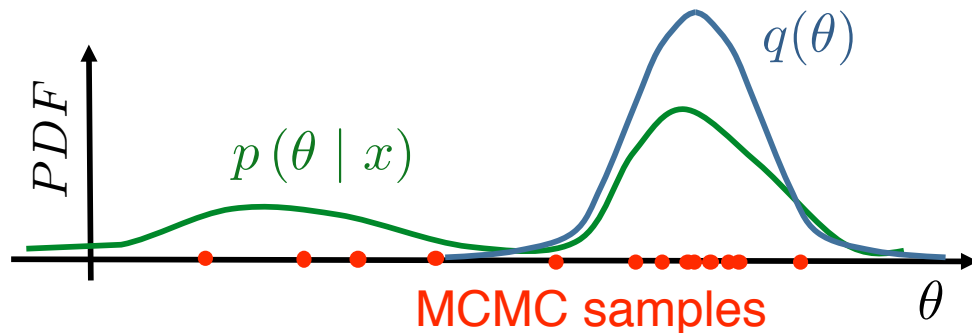
Approximate $p(\theta | x) \approx q(\theta) \in \mathcal{Q}$

- Biased
- Faster and more scalable

MCMC

Samples from unnormalized $p(\theta | x)$

- Unbiased
- Need a lot of samples



Variational inference

Probabilistic model: $p(x, \theta) = p(x \mid \theta)p(\theta)$

Main idea: find posterior approximation $p(\theta \mid x) \approx q(\theta) \in \mathcal{Q}$, using the following criterion function:

$$F(q) := KL(q(\theta) \parallel p(\theta \mid x)) \rightarrow \min_{q(\theta) \in \mathcal{Q}}$$



Kullback-Leibler divergence

a good mismatch measure between
two distributions over the **same domain**

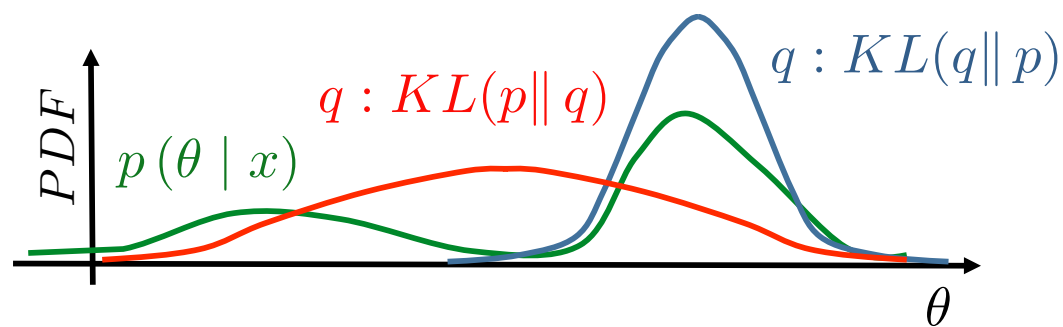
Kullback-Leibler divergence

A good mismatch measure between two distributions over the **same domain**

$$KL(q(\theta) \| p(\theta | x)) = \int q(\theta) \log \frac{q(\theta)}{p(\theta | x)} d\theta$$

Properties:

- $KL(q \| p) \geq 0$
- $KL(q \| p) = 0 \Leftrightarrow q = p$
- $KL(q \| p) \neq KL(p \| q)$



Variational inference

Probabilistic model: $p(x, \theta) = p(x \mid \theta)p(\theta)$

Main idea: find posterior approximation $p(\theta \mid x) \approx q(\theta) \in \mathcal{Q}$, using the following criterion function:

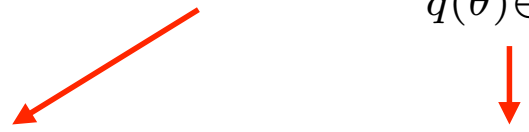
$$F(q) := KL(q(\theta) \parallel p(\theta \mid x)) \rightarrow \min_{q(\theta) \in \mathcal{Q}}$$

Variational inference

Probabilistic model: $p(x, \theta) = p(x \mid \theta)p(\theta)$

Main idea: find posterior approximation $p(\theta \mid x) \approx q(\theta) \in \mathcal{Q}$, using the following criterion function:

$$F(q) := KL(q(\theta) \parallel p(\theta \mid x)) \rightarrow \min_{q(\theta) \in \mathcal{Q}}$$



We could not compute the posterior in the first place

How to perform an optimization w.r.t. a distribution?

Mathematical magic

$$\begin{aligned}\log p(x) &= \int q(\theta) \log p(x) d\theta = \int q(\theta) \log \frac{p(x, \theta)}{p(\theta | x)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta) q(\theta)}{p(\theta | x) q(\theta)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta + \int q(\theta) \log \frac{q(\theta)}{p(\theta | x)} d\theta =\end{aligned}$$

Mathematical magic

$$\begin{aligned}\log p(x) &= \int q(\theta) \log p(x) d\theta = \int q(\theta) \log \frac{p(x, \theta)}{p(\theta | x)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta) q(\theta)}{p(\theta | x) q(\theta)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta + \int q(\theta) \log \frac{q(\theta)}{p(\theta | x)} d\theta = \\ &= \boxed{\mathcal{L}(q(\theta))} + \boxed{KL(q(\theta) \| p(\theta | x))}\end{aligned}$$

Evidence lower bound (ELBO)

KL-divergence we need for VI

ELBO = Evidence Lower Bound

$$\log p(x) = \mathcal{L}(q(\theta)) + KL(q(\theta) \parallel p(\theta \mid x))$$

Evidence:

$$p(\theta \mid x) = \frac{p(x \mid \theta)p(\theta)}{p(x)} = \frac{p(x \mid \theta)p(\theta)}{\int p(x \mid \theta)p(\theta)d\theta} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

Evidence of the probabilistic model shows the total probability of observing the data.

Lower Bound: KL is non-negative $\rightarrow \log p(x) \geq \mathcal{L}(q(\theta))$

Variational inference

Optimization problem with intractable posterior distribution:

$$F(q) := KL(q(\theta) \| p(\theta | x)) \rightarrow \min_{q(\theta) \in \mathcal{Q}}$$

Variational inference

Optimization problem with intractable posterior distribution:

$$F(q) := KL(q(\theta) \| p(\theta | x)) \rightarrow \min_{q(\theta) \in \mathcal{Q}}$$

Let's use our magic:


$$\log p(x) = \mathcal{L}(q(\theta)) + KL(q(\theta) \| p(\theta | x))$$

Variational inference

Optimization problem with intractable posterior distribution:

$$F(q) := KL(q(\theta) \| p(\theta | x)) \rightarrow \min_{q(\theta) \in \mathcal{Q}}$$

Let's use our magic:

$$\log p(x) = \mathcal{L}(q(\theta)) + KL(q(\theta) \| p(\theta | x))$$


The diagram illustrates the dependency of the terms in the equation on the variational distribution q . A red arrow points from the text "does not depend on q " to the $\log p(x)$ term. Two green arrows point from the text "depend on q " to the $\mathcal{L}(q(\theta))$ and $KL(q(\theta) \| p(\theta | x))$ terms respectively.

↑
does not depend on q

← ←
depend on q


Variational inference

Optimization problem with intractable posterior distribution:

$$F(q) := KL(q(\theta) \| p(\theta | x)) \rightarrow \min_{q(\theta) \in \mathcal{Q}}$$

Let's use our magic:

$$\log p(x) = \mathcal{L}(q(\theta)) + KL(q(\theta) \| p(\theta | x))$$



does not depend on q depend on q

$$KL(q(\theta) \| p(\theta | x)) \rightarrow \min_{q(\theta) \in \mathcal{Q}} \Leftrightarrow \mathcal{L}(q(\theta)) \rightarrow \max_{q(\theta) \in \mathcal{Q}}$$

Variational inference

Final optimisation problem:

$$\mathcal{L}(q(\theta)) = \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta \rightarrow \max_{q(\theta) \in \mathcal{Q}}$$

Variational inference: ELBO interpretation

Final optimisation problem:

$$\begin{aligned}\mathcal{L}(q(\theta)) &= \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta = \int q(\theta) \log \frac{p(x | \theta)p(\theta)}{q(\theta)} d\theta = \\ &= \int q(\theta) \log p(x | \theta) d\theta + \int q(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta =\end{aligned}$$

Variational inference: ELBO interpretation

Final optimisation problem:

$$\begin{aligned}\mathcal{L}(q(\theta)) &= \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta = \int q(\theta) \log \frac{p(x | \theta)p(\theta)}{q(\theta)} d\theta = \\ &= \int q(\theta) \log p(x | \theta) d\theta + \int q(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta = \\ &= \underbrace{\mathbb{E}_{q(\theta)} \log p(x | \theta)}_{\text{data term}} - \underbrace{KL(q(\theta) \| p(\theta))}_{\text{regularizer}}\end{aligned}$$

Variational inference: ELBO interpretation

Final optimisation problem:

$$\begin{aligned}\mathcal{L}(q(\theta)) &= \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta = \int q(\theta) \log \frac{p(x | \theta) p(\theta)}{q(\theta)} d\theta = \\ &= \int q(\theta) \log p(x | \theta) d\theta + \int q(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta = \\ &= \underbrace{\mathbb{E}_{q(\theta)} \log p(x | \theta)}_{\text{data term}} - \underbrace{KL(q(\theta) \| p(\theta))}_{\text{regularizer}}\end{aligned}$$

this is not the same KL-divergence!

Variational inference

Final optimisation problem:

$$\mathcal{L}(q(\theta)) = \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta \rightarrow \max_{q(\theta) \in \mathcal{Q}}$$

How to perform an optimization w.r.t. a distribution?

Variational inference

Final optimisation problem:

$$\mathcal{L}(q(\theta)) = \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta \rightarrow \max_{q(\theta) \in \mathcal{Q}}$$

How to perform an optimization w.r.t. a distribution?

Parametric approximation

Parametric family

$$q(\theta) = q(\theta \mid \lambda)$$

Parametric approximation

Parametric family of variational distributions:

$$q(\theta) = q(\theta \mid \lambda), \quad \lambda \text{ — some parameters}$$

Why is it a restriction? We choose a family of some fixed form:

- It may be too simple and insufficient to model the data
- If it is complex enough then there is no guaranty we can train it well to fit the data

Parametric approximation

Parametric family of variational distributions:

$$q(\theta) = q(\theta \mid \lambda), \quad \lambda \text{ — some parameters}$$

Variational inference transforms to parametric optimization problem:

$$\mathcal{L}(q(\theta \mid \lambda)) = \int q(\theta \mid \lambda) \log \frac{p(x, \theta)}{q(\theta \mid \lambda)} d\theta \rightarrow \max_{\lambda}$$

If we're able to calculate derivatives of ELBO w.r.t. λ then we can solve this problem using some numerical optimization solver.

Inference methods: summary

Probabilistic model: $p(x, \theta)$

We want to compute: $p(\theta \mid x)$

Approximation		Inference
Exact	$p(\theta \mid x)$	Full Bayesian inference
Parametric	$p(\theta \mid x) \approx q(\theta) = q(\theta \mid \lambda)$	Parametric VI
Delta function	$p(\theta \mid x) \approx \delta(\theta - \theta_{MP})$	MP inference
No prior	θ_{ML}	MLE