

# Introduction to Bayesian methods

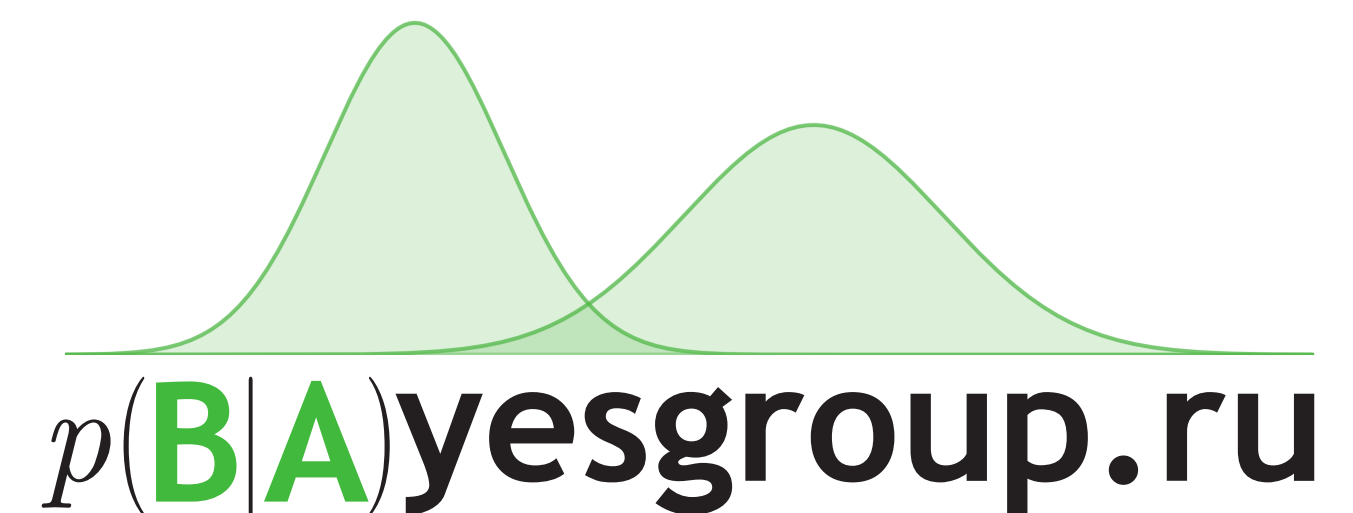
Dmitry P. Vetrov

Higher School of Economics, Samsung AI Center  
Moscow, Russia



NATIONAL RESEARCH  
UNIVERSITY

**SAMSUNG**  
**Research**



Special thanks to Ekaterina Lobacheva for assistance with preparation of slides

# Outline

- Bayesian framework
- Bayesian ML models
- Full Bayesian inference and conjugate distributions
- Approximate Bayesian inference

# How to work with distributions?

$$\text{Conditional} = \frac{\text{Joint}}{\text{Marginal}}, \quad p(x|y) = \frac{p(x, y)}{p(y)}$$

## Product rule

any joint distribution can be expressed as a product of one-dimensional conditional distributions

$$p(x, y, z) = p(x|y, z)p(y|z)p(z)$$

## Sum rule

any marginal distribution can be obtained from the joint distribution by integrating out

$$p(y) = \int p(x, y)dx$$

# Example

- We have a joint distribution over three groups of variables  $p(x, y, z)$
- We observe  $x$  and are interested in predicting  $y$
- Values of  $z$  are unknown and irrelevant to us
- How to estimate  $p(y|x)$  from  $p(x, y, z)$ ?

# Example

- We have a joint distribution over three groups of variables  $p(x, y, z)$
- We observe  $x$  and are interested in predicting  $y$
- Values of  $z$  are unknown and irrelevant to us
- How to estimate  $p(y|x)$  from  $p(x, y, z)$ ?

$$p(y|x) = \frac{p(x, y)}{p(x)} = \frac{\int p(x, y, z) dz}{\int p(x, y, z) dz dy}$$

Sum rule and product rule allow to obtain arbitrary conditional distributions from the joint one

# Bayes theorem

Bayes theorem (follows from product and sum rules):

$$p(y|x) = \frac{p(x, y)}{p(x)} = \frac{p(x|y)p(y)}{p(x)} = \frac{p(x|y)p(y)}{\int p(x|y)p(y)dy}$$

Bayes theorem defines the rule for uncertainty conversion when new information arrives:

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

# Statistical inference

**Problem:** given i.i.d. data  $X = (x_1, \dots, x_n)$  from distribution  $p(x|\theta)$  one needs to estimate  $\theta$

**Frequentist framework:** use maximum likelihood estimation (MLE)

$$\theta_{ML} = \arg \max p(X|\theta) = \arg \max \prod_{i=1}^n p(x_i|\theta) = \arg \max \sum_{i=1}^n \log p(x_i|\theta)$$

**Bayesian framework:** encode uncertainty about  $\theta$  in a prior  $p(\theta)$  and apply Bayesian inference

$$p(\theta|X) = \frac{\prod_{i=1}^n p(x_i|\theta) p(\theta)}{\int \prod_{i=1}^n p(x_i|\theta) p(\theta) d\theta}$$

# Example: coin tossing

- We have a coin which may be fair or not
- The task is to estimate a probability  $\theta$  of landing heads up
- Data: 2 tosses with a result (H,H)



Head (H)



Tail (T)



# Example: coin tossing

- We have a coin which may be fair or not
- The task is to estimate a probability  $\theta$  of landing heads up
- Data: 2 tosses with a result (H,H)



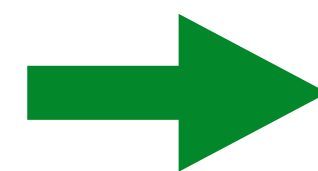
Head (H)

Tail (T)

## Frequentist framework:

In all experiments the coin  
landed heads up

$$\theta_{ML} = 1$$



The coin is not fair and  
always lands heads up

# Example: coin tossing

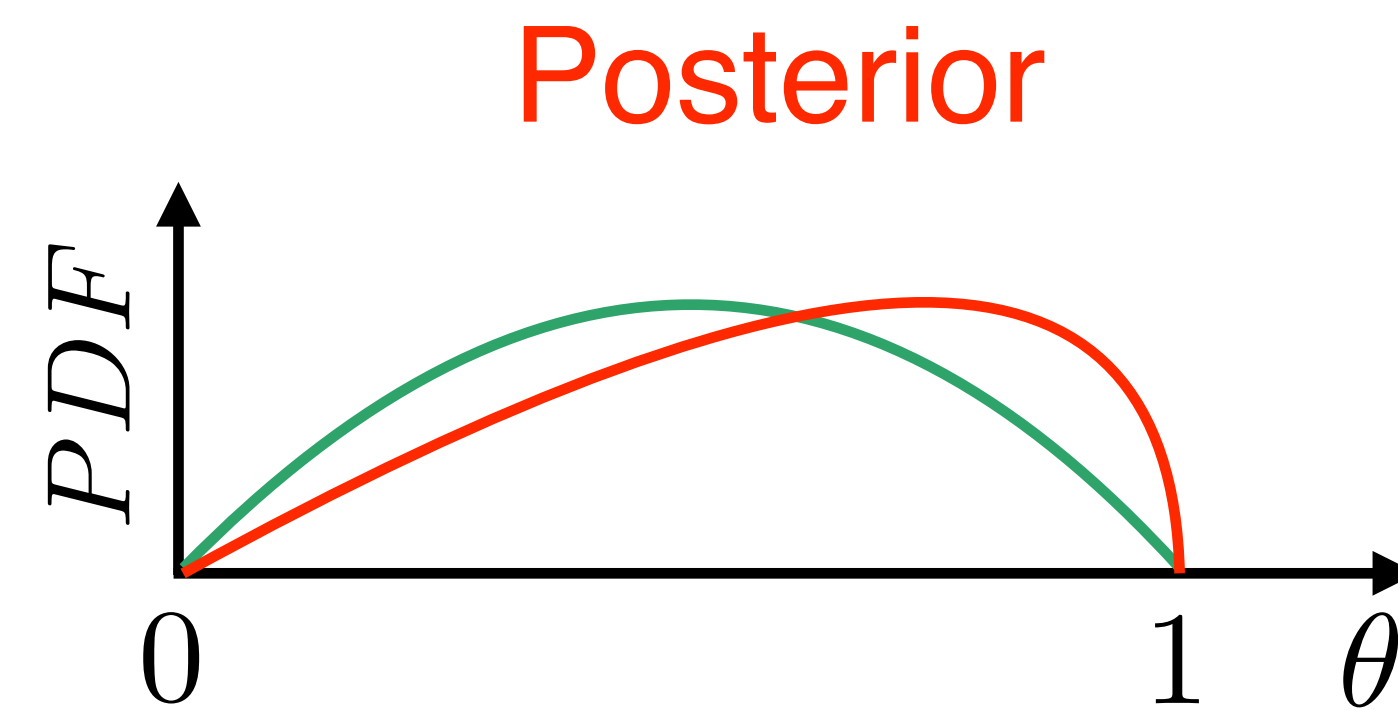
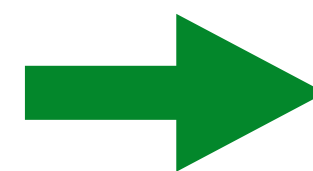
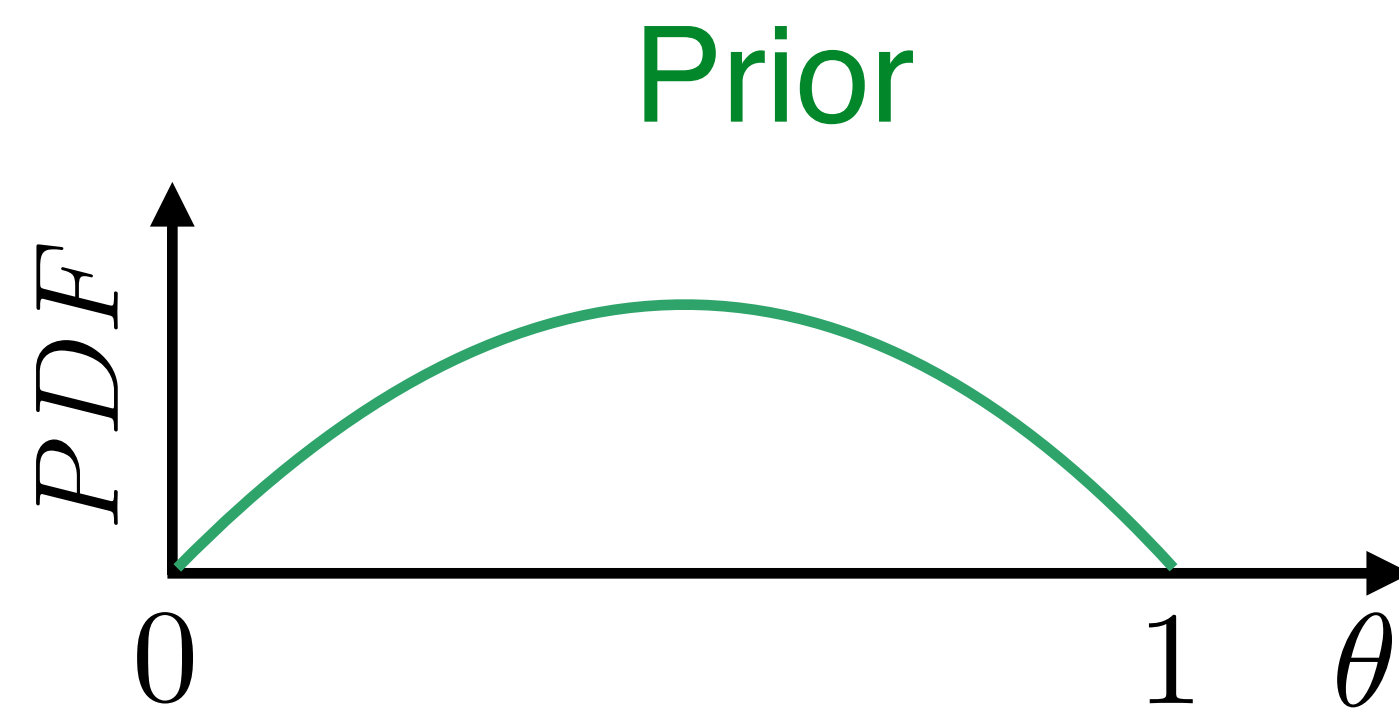
- We have a coin which may be fair or not
- The task is to estimate a probability  $\theta$  of landing heads up
- Data: 2 tosses with a result (H,H)



Head (H)

Tail (T)

## Bayesian framework:



# Example: coin tossing

- We have a coin which may be fair or not
- The task is to estimate a probability  $\theta$  of landing heads up
- Data: 1000 tosses with a result (H,H,T,...) — 489 tails and 511 heads



Head (H)



Tail (T)

# Example: coin tossing

- We have a coin which may be fair or not
- The task is to estimate a probability  $\theta$  of landing heads up
- Data: 1000 tosses with a result (H,H,T,...) — 489 tails and 511 heads

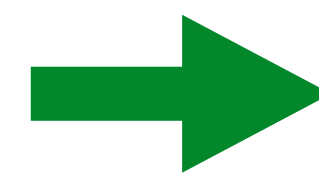


Head (H)

Tail (T)

## Both frameworks:

Sufficient amount of data  
matches our expectations



The coin is fair

# Frequentist vs. Bayesian frameworks

	Frequentist	Bayesian
Variables	random and deterministic	everything is random
Applicability	$n \gg d$	$\forall n$

- The number of tunable parameters in modern ML models is comparable with the sizes of training data
- Frequentist framework is a limit case of Bayesian one:

$$\lim_{n/d \rightarrow \infty} p(\theta | x_1, \dots, x_n) = \delta(\theta - \theta_{ML})$$

# Advantages of Bayesian framework

- We can encode our prior knowledge or desired properties of the final solution into a prior distribution
- Prior is a form of regularization
- Additionally to the point estimate of  $\theta$  posterior contains information about the uncertainty of the estimate

Bayesian framework just provides an alternative point of view, it DOES NOT contradict or deny frequentist framework

# Probabilistic ML model

For each object in the data:

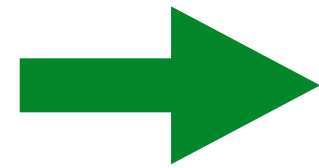
- $x$  — set of observed variables (features)
- $y$  — set of hidden / latent variables (class label / hidden representation, etc.)

Model:

- $\theta$  — model parameters (e.g. weights of the linear model)

# Discriminative probabilistic ML model

Models  $p(y, \theta | x)$



Cannot generate new objects —  
needs  $x$  as an input

Usually assumes that prior over  $\theta$  does not depend on  $x$ :

$$p(y, \theta | x) = p(y | x, \theta)p(\theta)$$

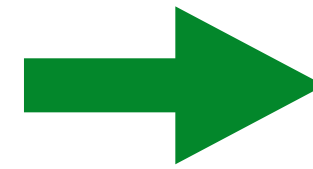
Examples:

- Classification or regression task (hidden space is much easier than the observed one)
- Machine translation (hidden and observed spaces have the same complexity)



# Generative probabilistic ML model

Models joint distribution  
 $p(x, y, \theta) = p(x, y | \theta)p(\theta)$



Can generate new objects,  
i.e. pairs  $(x, y)$

May be quite difficult to train since the observed space is usually much more complicated than the hidden one

Examples:

- Generation of text, speech, images, etc.

# Training Bayesian ML models

We are given training data  $(X_{tr}, Y_{tr})$  and a discriminative model  $p(y, \theta | x)$

**Training stage** — Bayesian inference over  $\theta$ :

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

**Result:** ensemble of algorithms rather than a single one  $\theta_{ML}$

- Ensemble usually outperforms single best model
- Posterior capture all dependencies from the training data that the model could extract and may be used as a new prior later

# Predictions of Bayesian ML models

## Testing stage:

- From training we have a posterior distribution  $p(\theta | X_{tr}, Y_{tr})$
- New data point  $x$  arrives
- We need to compute the predictive distribution on its hidden value  $y$

Ensembling w.r.t. posterior over the parameters  $\theta$ :

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

# Bayesian ML models

**Training stage:**

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

**Testing stage:**

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

# Bayesian ML models

**Training stage:**

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

**Testing stage:**

May be intractable

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

When are the integrals tractable?  
What can we do if they are intractable?

# Conjugate distributions

Distribution  $p(\theta)$  and  $p(x | \theta)$  are conjugate iff  $p(\theta | x)$  belongs to the same parametric family as  $p(\theta)$ :

$$p(\theta) \in \mathcal{A}(\alpha), \quad p(x | \theta) \in \mathcal{B}(\theta) \quad \longrightarrow \quad p(\theta | x) \in \mathcal{A}(\alpha')$$

# Conjugate distributions

Distribution  $p(\theta)$  and  $p(x | \theta)$  are conjugate iff  $p(\theta | x)$  belongs to the same parametric family as  $p(\theta)$ :

$$p(\theta) \in \mathcal{A}(\alpha), \quad p(x | \theta) \in \mathcal{B}(\theta) \quad \longrightarrow \quad p(\theta | x) \in \mathcal{A}(\alpha')$$

**Intuition:**

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

# Conjugate distributions

Distribution  $p(\theta)$  and  $p(x | \theta)$  are conjugate iff  $p(\theta | x)$  belongs to the same parametric family as  $p(\theta)$ :

$$p(\theta) \in \mathcal{A}(\alpha), \quad p(x | \theta) \in \mathcal{B}(\theta) \quad \longrightarrow \quad p(\theta | x) \in \mathcal{A}(\alpha')$$

**Intuition:**

$$p(\theta | x) = \frac{\boxed{p(x | \theta)p(\theta)}}{\int p(x | \theta)p(\theta)d\theta} \longleftarrow \text{conjugate}$$

- Denominator is tractable since any distribution in  $\mathcal{A}$  is normalized



# Conjugate distributions

Distribution  $p(\theta)$  and  $p(x | \theta)$  are conjugate iff  $p(\theta | x)$  belongs to the same parametric family as  $p(\theta)$ :

$$p(\theta) \in \mathcal{A}(\alpha), \quad p(x | \theta) \in \mathcal{B}(\theta) \quad \longrightarrow \quad p(\theta | x) \in \mathcal{A}(\alpha')$$

## Intuition:

$$p(\theta | x) = \frac{p(x | \theta)p(\theta)}{\int p(x | \theta)p(\theta)d\theta} \propto p(x | \theta)p(\theta)$$

- Denominator is tractable since any distribution in  $\mathcal{A}$  is normalized
- All we need is to compute  $\alpha'$

# Full Bayesian inference

**Training stage:**

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

**Testing stage:**

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

Integrals are tractable if prior and likelihood are conjugate

# Full Bayesian inference

- Easy to use - analytical formulas for training and testing stages
- Strong assumptions on the model - conjugacy of prior and likelihood
  - Choose conjugate prior
  - Only simple models (not flexible enough for most of the cases)

# Example: coin tossing

- We have a coin which may be fair or not
- The task is to estimate a probability  $\theta$  of landing heads up
- Data:  $X = (x_1, \dots, x_n)$ ,  $x \in \{0, 1\}$



Head (H)



Tail (T)

## Probabilistic model:

$$p(x, \theta) = p(x | \theta)p(\theta)$$

# Example: coin tossing

- We have a coin which may be fair or not
- The task is to estimate a probability  $\theta$  of landing heads up
- Data:  $X = (x_1, \dots, x_n)$ ,  $x \in \{0, 1\}$



Head (H)



Tail (T)

## Probabilistic model:

$$p(x, \theta) = p(x | \theta)p(\theta)$$

**Likelihood:**  $Bern(x | \theta) = \theta^x (1 - \theta)^{1-x}$

# Example: coin tossing

- We have a coin which may be fair or not
- The task is to estimate a probability  $\theta$  of landing heads up
- Data:  $X = (x_1, \dots, x_n)$ ,  $x \in \{0, 1\}$



## Probabilistic model:

$$p(x, \theta) = p(x | \theta)p(\theta)$$

**Likelihood:**  $Bern(x | \theta) = \theta^x (1 - \theta)^{1-x}$

**Prior: ???**

# Example: coin tossing

How to choose a prior?

- Correct domain:  $\theta \in [0, 1]$
- Include prior knowledge: a coin is most likely fair
- Inference complexity: use conjugate prior

# Example: coin tossing

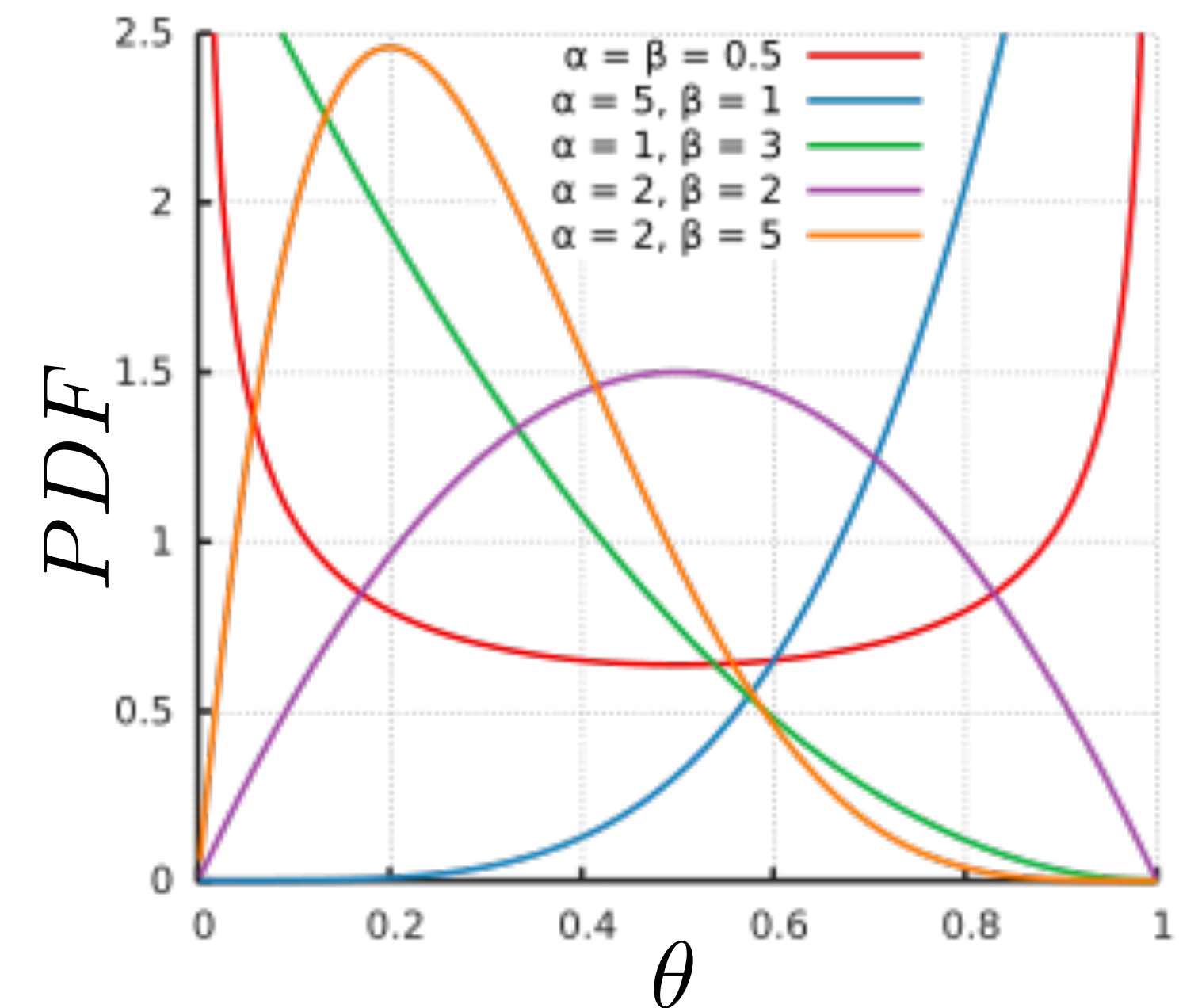
How to choose a prior?

- Correct domain:  $\theta \in [0, 1]$
- Include prior knowledge: a coin is most likely fair
- Inference complexity: use conjugate prior

Beta distribution matches all requirements:

$$Beta(\theta \mid a, b) = \frac{1}{B(a, b)} \theta^{a-1} (1 - \theta)^{b-1}$$

Beta distribution





# What to do if there is no conjugacy?

Simplest way — approximate posterior with delta function in  $\theta_{MP}$ :

$$\theta_{MP} = \arg \max_{\theta} p(\theta | X_{tr}, Y_{tr}) = \arg \max_{\theta} p(Y_{tr} | X_{tr}, \theta) p(\theta)$$

# What to do if there is no conjugacy?

Simplest way — approximate posterior with delta function in  $\theta_{MP}$ :

$$\theta_{MP} = \arg \max_{\theta} p(\theta | X_{tr}, Y_{tr}) = \arg \max_{\theta} p(Y_{tr} | X_{tr}, \theta) p(\theta)$$

On the testing stage:

$$p(y | x, X_{tr}, Y_{tr}) = \int p(y | x, \theta) p(\theta | X_{tr}, Y_{tr}) d\theta \approx p(y|x, \theta_{MP})$$

# What to do if there is no conjugacy?

Simplest way — approximate posterior with delta function in  $\theta_{MP}$ :

$$\theta_{MP} = \arg \max p(\theta | X_{tr}, Y_{tr}) = \arg \max p(Y_{tr} | X_{tr}, \theta) p(\theta)$$

We do not need to calculate  
the normalisation constant

On the testing stage:

$$p(y | x, X_{tr}, Y_{tr}) = \int p(y | x, \theta) p(\theta | X_{tr}, Y_{tr}) d\theta \approx p(y|x, \theta_{MP})$$

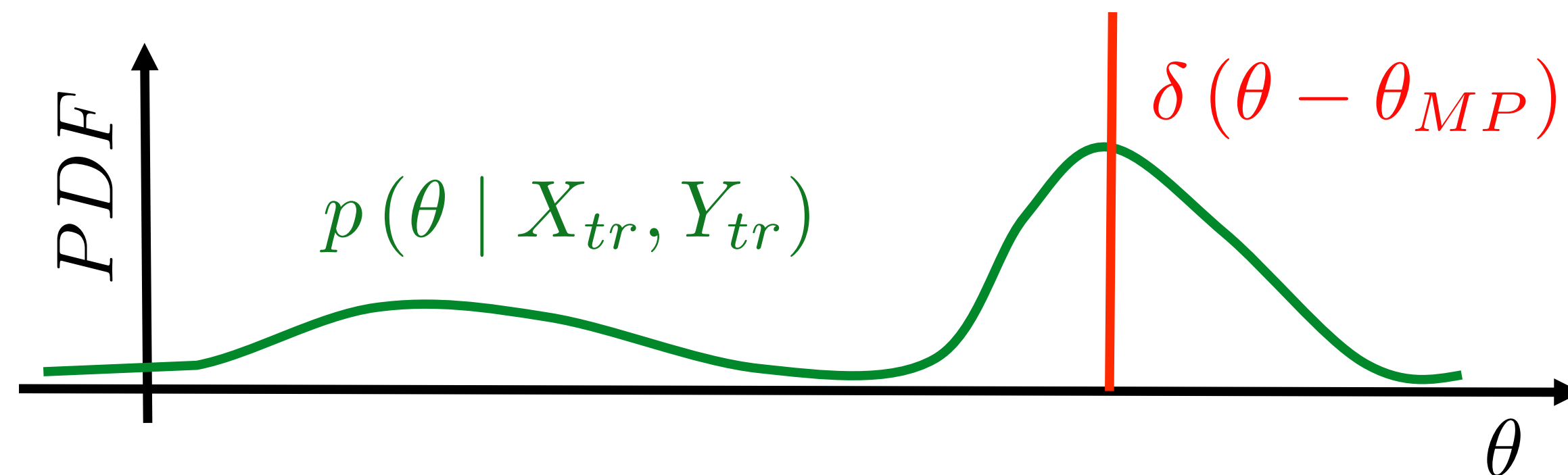
# What to do if there is no conjugacy?

Simplest way — approximate posterior with delta function in  $\theta_{MP}$ :

$$\theta_{MP} = \arg \max p(\theta | X_{tr}, Y_{tr}) = \arg \max p(Y_{tr} | X_{tr}, \theta) p(\theta)$$

On the testing stage:

$$p(y | x, X_{tr}, Y_{tr}) = \int p(y | x, \theta) p(\theta | X_{tr}, Y_{tr}) d\theta \approx p(y|x, \theta_{MP})$$



# What to do if there is no conjugacy?

Simplest way — approximate posterior with delta function in  $\theta_{MP}$ :

$$\theta_{MP} = \arg \max_{\theta} p(\theta | X_{tr}, Y_{tr}) = \arg \max_{\theta} p(Y_{tr} | X_{tr}, \theta) p(\theta)$$

On the testing stage:

$$p(y | x, X_{tr}, Y_{tr}) = \int p(y | x, \theta) p(\theta | X_{tr}, Y_{tr}) d\theta \approx p(y|x, \theta_{MP})$$

\* Not the same as  $\theta_{ML}$  — here we use prior

# What to do if there is no conjugacy?

Simplest way — approximate posterior with delta function in  $\theta_{MP}$ :

$$\theta_{MP} = \arg \max_{\theta} p(\theta | X_{tr}, Y_{tr}) = \arg \max_{\theta} p(Y_{tr} | X_{tr}, \theta) p(\theta)$$

On the testing stage:

$$p(y | x, X_{tr}, Y_{tr}) = \int p(y | x, \theta) p(\theta | X_{tr}, Y_{tr}) d\theta \approx p(y|x, \theta_{MP})$$

More advanced techniques are needed!

# Approximate inference

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

## Variational Inference

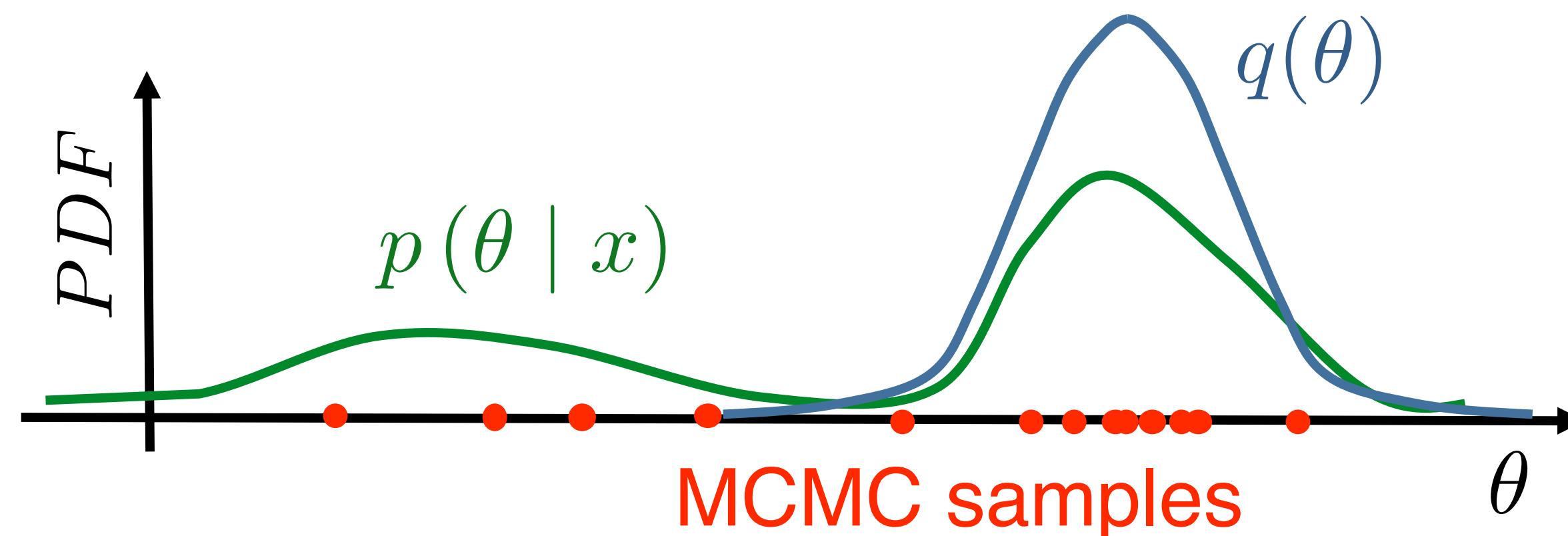
Approximate  $p(\theta | x) \approx q(\theta) \in \mathcal{Q}$

- Biased
- Faster and more scalable

## MCMC

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

- Unbiased
- Need a lot of samples



# Variational inference

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

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

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



**Kullback-Leibler divergence**

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



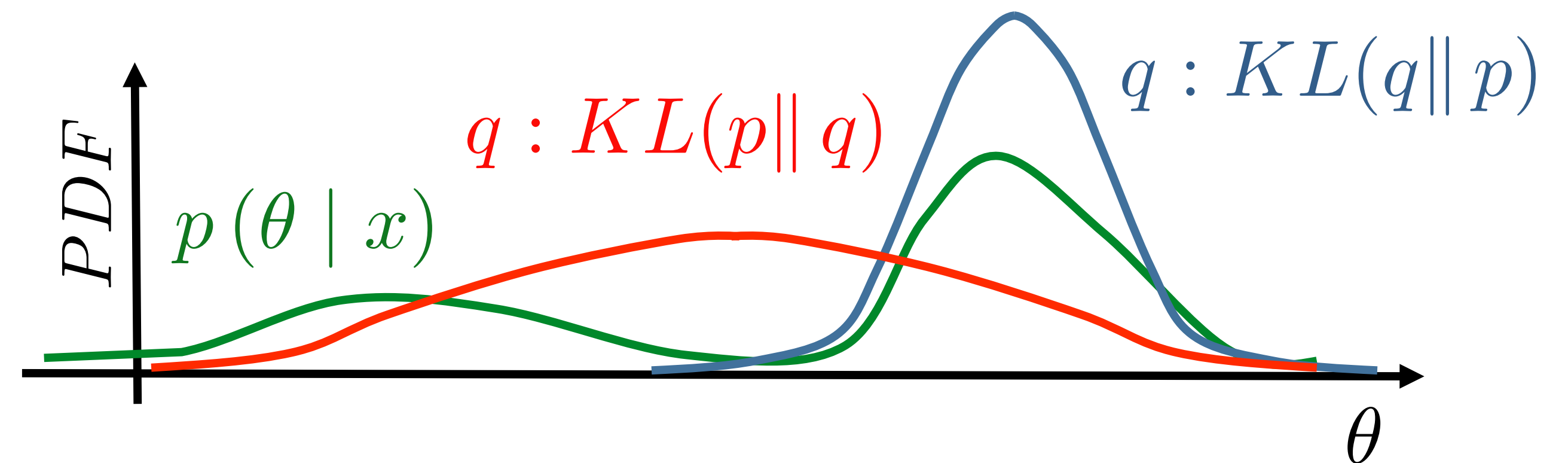
# Kullback-Leibler divergence

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

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

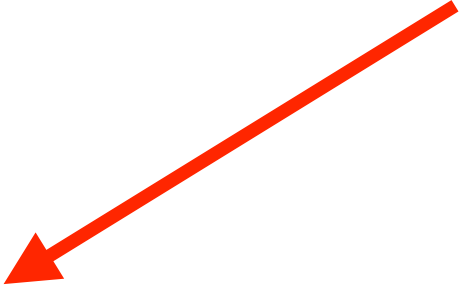
$$F(q) := KL(q(\theta) || p(\theta | x)) \rightarrow \min_{q(\theta) \in \mathcal{Q}}$$

# Variational inference

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

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

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



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

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

Evidence:

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

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:

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

# Variational inference

Optimization problem with intractable posterior distribution:

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

Let's use our magic:

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



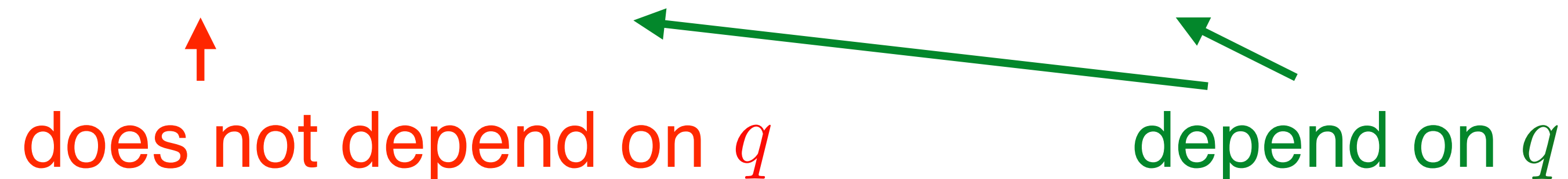
# Variational inference

Optimization problem with intractable posterior distribution:

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

Let's use our magic:

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



↑  
does not depend on  $q$

← ←  
depend on  $q$

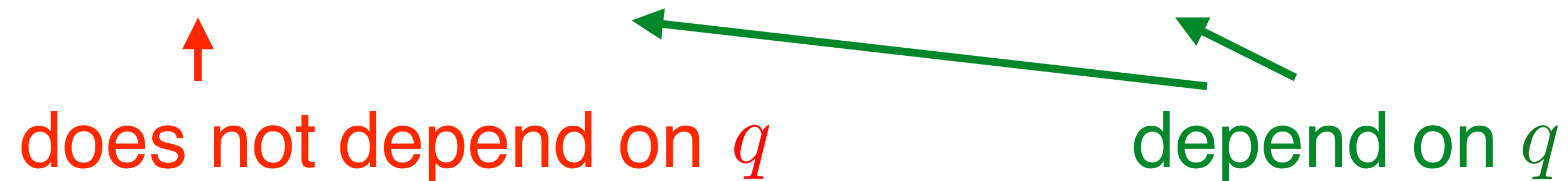
# Variational inference

Optimization problem with intractable posterior distribution:

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

Let's use our magic:

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



does not depend on  $q$                       depend on  $q$

$$KL(q(\theta) \parallel p(\theta \mid x)) \rightarrow \min_{q(\theta) \in \mathcal{Q}} \iff \mathcal{L}(q(\theta)) \rightarrow \max_{q(\theta) \in \mathcal{Q}}$$

# 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:

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

# 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

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

# 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.  $\lambda$  then we can solve this problem using some numerical optimization solver.

# Inference methods: summary

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

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

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