

# Variational Inference

Cédric Archambeau  
[cedrica@amazon.com](mailto:cedrica@amazon.com)

StatML.io, Oxford  
October 2021

# Outline

- ① Approximate Bayesian inference
- ② Variational inference
  - ▶ Mean field
  - ▶ Relation to Expectation-Maximisation
  - ▶ Structured variational inference
- ③ Stochastic and extensions
- ④ Case studies:
  - ▶ Gaussian variational approximations
  - ▶ Sparse Gaussian processes
  - ▶ Variational auto-encoders

# Bayesian statistics

$$\underbrace{p(\Theta|X)}_{\text{posterior}} = \frac{\overbrace{p(X|\Theta)}^{\text{likelihood}} \overbrace{p(\Theta)}^{\text{prior}}}{\underbrace{p(X)}_{\text{evidence}}},$$

$$p(X) = \int p(X, \Theta) d\Theta.$$

- The likelihood is the noise model.
- The prior encodes constraints (if any) on the parameters  $\Theta$ .
- Structure is added to the model through **latent variables**  $Z$ :  $p(X, Z|\Theta)$

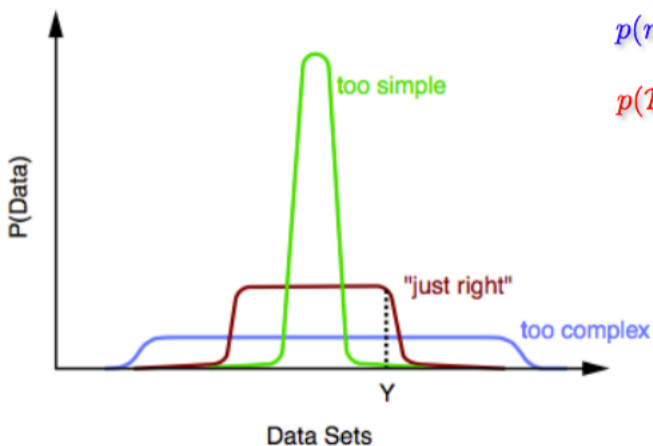
# Bayesian statistics

$$\underbrace{p(\Theta|X)}_{\text{posterior}} = \frac{\overbrace{p(X|\Theta)}^{\text{likelihood}} \overbrace{p(\Theta)}^{\text{prior}}}{\underbrace{p(X)}_{\text{evidence}}},$$

$$p(X) = \int p(X, \Theta) d\Theta.$$

- The likelihood is the noise model.
- The prior encodes constraints (if any) on the parameters  $\Theta$ .
- Structure is added to the model through **latent variables**  $Z$ :  $p(X, Z|\Theta)$
- Predictions are averaged over **all** possible models:  $p(x_*|X) = \int p(x_*|\Theta) p(\Theta|X) d\Theta$ .
- The goal is to maximise the marginal likelihood or **evidence**  $p(X|m)$ .

## What is great about the Bayesian paradigm?



$$p(m|\mathcal{D}) = \frac{p(\mathcal{D}|m)p(m)}{p(\mathcal{D})}$$

$$p(\mathcal{D}|m) = \sum_{\theta} p(\mathcal{D}|\theta, m)p(\theta|m)$$

## What is not so great with Bayesian paradigm?

Posterior inference:

$$p(\Theta|X, m) \propto \int p(X, Z, \Theta|m) dZ.$$

Model averaging:

$$p(x_*|X, m) = \int p(x_*|\Theta, m) p(\Theta|X, m) d\Theta.$$

Evidence maximisation:

$$p(X|m) = \int p(X, \Theta|m) d\Theta.$$

## Variational lower bound or evidence lower bound (ELBO)

$$\ln p(\mathbf{X}|m) \geq \ln p(\mathbf{X}|m) - \text{KL}(q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}) \| p(\mathbf{Z}, \boldsymbol{\Theta} | \mathbf{X}, m)) \triangleq -\mathcal{F}(\mathbf{w}).$$

## Variational lower bound or evidence lower bound (ELBO)

$$\ln p(\mathbf{X}|\mathbf{m}) \geq \ln p(\mathbf{X}|\mathbf{m}) - \text{KL}(q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}) \| p(\mathbf{Z}, \boldsymbol{\Theta}|\mathbf{X}, \mathbf{m})) \triangleq -\mathcal{F}(\mathbf{w}).$$

- The lower bound to the log marginal likelihood is obtained by applying **Jensen's** inequality:

$$\begin{aligned}\ln p(\mathbf{X}|\mathbf{m}) &= \ln \iint p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\Theta}|\mathbf{m}) d\mathbf{Z} d\boldsymbol{\Theta} \\ &= \ln \iint q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}) \frac{p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\Theta}|\mathbf{m})}{q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta})} d\mathbf{Z} d\boldsymbol{\Theta} \\ &\geq \iint q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}) \ln \frac{p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\Theta}|\mathbf{m})}{q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta})} d\mathbf{Z} d\boldsymbol{\Theta} \\ &= \ln p(\mathbf{X}|\mathbf{m}) + \iint q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}) \ln \frac{p(\mathbf{Z}, \boldsymbol{\Theta}|\mathbf{X}, \mathbf{m})}{q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta})} d\mathbf{Z} d\boldsymbol{\Theta}.\end{aligned}$$

## Variational lower bound or evidence lower bound (ELBO)

$$\ln p(\mathbf{X}|\mathbf{m}) \geq \ln p(\mathbf{X}|\mathbf{m}) - \text{KL}(q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}) \| p(\mathbf{Z}, \boldsymbol{\Theta}|\mathbf{X}, \mathbf{m})) \triangleq -\mathcal{F}(\mathbf{w}).$$

- The lower bound to the log marginal likelihood is obtained by applying **Jensen's** inequality:

$$\begin{aligned} \ln p(\mathbf{X}|\mathbf{m}) &= \ln \iint p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\Theta}|\mathbf{m}) d\mathbf{Z} d\boldsymbol{\Theta} \\ &= \ln \iint q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}) \frac{p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\Theta}|\mathbf{m})}{q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta})} d\mathbf{Z} d\boldsymbol{\Theta} \\ &\geq \iint q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}) \ln \frac{p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\Theta}|\mathbf{m})}{q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta})} d\mathbf{Z} d\boldsymbol{\Theta} \\ &= \ln p(\mathbf{X}|\mathbf{m}) + \iint q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}) \ln \frac{p(\mathbf{Z}, \boldsymbol{\Theta}|\mathbf{X}, \mathbf{m})}{q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta})} d\mathbf{Z} d\boldsymbol{\Theta}. \end{aligned}$$

- The analytically intractable integration problem is replaced by an **optimisation** problem!

## Other forms of the ELBO

$$-\mathcal{F}(\mathbf{w}) = \iint q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}) \ln \frac{p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\Theta} | m)}{q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta})} d\mathbf{Z} d\boldsymbol{\Theta}.$$

## Other forms of the ELBO

$$-\mathcal{F}(\mathbf{w}) = \iint q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}) \ln \frac{p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\Theta} | m)}{q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta})} d\mathbf{Z} d\boldsymbol{\Theta}.$$

- Free energy interpretation:

$$+\mathcal{F}(\mathbf{w}) = \underbrace{-\mathbb{E}(\ln p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\Theta} | m))}_{\text{energy}} - \underbrace{\mathbb{H}(q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta}))}_{\text{entropy}}. \quad (1)$$

## Other forms of the ELBO

$$-\mathcal{F}(w) = \iint q_w(Z, \Theta) \ln \frac{p(X, Z, \Theta | m)}{q_w(Z, \Theta)} dZ d\Theta.$$

- Free energy interpretation:

$$+\mathcal{F}(w) = \underbrace{-\mathbb{E}(\ln p(X, Z, \Theta | m))}_{\text{energy}} - \underbrace{H(q_w(Z, \Theta))}_{\text{entropy}}. \quad (1)$$

- Penalized model fit interpretation:

$$-\mathcal{F}(w) = \underbrace{\mathbb{E}(\ln p(X | Z, \Theta, m))}_{\text{model fit}} - \underbrace{\text{KL}(q_w(Z, \Theta) \| p(Z, \Theta | m))}_{\text{penalty}}. \quad (2)$$

# Definitions

The differential **entropy** measures the randomness of a random variable:

$$H(p) = - \int p(x) \ln p(x) dx.$$

The **Kullback-Leibler divergence** or relative entropy measures how two probability densities differ:

$$KL(q||p) = - \int q(x) \ln \frac{p(x)}{q(x)} dx \geq 0.$$

The KL is asymmetric (thus not a distance) and only zero if  $q(x) = p(x)$  for all  $x$ .

# Variational Inference

## Mean field variational inference [Bea03]

- A tractable solution is found by assuming  $q_w$  factorises given the data:

$$q_w(Z, \Theta) = \prod_n q(z_n; w_n) \times \prod_m q(\theta_m; w_m).$$

## Mean field variational inference [Bea03]

- A tractable solution is found by assuming  $q_w$  factorises given the data:

$$q_w(Z, \Theta) = \prod_n q(z_n; w_n) \times \prod_m q(\theta_m; w_m).$$

- The ELBO is given by

$$-\mathcal{F}(w) = \sum_n \mathbb{E}(\ln p(x_n | z_n, \Theta)) - \sum_n \text{KL}(q(z_n; w_n) \| p(z_n)) - \sum_m \text{KL}(q(\theta_m; w_m) \| p(\theta_m)).$$

- Variational inference (or variational Bayes or variational EM) alternates between the following two steps:

$$q(z_n; w_n) \propto e^{\mathbb{E}_{\neg z_n}(\ln p(x_n, z_n | \Theta))}, \quad q(\theta_m; w_m) \propto e^{\mathbb{E}_{\neg \theta_m}(\ln p(X, Z | \Theta))} p(\theta_m).$$

## Mean field variational inference [Bea03]

- A tractable solution is found by assuming  $q_w$  factorises given the data:

$$q_w(Z, \Theta) = \prod_n q(z_n; w_n) \times \prod_m q(\theta_m; w_m).$$

- The ELBO is given by

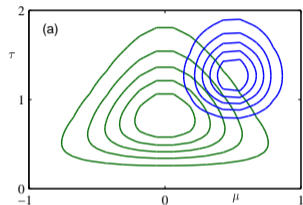
$$-\mathcal{F}(w) = \sum_n \mathbb{E}(\ln p(x_n | z_n, \Theta)) - \sum_n \text{KL}(q(z_n; w_n) \| p(z_n)) - \sum_m \text{KL}(q(\theta_m; w_m) \| p(\theta_m)).$$

- Variational inference (or variational Bayes or variational EM) alternates between the following two steps:

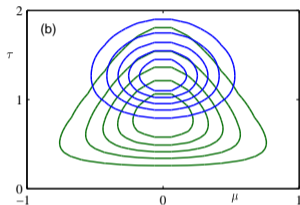
$$q(z_n; w_n) \propto e^{\mathbb{E}_{\neg z_n}(\ln p(x_n, z_n | \Theta))}, \quad q(\theta_m; w_m) \propto e^{\mathbb{E}_{\neg \theta_m}(\ln p(X, Z | \Theta))} p(\theta_m).$$

- The algorithm iteratively and **monotonically** maximises the ELBO, converging to a **local** maximum of the bound (not the evidence!)

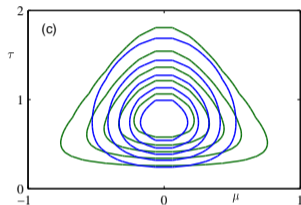
## Variational inference in action



Iteration 1



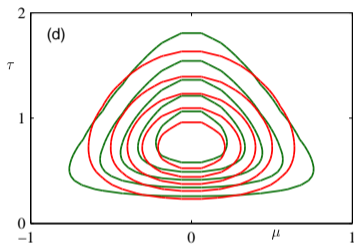
Iteration 2



Iteration 3

(Image credit [Bis06])

## What is lost?



Gaussian-Gamma

(Image credit [Bis06])

## How to make predictions?

- The predictive distribution is approximated by plugging in the approximate posterior  $q_w$ :

$$p(x_*|X) \approx \iint p(x_*|z_*, \Theta) q(z_*; w_*) q(\Theta; \{w_m\}_m) dz_* d\Theta.$$

- When analytically intractable, one can use Monte Carlo integration or heuristics based on statistics under the approximate posterior:

$$p(x_*|X) \approx p(x_*|\mathbb{E}(z_*), \mathbb{E}(\Theta)).$$

## Relation to expectation-maximisation (EM) [NH93]

$$\begin{aligned} -\mathcal{F}(\mathbf{w}) &= \ln p(\mathbf{X}|\Theta) - \text{KL}(q(Z) \| p(Z|\mathbf{X}, \Theta)), \\ -\mathcal{F}(\mathbf{w}) &= \mathbb{E}(\ln p(\mathbf{X}, Z|\Theta)) + H(q(Z)). \end{aligned}$$

## Relation to expectation-maximisation (EM) [NH93]

$$\begin{aligned}-\mathcal{F}(\mathbf{w}) &= \ln p(\mathbf{X}|\Theta) - \text{KL}(q(Z)||p(Z|\mathbf{X}, \Theta)), \\ -\mathcal{F}(\mathbf{w}) &= \mathbb{E}(\ln p(\mathbf{X}, Z|\Theta)) + H(q(Z)).\end{aligned}$$

- Expectation step:  $q(Z) \leftarrow p(Z|\mathbf{X}, \Theta^{old})$ .
- Maximisation step:  $\Theta^{new} = \arg \max_{\Theta} \mathbb{E}(\ln p(\mathbf{X}, Z|\Theta))$ .

## Relation to expectation-maximisation (EM) [NH93]

$$\begin{aligned} -\mathcal{F}(\mathbf{w}) &= \ln p(\mathbf{X}|\Theta) - \text{KL}(q(\mathbf{Z})\|p(\mathbf{Z}|\mathbf{X}, \Theta)), \\ -\mathcal{F}(\mathbf{w}) &= \mathbb{E}(\ln p(\mathbf{X}, \mathbf{Z}|\Theta)) + H(q(\mathbf{Z})). \end{aligned}$$

- Expectation step:  $q(\mathbf{Z}) \leftarrow p(\mathbf{Z}|\mathbf{X}, \Theta^{old})$ .
- Maximisation step:  $\Theta^{new} = \arg \max_{\Theta} \mathbb{E}(\ln p(\mathbf{X}, \mathbf{Z}|\Theta))$ .
- EM guarantees monotonic increase of the bound by construction.
- EM converges to local optimum of the log likelihood [Wu83].
- Approximate EM if  $q$  approximates the posterior [HZW03].

# EM in pictures

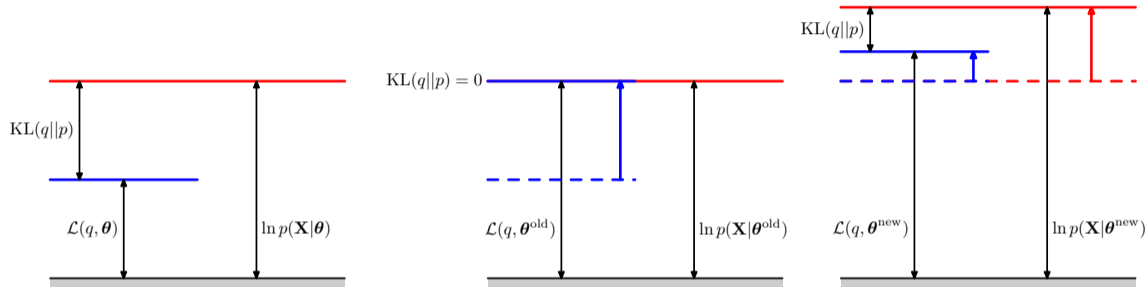


Image credit: [Bis06].

## Structured variational inference [SJ95, Wie00]

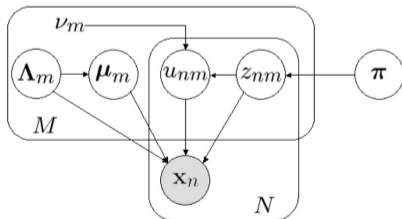
$$\arg \min_w \text{KL}(q_w(Z, \Theta) \| p(Z, \Theta | X, m))$$

- Mean field considers a fully factorised form to find a tractable solution.
- Structured variational inference avoids factorising when possible or imposes an approximate posterior of a predefined specific form.

## Example: mixture of Student- $t$ distributions [AV07]

$$p(\mathbf{x}|\Theta) = \sum_m \pi_m \text{Student}(\mathbf{x}|\boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m, \nu_m),$$

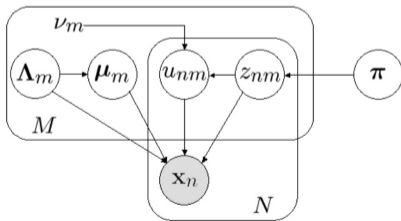
$$\text{Student}(\mathbf{x}|\boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m, \nu_m) = \int_{-\infty}^{+\infty} \text{Gaussian}(\mathbf{x}|\boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m \mathbf{u}_m) \text{Gamma}(\mathbf{u}_m|\nu_m/2, \nu_m/2) d\mathbf{u}_m.$$



## Example: mixture of Student- $t$ distributions [AV07]

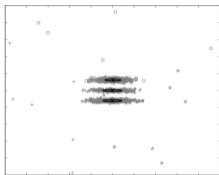
$$p(\mathbf{x}|\Theta) = \sum_m \pi_m \text{Student}(\mathbf{x}|\boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m, \nu_m),$$

$$\text{Student}(\mathbf{x}|\boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m, \nu_m) = \int_{-\infty}^{+\infty} \text{Gaussian}(\mathbf{x}|\boldsymbol{\mu}_m, \boldsymbol{\Lambda}_m \mathbf{u}_m) \text{Gamma}(\mathbf{u}_m|\nu_m/2, \nu_m/2) d\mathbf{u}_m.$$

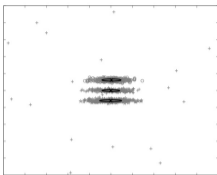


$$q(u_n, z_n) = \prod_m q(u_{nm})q(z_{nm}) \quad (\text{SMM1})$$

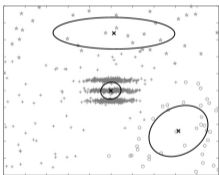
$$q(u_n, z_n) = \prod_m q(u_{nm}, z_{nm}) \quad (\text{SMM2})$$



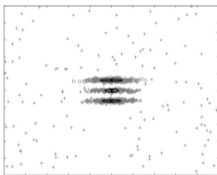
(a) Type-1 SMM, 2% of outliers.



(b) Type-2 SMM, 2% of outliers.

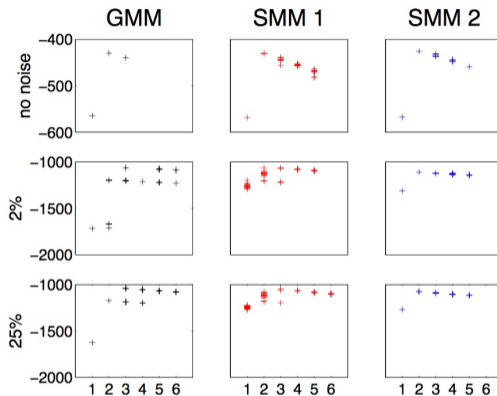


(c) Type-1 SMM, 15% of outliers.



(d) Type-2 SMM, 15% of outliers.

Robustness against outliers.



Old Faithful Geyser data – model selection with ELBO.

# Stochastic Variational Inference and Other Variants

## Mean field variational inference (MVI)

$$-\mathcal{F}(\mathbf{w}) = \sum_n \underbrace{\mathbb{E}(\ln p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\Theta}))}_{=\ell_n(\mathbf{w})} - \sum_n \text{KL}(q(\mathbf{z}_n; \mathbf{w}_n) \| p(\mathbf{z}_n)) - \sum_m \text{KL}(q(\boldsymbol{\theta}_m; \mathbf{w}_m) \| p(\boldsymbol{\theta}_m)) .$$

## Mean field variational inference (MVI)

$$-\mathcal{F}(w) = \sum_n \underbrace{\mathbb{E}(\ln p(x_n | z_n, \Theta))}_{=\ell_n(w)} - \sum_n \text{KL}(q(z_n; w_n) \| p(z_n)) - \sum_m \text{KL}(q(\theta_m; w_m) \| p(\theta_m)).$$

MVI can be rewritten as **batch gradient ascent**:

$$w_n \leftarrow \arg \max_{w_n} \ell_n(w) - \text{KL}(q(z_n; w_n) \| p(z_n)), \quad (\text{VE} - \text{step})$$

$$w_m \leftarrow \arg \max_{w_m} \sum_n \ell_n(w) - \text{KL}(q(\theta_m; w_m) \| p(\theta_m)). \quad (\text{VM} - \text{step})$$

- Monotonic increase of the bound; converges to local maximum of ELBO
- Priors are conjugate to the likelihood; updates are similar to Gibbs sampling.
- Not suitable for large data sets!

Noisy, but unbiased estimates of the gradient wrt  $w_m$

$$-\mathcal{F}(w) = \sum_n \ell_n(w) - \sum_n \text{KL}(q_w(z_n) \| p(z_n)) - \sum_m \text{KL}(q_w(\theta_m) \| p(\theta_m)).$$

Noisy, but unbiased estimates of the gradient wrt  $w_m$

$$-\mathcal{F}(\mathbf{w}) = \sum_n \ell_n(\mathbf{w}) - \sum_n \text{KL}(q_{\mathbf{w}}(z_n) \| p(z_n)) - \sum_m \text{KL}(q_{\mathbf{w}}(\boldsymbol{\theta}_m) \| p(\boldsymbol{\theta}_m)).$$

$$-\frac{\partial \mathcal{F}(\mathbf{w})}{\partial w_m} = \frac{\partial}{\partial w_m} \left( \sum_n \ell_n(\mathbf{w}) - \text{KL}(q_{\mathbf{w}}(\boldsymbol{\theta}_m) \| p(\boldsymbol{\theta}_m)) \right)$$

Noisy, but unbiased estimates of the gradient wrt  $w_m$

$$-\mathcal{F}(w) = \sum_n \ell_n(w) - \sum_n \text{KL}(q_w(z_n) \| p(z_n)) - \sum_m \text{KL}(q_w(\theta_m) \| p(\theta_m)).$$

$$\begin{aligned} -\frac{\partial \mathcal{F}(w)}{\partial w_m} &= \frac{\partial}{\partial w_m} \left( \sum_n \ell_n(w) - \text{KL}(q_w(\theta_m) \| p(\theta_m)) \right) \\ &= \sum_n \frac{\partial}{\partial w_m} \left( \ell_n(w) - \frac{\text{KL}(q_w(\theta_m) \| p(\theta_m))}{N} \right). \end{aligned}$$

## Stochastic variational inference (SVI) [HBB10]

We use stochastic gradient descent in the variational M-step:

$$\mathbf{w}_m \leftarrow \mathbf{w}_m + \rho_t N \frac{\partial}{\partial \mathbf{w}_m} \left( \ell_n(\mathbf{w}) - \frac{\text{KL}(q(\boldsymbol{\theta}_m; \mathbf{w}_m) \| p(\boldsymbol{\theta}_m))}{N} \right),$$

where  $\sum_t \rho_t = \infty$  and  $\sum_t \rho_t^2 < \infty$ .

## Stochastic variational inference (SVI) [HBB10]

We use stochastic gradient descent in the variational M-step:

$$\mathbf{w}_m \leftarrow \mathbf{w}_m + \rho_t N \frac{\partial}{\partial \mathbf{w}_m} \left( \ell_n(\mathbf{w}) - \frac{\text{KL}(q(\boldsymbol{\theta}_m; \mathbf{w}_m) \| p(\boldsymbol{\theta}_m))}{N} \right),$$

where  $\sum_t \rho_t = \infty$  and  $\sum_t \rho_t^2 < \infty$ .

- **Stochastic approximation** of the gradient [RM51]:
  - ▶ Small memory footprint; **sequential** method.
  - ▶ Requires adjusting the learning rate  $\rho_t$ .
  - ▶ Monotonic increase of bound is lost (no sanity check)

# Stochastic variational inference (SVI) [HBB10]

We use stochastic gradient descent in the variational M-step:

$$\mathbf{w}_m \leftarrow \mathbf{w}_m + \rho_t N \frac{\partial}{\partial \mathbf{w}_m} \left( \ell_n(\mathbf{w}) - \frac{\text{KL}(q(\boldsymbol{\theta}_m; \mathbf{w}_m) \| p(\boldsymbol{\theta}_m))}{N} \right),$$

where  $\sum_t \rho_t = \infty$  and  $\sum_t \rho_t^2 < \infty$ .

- **Stochastic approximation** of the gradient [RM51]:
  - ▶ Small memory footprint; **sequential** method.
  - ▶ Requires adjusting the learning rate  $\rho_t$ .
  - ▶ Monotonic increase of bound is lost (no sanity check)
- SVI corresponds to **stochastic natural gradients** wrt  $q_{\mathbf{w}_m}$  for exponential family distributions [HBWP13].

## Incremental variational inference (IVI) [AE15]

$$-\mathcal{F}(\mathbf{w}) = \underbrace{\sum_n \ell_n(\mathbf{w})}_{=\ell_N(\mathbf{w})} - \sum_n \text{KL}(q(\mathbf{z}_n; \mathbf{w}_n) \| p(\mathbf{z}_n)) - \sum_m \text{KL}(q(\boldsymbol{\theta}_m; \mathbf{w}_m) \| p(\boldsymbol{\theta}_m)) .$$

## Incremental variational inference (IVI) [AE15]

$$-\mathcal{F}(\mathbf{w}) = \underbrace{\sum_n \ell_n(\mathbf{w})}_{=\ell_N(\mathbf{w})} - \sum_n \text{KL}(q(\mathbf{z}_n; \mathbf{w}_n) \| p(\mathbf{z}_n)) - \sum_m \text{KL}(q(\boldsymbol{\theta}_m; \mathbf{w}_m) \| p(\boldsymbol{\theta}_m)).$$

Let  $\mathbf{s}(X, Z) = \sum_n \mathbf{s}_n(\mathbf{x}_n, \mathbf{z}_n)$  be the vector of sufficient statistics:

$$\mathbf{w}_m \leftarrow \arg \max_{\mathbf{w}_m} \ell_N(\mathbf{s}, \mathbf{w}) - \ell_n(\mathbf{s}_n, \mathbf{w}) + \ell_n(\mathbf{s}_n^*, \mathbf{w}) - \text{KL}(q(\boldsymbol{\theta}_m; \mathbf{w}_m) \| p(\boldsymbol{\theta}_m)).$$

where  $\mathbf{s}_n^*(\mathbf{x}_n, \mathbf{z}_n)$  is the new vector of sufficient statistics.

## Incremental variational inference (IVI) [AE15]

$$-\mathcal{F}(\mathbf{w}) = \underbrace{\sum_n \ell_n(\mathbf{w})}_{=\ell_N(\mathbf{w})} - \sum_n \text{KL}(q(\mathbf{z}_n; \mathbf{w}_n) \| p(\mathbf{z}_n)) - \sum_m \text{KL}(q(\boldsymbol{\theta}_m; \mathbf{w}_m) \| p(\boldsymbol{\theta}_m)).$$

Let  $\mathbf{s}(X, Z) = \sum_n \mathbf{s}_n(\mathbf{x}_n, \mathbf{z}_n)$  be the vector of sufficient statistics:

$$\mathbf{w}_m \leftarrow \arg \max_{\mathbf{w}_m} \ell_N(\mathbf{s}, \mathbf{w}) - \ell_n(\mathbf{s}_n, \mathbf{w}) + \ell_n(\mathbf{s}_n^*, \mathbf{w}) - \text{KL}(q(\boldsymbol{\theta}_m; \mathbf{w}_m) \| p(\boldsymbol{\theta}_m)).$$

where  $\mathbf{s}_n^*(\mathbf{x}_n, \mathbf{z}_n)$  is the new vector of sufficient statistics.

- **Sequential** like SVI, but maintains a batch estimate of  $\mathbf{s}(X, Z)$ .
- Needs to store the sufficient statistics.
- No parameters to tune.
- **Monotonic** increase of bound is recovered!
- Can be interpreted as **stochastic average gradient descent** [SLB13].

## Relation to incremental EM

- MVI updates can be re-written in terms of the sufficient statistics:

$$q(z_n; w_n) \propto e^{\mathbb{E}_{z_n}(\ln p(s_n|\Theta))}, \quad q(\theta_m; w_m) \propto e^{\mathbb{E}_{\theta_m}(\ln p(s|\Theta))} p(\theta_m).$$

## Relation to incremental EM

- MVI updates can be re-written in terms of the sufficient statistics:

$$q(z_n; w_n) \propto e^{\mathbb{E}_{\neg z_n}(\ln p(s_n | \Theta))}, \quad q(\theta_m; w_m) \propto e^{\mathbb{E}_{\neg \theta_m}(\ln p(s | \Theta))} p(\theta_m).$$

- IVI updates can be re-written in a similar fashion as in incremental EM [NH93]:

$$q(z_n; w_n) \propto e^{\mathbb{E}_{\neg z_n}(\ln p(s_n^* | \Theta))}, \quad q(\theta_m; w_m) \propto e^{\mathbb{E}_{\neg \theta_m}(\ln p(s - s_n + s_n^* | \Theta))} p(\theta_m).$$

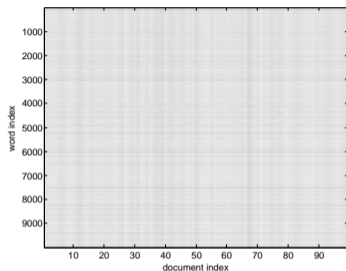
# Topic models

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

- Organise and browse large document collections.
- Capture underlying semantic structure in an unsupervised way.
- Extremely popular (e.g., more than 22k citations in Google Scholar)

# Latent Dirichlet allocation (LDA) [DMB03]



Observations are word counts per document. LDA assumes an admixture model:

$$\mathbf{X} \in \mathbb{N}^{V \times D}.$$

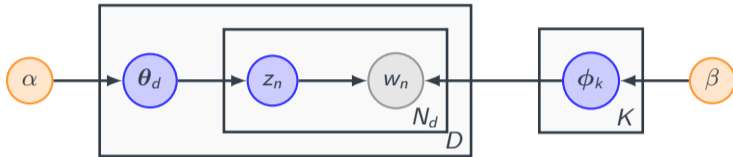
LDA infers a low-rank approximation of the matrix of counts:

$$\mathbb{E}(\mathbf{X}) \approx \mathbf{\Phi} \mathbf{\Theta}^\top,$$

$$\mathbf{x}_d \sim \text{Multinomial}(\mathbf{\Phi} \mathbf{\theta}_d, N_d)$$

where  $\mathbf{\Phi} \in \mathbb{R}_+^{V \times K}$ ,  $\mathbf{\Theta} \in \mathbb{R}_+^{D \times K}$  and  $K$  is small.

## Graphical model



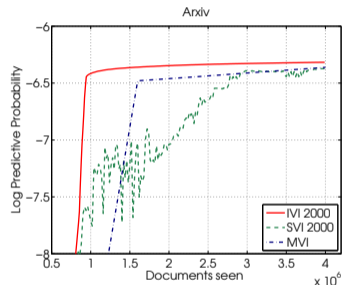
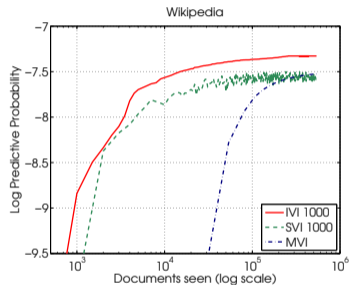
$$\theta_d \sim \text{Dirichlet}(\alpha \mathbf{1}_K),$$

$$\phi_k \sim \text{Dirichlet}(\beta \mathbf{1}_V),$$

$$z_n | \theta_d \sim \text{Categorical}(\theta_d),$$

$$w_n | z_n, \{\phi_k\}_{k=1}^K \sim \text{Categorical}(\phi_{z_n}).$$

# Log-predictive probability for LDA as a function of the number of processed documents



IVI converges faster and to a higher value on all considered datasets. ( $K=100$ ,  $\alpha_0 = 0.5$  and  $\beta_0 = 0.05$ )

Yet another form of the ELBO based on the score function

$$-\mathcal{F}(\mathbf{w}) = \mathbb{E} (\ln p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\Theta} | m)) + \mathbb{H} (q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta})) .$$

## Yet another form of the ELBO based on the score function

$$-\mathcal{F}(\mathbf{w}) = \mathbb{E} (\ln p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\Theta} | m)) + \mathbb{H} (q_{\mathbf{w}}(\mathbf{Z}, \boldsymbol{\Theta})) .$$

Write the gradient in terms of the score function:

$$\begin{aligned} -\frac{\partial \mathcal{F}(\mathbf{w})}{\partial \mathbf{w}_n} &= \mathbb{E} \left( \frac{\partial \ln q(\mathbf{z}_n; \mathbf{w}_n)}{\partial \mathbf{w}_n} (\ln p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\Theta}) - \ln q(\mathbf{z}_n; \mathbf{w}_n)) \right) \\ &\approx \frac{1}{K} \sum_{k=1}^K \left( \frac{\partial \ln q(\mathbf{z}_n^{(k)}; \mathbf{w}_n)}{\partial \mathbf{w}_n} (\ln p(\mathbf{x}_n, \mathbf{z}_n^{(k)} | \boldsymbol{\Theta}) - \ln q(\mathbf{z}_n^{(k)}; \mathbf{w}_n)) \right) , \end{aligned}$$

where  $\mathbf{z}_n^{(k)} \sim q(\mathbf{z}_n^{(k)}; \mathbf{w}_n)$ .

# Black-box variational inference [RGB14]

$$\mathbf{w}_n \leftarrow \mathbf{w}_n + \frac{\lambda_t}{K} \sum_{k=1}^K \left( \frac{\partial \ln q(\mathbf{z}_n^{(k)}; \mathbf{w}_n)}{\partial \mathbf{w}_n} \left( \ln p(\mathbf{x}_n, \mathbf{z}_n^{(k)} | \Theta) - \ln q(\mathbf{z}_n^{(k)}; \mathbf{w}_n) \right) \right),$$

where  $\sum_t \lambda_t = \infty$  and  $\sum_t \lambda_t^2 < \infty$ .

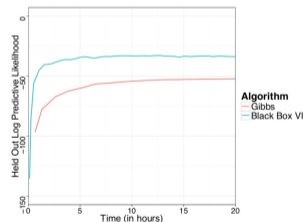


Figure 1: Comparison between Metropolis-Hastings within Gibbs and Black Box Variational Inference. In the x axis is time and in the y axis is the predictive likelihood of the test set. Black Box Variational Inference reaches better predictive likelihoods faster than Gibbs sampling. The Gibbs sampler's progress slows considerably after 5 hours.

# Black-box variational inference [RGB14]

$$w_n \leftarrow w_n + \frac{\lambda_t}{K} \sum_{k=1}^K \left( \frac{\partial \ln q(z_n^{(k)}; w_n)}{\partial w_n} \left( \ln p(x_n, z_n^{(k)} | \Theta) - \ln q(z_n^{(k)}; w_n) \right) \right),$$

where  $\sum_t \lambda_t = \infty$  and  $\sum_t \lambda_t^2 < \infty$ .

- Remove conjugacy requirement
- Variance reduction techniques:
  - ▶ Rao-Blackwellization
  - ▶ Control variates
- Can be scaled up with SVI

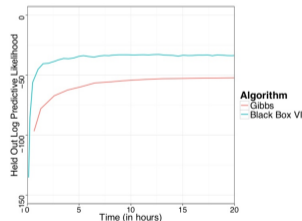
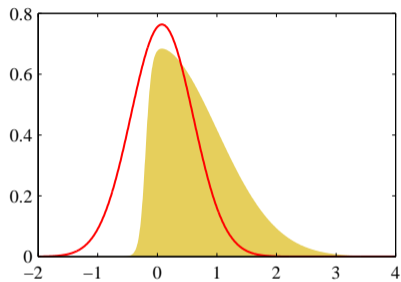
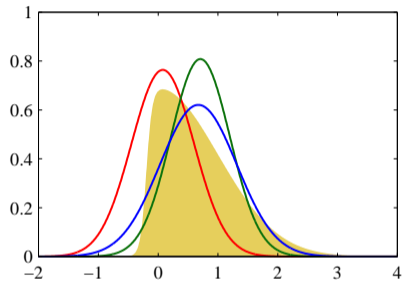


Figure 1: Comparison between Metropolis-Hastings within Gibbs and Black Box Variational Inference. In the x axis is time and in the y axis is the predictive likelihood of the test set. Black Box Variational Inference reaches better predictive likelihoods faster than Gibbs sampling. The Gibbs sampler's progress slows considerably after 5 hours.

## Other approximate inference methods



Laplace approximation.



$KL(q||p)$  vs.  $KL(p||q)$ . [Min01]

(Image credit: [Bis06].)

## Further reading

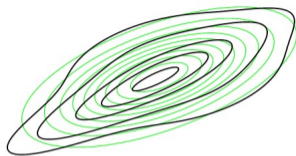
Christopher Bishop (2006): [Pattern Recognition and Machine Learning](#). [Bis06]

Kevin Murphy (2012): [Machine Learning: a Probabilistic Perspective](#). [Mur12]

David Blei, et al. (2017): [Variational Inference: a Review for Statisticians](#). [BKM17]

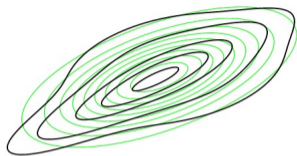
# Case Studies

## Gaussian variational approximation (GVA) [OA09]



$$q(\mathbf{z}) = \text{Gaussian}(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}), \quad -\mathcal{F}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\mathbb{E}(\ln p(\mathbf{x}, \mathbf{z})) - \underbrace{\left( \frac{N}{2} \ln 2\pi e + \frac{1}{2} \ln |\boldsymbol{\Sigma}| \right)}_{\text{entropy}}.$$

## Gaussian variational approximation (GVA) [OA09]



$$q(\mathbf{z}) = \text{Gaussian}(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}), \quad -\mathcal{F}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\mathbb{E}(\ln p(\mathbf{x}, \mathbf{z})) - \underbrace{\left( \frac{N}{2} \ln 2\pi e + \frac{1}{2} \ln |\boldsymbol{\Sigma}| \right)}_{\text{entropy}}.$$

### Interpretation of GVA:

$$\begin{aligned} 0 &= \nabla_{\boldsymbol{\mu}} \mathbb{E}(\ln p(\mathbf{x}, \mathbf{z})) = \mathbb{E}(\nabla_{\mathbf{x}} \ln p(\mathbf{x}, \mathbf{z})), \\ \boldsymbol{\Sigma}^{-1} &= -2 \nabla_{\boldsymbol{\Sigma}} \mathbb{E}(\ln p(\mathbf{x}, \mathbf{z})) = -\mathbb{E}(\nabla_{\mathbf{x}} \nabla_{\mathbf{x}} \ln p(\mathbf{x}, \mathbf{z})). \end{aligned}$$

## Application to Gaussian process (GP) models

Let's consider a factorized likelihood and assume a GP prior:

$$p(\mathbf{x}, \mathbf{z}) \propto e^{-\sum_n V(\mathbf{x}_n, \mathbf{z}_n) - \frac{1}{2} \mathbf{z}^\top \mathbf{K}^{-1} \mathbf{z}}.$$

## Application to Gaussian process (GP) models

Let's consider a factorized likelihood and assume a GP prior:

$$p(\mathbf{x}, \mathbf{z}) \propto e^{-\sum_n V(\mathbf{x}_n, \mathbf{z}_n) - \frac{1}{2} \mathbf{z}^\top \mathbf{K}^{-1} \mathbf{z}}.$$

The ELBO is given by

$$-\mathcal{F}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\sum_n \mathbb{E}(V(\mathbf{x}_n, \mathbf{z}_n)) + \frac{1}{2} \text{tr} \{ \mathbf{K}^{-1} \boldsymbol{\Sigma} \} + \frac{1}{2} \boldsymbol{\mu}^\top \mathbf{K}^{-1} \boldsymbol{\mu} - \frac{1}{2} \ln |\boldsymbol{\Sigma}| + \text{const.}$$

## Application to Gaussian process (GP) models

Let's consider a factorized likelihood and assume a GP prior:

$$p(\mathbf{x}, \mathbf{z}) \propto e^{-\sum_n V(\mathbf{x}_n, \mathbf{z}_n) - \frac{1}{2} \mathbf{z}^\top \mathbf{K}^{-1} \mathbf{z}}.$$

The ELBO is given by

$$-\mathcal{F}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\sum_n \mathbb{E}(V(\mathbf{x}_n, \mathbf{z}_n)) + \frac{1}{2} \text{tr} \{ \mathbf{K}^{-1} \boldsymbol{\Sigma} \} + \frac{1}{2} \boldsymbol{\mu}^\top \mathbf{K}^{-1} \boldsymbol{\mu} - \frac{1}{2} \ln |\boldsymbol{\Sigma}| + \text{const.}$$

Applying GVA only requires  $\mathcal{O}(N)$  variational parameters:

$$\boldsymbol{\mu} = \mathbf{K} \boldsymbol{\nu}, \quad \boldsymbol{\nu} = \left( \dots - \frac{\partial \mathbb{E} V(\mathbf{x}_n, \mathbf{z}_n)}{\partial \mu_n} \dots \right),$$

$$\boldsymbol{\Sigma}^{-1} = (\mathbf{K}^{-1} + \text{diag}\{\boldsymbol{\lambda}\}), \quad \boldsymbol{\lambda} = \left( \dots 2 \frac{\partial \mathbb{E} V(\mathbf{x}_n, \mathbf{z}_n)}{\partial \Sigma_{nn}} \dots \right).$$

## Sparse Gaussian processes [Tit09]

Let us introduce inducing points  $u$ :

$$\ln p(x) \geq \iint q(z, u) \ln \frac{\overbrace{p(x|z) p(z|u) p(u)}^{\text{GP prior}}}{\underbrace{p(z|u) q(u)}_{=q(z,u)}} dz du.$$

## Sparse Gaussian processes [Tit09]

Let us introduce inducing points  $u$ :

$$\ln p(\mathbf{x}) \geq \iint q(\mathbf{z}, \mathbf{u}) \ln \underbrace{\frac{p(\mathbf{x}|\mathbf{z}) \overbrace{p(\mathbf{z}|\mathbf{u})p(\mathbf{u})}^{\text{GP prior}}}{\underbrace{p(\mathbf{z}|\mathbf{u})q(\mathbf{u})}_{=q(\mathbf{z},\mathbf{u})}}}_{=q(\mathbf{z},\mathbf{u})} d\mathbf{z} d\mathbf{u}.$$

The ELBO is given by

$$-\mathcal{F}(\mathbf{u}) = -\frac{N}{2} \ln 2\pi - \underbrace{\frac{1}{2} \ln |\mathbf{Q}_{\mathbf{z}\mathbf{z}} + \sigma^2 \mathbf{I}_n|}_{\text{complexity}} - \underbrace{\frac{1}{2} \mathbf{x}^\top (\mathbf{Q}_{\mathbf{z}\mathbf{z}} + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{x}}_{\text{data fit}} - \underbrace{\frac{1}{2\sigma^2} \text{tr}\{\mathbf{K}_{\mathbf{z}\mathbf{z}} - \mathbf{Q}_{\mathbf{z}\mathbf{z}}\}}_{\text{approx quality}},$$

where  $\mathbf{Q}_{\mathbf{z}\mathbf{z}} = \mathbf{K}_{\mathbf{z}\mathbf{u}} \mathbf{K}_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{K}_{\mathbf{u}\mathbf{z}}$ .

# Sparse Gaussian processes [Tit09]

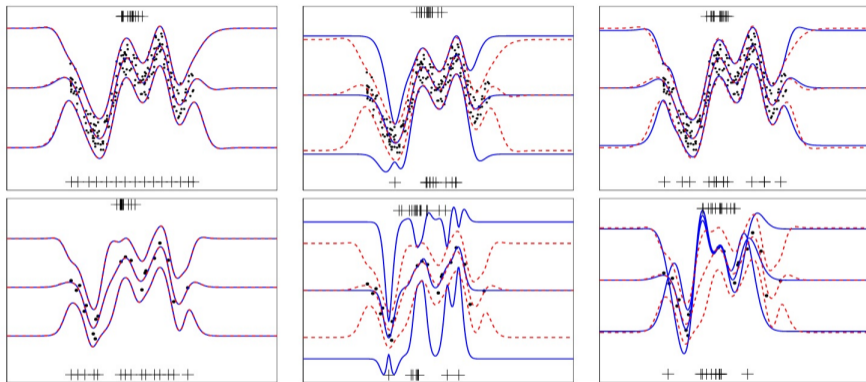


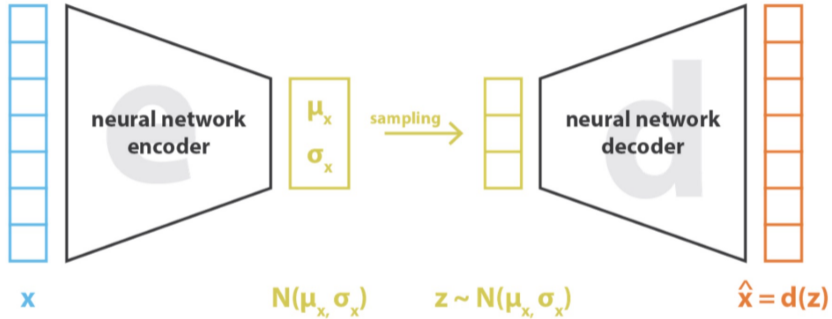
Figure 1: The first row corresponds to 200 training points and the second row to 20 training points. The first column shows the prediction (blue solid lines) obtained by maximizing  $F_V$  over the 15 pseudo-inputs and the hyperparameters. The full GP prediction is shown with red dashed lines. Initial locations of the pseudo-inputs are shown on the top as crosses, while final positions are given on the bottom as crosses. The second column shows the predictive distributions found by PP and similarly the third column for SPGP.

# Deep generative models



Face images generated with a Variational Autoencoder (source: [Wojciech Mormul on Github](#)).

# Variational auto-encoders (VAE)



(Image credit: Joseph Rocca)

## Variational auto-encoders (VAE) [KW14]

$$-\mathcal{F}(\mathbf{v}, \mathbf{w}) = \sum_n \mathbb{E} \left( \ln \underbrace{p(\mathbf{x}_n | \mathbf{z}_n)}_{\text{decoder}} \right) - \sum_n \text{KL} \left( \underbrace{q(\mathbf{z}_n | \mathbf{x}_n)}_{\text{encoder}} \| p(\mathbf{z}_n) \right).$$

where

$$p(\mathbf{x}_n | \mathbf{z}_n) = \text{Gaussian}(\mathbf{f}_{\mathbf{v}}(\mathbf{z}_n), \sigma^2 \mathbf{I}),$$

$$q(\mathbf{z}_n | \mathbf{x}_n) = \text{Gaussian}(\mathbf{g}_{\mathbf{w}}(\mathbf{x}_n), \mathbf{H}_{\mathbf{w}}(\mathbf{x}_n)),$$

$$p(\mathbf{z}_n) = \text{Gaussian}(\mathbf{0}, \mathbf{I}).$$

## Variational auto-encoders (VAE) [KW14]

$$-\mathcal{F}(\mathbf{v}, \mathbf{w}) = \sum_n \mathbb{E} \left( \ln \underbrace{p(\mathbf{x}_n | \mathbf{z}_n)}_{\text{decoder}} \right) - \sum_n \text{KL} \left( \underbrace{q(\mathbf{z}_n | \mathbf{x}_n)}_{\text{encoder}} \parallel p(\mathbf{z}_n) \right).$$

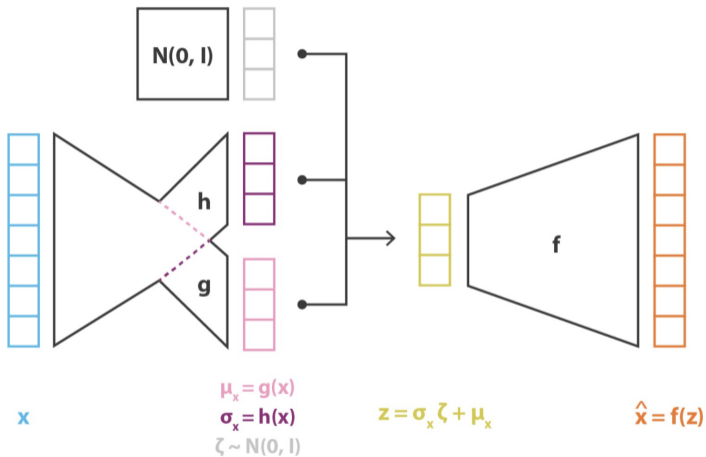
where

$$p(\mathbf{x}_n | \mathbf{z}_n) = \text{Gaussian}(\mathbf{f}_{\mathbf{v}}(\mathbf{z}_n), \sigma^2 \mathbf{I}),$$

$$q(\mathbf{z}_n | \mathbf{x}_n) = \text{Gaussian}(\mathbf{g}_{\mathbf{w}}(\mathbf{x}_n), \mathbf{H}_{\mathbf{w}}(\mathbf{x}_n)), \quad p(\mathbf{z}_n) = \text{Gaussian}(0, \mathbf{I}).$$

- What choice for  $\mathbf{f}_{\mathbf{v}}$ ,  $\mathbf{g}_{\mathbf{w}}$  and  $\mathbf{H}_{\mathbf{w}}$ ?
- Can we train de neural networks by back-propagation?

# VAE architecture and reparametrization trick



(Image credit: Joseph Rocca)

# References I



Cedric Archambeau and Beyza Ermiş.

Incremental Variational Inference for Latent Dirichlet Allocation.

ArXiv e-prints, 2015.



Cedric Archambeau and Michel Verleysen.

Robust bayesian clustering.

Neural Networks, 20(1):129 – 138, 2007.



M.J. Beal.

Variational Algorithms for Approximate Bayesian Inference.

PhD thesis, Gatsby Computational Neuroscience Unit, University College London, 2003.



Christopher M. Bishop.

Pattern Recognition and Machine Learning.

Springer, 2006.



David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe.

Variational inference: A review for statisticians.

Journal of the American Statistical Association, 112(518):859–877, 2017.

## References II



Michael I. Jordan David M. Blei, Andrew Y. Ng.

Latent dirichlet allocation.

Journal of Machine Learning Research, 3:993–1022, 2003.



Matthew Hoffman, Francis R. Bach, and David M. Blei.

Online learning for latent dirichlet allocation.

In J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, editors, Advances in Neural Information Processing Systems 23, pages 856–864. 2010.



Matthew D. Hoffman, David M. Blei, Chong Wang, and John Paisley.

Stochastic variational inference.

Journal of Machine Learning Research, 14:1303–1347, 2013.



Tom Heskes, Onno Zoeter, and Wim Wiegerinck.

Approximate expectation maximization.

In Advances in Neural Information Processing Systems 16, pages 353–360, 2003.

## References III



Diederik P. Kingma and Max Welling.

Auto-encoding variational bayes.

In Yoshua Bengio and Yann LeCun, editors, 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings, 2014.



Thomas Minka.

Expectation propagation for approximate bayesian inference.

In Daphne Koller Jack S. Breese, editor, Proceedings of the 17th Conference in Uncertainty in Artificial Intelligence, pages 362–369, 2001.



Kevin P. Murphy.

Machine learning: a probabilistic perspective.

MIT press, 2012.



Radford M. Neal and Geoffrey E. Hinton.

A new view of the em algorithm that justifies incremental and other variants.

In Learning in Graphical Models, pages 355–368. Kluwer Academic Publishers, 1993.

## References IV



Manfred Opper and Cedric Archambeau.  
The Gaussian variational approximations revisited.  
Neural Computation, 21(3):786 – 792, 2009.



Rajesh Ranganath, Sean Gerrish, and David Blei.  
Black Box Variational Inference.  
In Samuel Kaski and Jukka Corander, editors, Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics, volume 33 of Proceedings of Machine Learning Research, pages 814–822, Reykjavik, Iceland, 22–25 Apr 2014. PMLR.



Herbert Robbins and Sutton Monro.  
A stochastic approximation method.  
The Annals of Mathematical Statistics, 22(3):400–407, 1951.



Lawrence Saul and Michael I. Jordan.  
Exploiting tractable substructures in intractable networks.  
In Advances in Neural Information Processing Systems 8, pages 486–492. MIT Press, 1995.

# References V



M. Schmidt, N. Le Roux, and F. Bach.

Minimizing Finite Sums with the Stochastic Average Gradient.

ArXiv e-prints, 2013.



Michalis Titsias.

Variational Learning of Inducing Variables in Sparse Gaussian Processes.

In Proceedings of the 12th Conference in Uncertainty in Artificial Intelligence, 2009.



Wim Wiegerinck.

Variational approximations between mean field theory and the junction tree algorithm.

In Proceedings of the 16th Conference in Uncertainty in Artificial Intelligence, pages 626–633, 2000.



C. F. Jeff Wu.

On the convergence properties of the em algorithm.

The Annals of Statistics, 11(1):95–103, 1983.