

# Lecture: Variational inference

We have  $p(x)$ . We want "good"  $q(x) \approx p(x)$

## Entropy

$x, y$  - observations

$h(x)$  - information from observ.

$h(x, y) = h(x) + h(y)$ , if  $x, y$  - ind. r.v.

$p(x, y) = p(x) \cdot p(y)$ , ~ ||           

$$\log p(x, y) = \log p(x) + \log p(y)$$

$h(x) = -\log p(x)$  - information

$$\boxed{H(p) = -\sum_{x \in \mathcal{X}} p(x) \log p(x) - \text{entropy}} \\ = \mathbb{E}_p[-\log p(x)]$$

$$H(p) \rightarrow \max_p$$

Ex.

$x_1, \dots, x_k$  - possible values for r.v.  $X$   
 $\vec{p} = \{p_1, \dots, p_k\}$  - probabilities for  $x_1, \dots, x_k$   
 $\sum_{i=1}^k p_i = 1, \quad p_i > 0, \quad i = 1, \dots, k$

$$- \sum_{i=1}^k p_i \log p_i + \lambda (\sum p_i - 1) \rightarrow \max_{\vec{p}}$$

$$- \log p_i - 1 + \lambda = 0, \quad i = 1, \dots, k$$

$$\log p_i = (1 - \lambda) \Rightarrow \text{All } p_i \text{ are equal}$$

$$p_i = \exp(1 - \lambda)$$

$$H(p) = \int p(x) \log p(x) dx = \mathbb{E}_p \log p(x)$$

Ex.

$$\int p(x) dx = 1, \quad \int x p(x) dx = \mu, \quad \int (x - \mu)^2 p(x) dx = \sigma^2$$

$$L(p) = - \int p(x) \log p(x) dx + \lambda_1 (\int p(x) dx - 1) + \lambda_2 (\int x p(x) dx - \mu) + \lambda_3 (\int (x - \mu)^2 p(x) dx - \sigma^2) \rightarrow \max_p$$

$$\frac{\partial L(p)}{\partial p} = 0 \quad \text{Euler-Lagrange eq.}$$

$\frac{\partial p}{\partial x} \frac{dx}{\partial p}$   
Solution is Gaussian  $N(\mu, \sigma^2)$

$$L(p) = \int p(x) [-\log p(x) + \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2]$$

---

Mutual entropy

$$H(\vec{y} | \vec{x}) = - \iint p(\vec{x}, \vec{y}) \ln p(\vec{y} | \vec{x}) d\vec{y} d\vec{x}$$

$$H(\vec{y} | \vec{x}) + H(\vec{x}) = H(\vec{x}, \vec{y})$$

---

Kullback - Liebler divergence

$p, q$  - two distributions

with densities  $p(x), q(x)$

$KL(p || q)$  measures how different

are they

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

$$1. KL(p|q) \geq 0$$

$$\begin{aligned} & - \int p(x) \ln \frac{q(x)}{p(x)} dx \geq - \int p(x) \left( \frac{q(x)}{p(x)} - 1 \right) dx = \\ & = - \int q(x) dx + \int p(x) dx = 1 - 1 = 0 \quad \triangle \end{aligned}$$

2. Non-symmetric

Let's establish a connection between

KL divergence &

MLE (maximum likelihood estimate)

$p(x)$  is true

$q_{\theta}(x)$  - want to be close to  $p(x)$ ,

$$\theta \in \Theta \subseteq \mathbb{R}^d$$

$$KL(p \parallel q_{\theta}) \rightarrow \min_{\theta}$$

$$KL(p \parallel q_{\theta}) = - \int p(x) \ln \frac{q(x)}{p(x)} dx = - \mathbb{E}_p \ln \frac{q(x)}{p(x)} \approx$$
$$\approx \mathbb{E}_{p_e} \ln \frac{q(x)}{p(x)} = \quad (\text{p}_e\text{-empirical})$$

$$= - \frac{1}{n} \sum_{i=1}^n [\ln q_{\theta}(x_i) - \ln p(x_i)] \rightarrow \min_{\theta} \Leftrightarrow$$

$$L_{\theta}(\mathcal{D}) = \sum_{i=1}^n \ln q_{\theta}(x_i) \rightarrow \max_{\theta}$$

Mutual information

Def.  $I(x, y) = KL(p(x, y) \parallel p(x)p(y))$





$$\begin{aligned}
 KL(q(\vec{\theta}) \parallel p(\vec{\theta} | \mathcal{D})) &= \int q(\vec{\theta}) \ln \frac{q(\vec{\theta})}{p(\vec{\theta} | \mathcal{D})} d\vec{\theta} = \\
 &= \int q(\vec{\theta}) \ln \frac{q(\vec{\theta}) p(\mathcal{D})}{p(\vec{\theta}, \mathcal{D})} d\vec{\theta} = \\
 &= \int q(\vec{\theta}) \ln p(\mathcal{D}) d\vec{\theta} + \int q(\vec{\theta}) \ln q(\vec{\theta}) d\vec{\theta} - \\
 &\quad - \int q(\vec{\theta}) \ln p(\vec{\theta}, \mathcal{D}) d\vec{\theta} =
 \end{aligned}$$

$$= \ln p(\mathcal{D}) - \int q(\vec{\theta}) \ln \frac{p(\vec{\theta}, \mathcal{D})}{q(\vec{\theta})} d\vec{\theta}$$

$$ELBO(q) = \int q(\vec{\theta}) \ln \frac{p(\vec{\theta}, \mathcal{D})}{q(\vec{\theta})} d\vec{\theta}$$

1. We maximize  $ELBO(q)$  w.r.t.  $q$

$$2. \quad \ln p(\mathcal{D}) = ELBO(q) + KL(q \parallel p) \Rightarrow$$

$$ELBO(q) \leq \ln p(\mathcal{D})$$

- evidence lower bound

3. To calculate  $ELBO(q)$  we don't need  $p(\mathcal{D})$ , only  $p(\mathcal{D}|\Theta)$  and  $p(\Theta) \Rightarrow$   
 we look at unnormalized density  
 (c.t. Laplace approximation)

---

## Mean - Field Variational Bayes

We factorize  $\vec{\Theta}$  to  $\tau \{ \vec{\Theta}_1, \dots, \vec{\Theta}_\tau \}$

$$Q_{\text{MFVB}} = \{ q : q(\vec{\Theta}) = \prod_{j=1}^{\tau} q_j(\Theta_j) \}$$

$$ELBO(q) = \int q(\vec{\Theta}) \log p(\mathcal{D}|\vec{\Theta}) + \int q(\vec{\Theta}) \log \frac{1}{q(\vec{\Theta})} d\vec{\Theta} =$$

$$= \int \prod_{j=1}^{\tau} q_j(\Theta_j) \log p(\mathcal{D}|\vec{\Theta}) + \sum_{j=1}^{\tau} \int q_j(\Theta_j) \log \frac{1}{q_j(\Theta_j)} d\Theta_j$$



$$= \int q_j(\theta_j) \left[ \log p(\mathcal{Z}, \theta) \prod_{i \neq j} q_i(\theta_i) \right] d\theta + \\ + \sum_{i=1}^2 \int q_i(\theta_i) \log \frac{1}{q_i(\theta_i)} d\theta_i =$$

$$= \int q_j(\theta_j) \left[ \int \log p(\mathcal{Z}, \vec{\theta}) \prod_{i \neq j} q_i(\theta_i) d\vec{\theta}_{-j} \right] d\theta_j + \\ + \int q_j(\theta_j) \log \frac{1}{q_j(\theta_j)} d\theta_j + \text{const} \quad \textcircled{=}$$

independent of  $q_j$

$q_i$  - fixed, if  $i \neq j$

$q_j$  - what we optimize now

$$\textcircled{=} \int q_j(\theta_j) \underbrace{\mathbb{E}_{i \neq j} \log p(\mathcal{Z}, \vec{\theta})}_{\frac{1}{q_j(\theta_j)}} d\theta_j + \\ + \int q_j(\theta_j) \log \frac{1}{q_j(\theta_j)} d\theta_j + c =$$

$$= K \mathcal{L}(q_j(\theta_j) \mid \exp(\mathbb{E}_{i \neq j} \log p(\mathcal{Z} \mid \vec{\theta})))$$

$$\boxed{q^*(\theta_j) = \exp(\mathbb{E}_{i \neq j} \log p(\mathcal{Z} \mid \vec{\theta}))}$$

$$q_j(\theta_j) \propto \prod_{i=1}^n \frac{1}{\sigma_j} \exp\left(-\frac{1}{2\sigma_j^2} (\theta_j - \mu_j)^2\right)$$

Algorithm :

1. Initialize  $q_j^0(\theta_j), j=1 \dots \tau$
2. Update factors one by one  

$$q_j^t(\theta_j) \leftarrow \vec{q}^t(\vec{\theta})$$

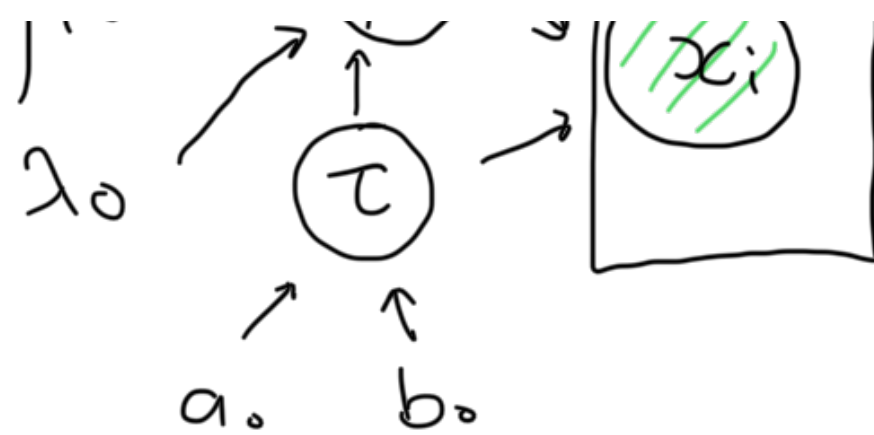
Ex.  $x_i \sim \mathcal{N}(\mu, 1/\tau)$   
 $\mathcal{D} = \{x_i\}_{i=1}^n, x_i - \text{i.i.d.}$

Likelihood  $p(\mathcal{D} | \mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{n/2} \exp\left(-\frac{\tau}{2} \sum_{i=1}^n (x_i - \mu)^2\right)$

Priors  $p(\mu | \tau) = \mathcal{N}(\mu | \mu_0, (\lambda_0 \tau)^{-1})$

$p(\tau) = \Gamma(\tau | a_0, b_0)$

$\mu_0 \rightarrow (\mu) \sim \overbrace{\text{curve}}^n$  analytical model



$$\Theta = \{\mu, \tau\} - ?$$

$p(\mu, \tau | \mathcal{D})$  - not tractable

$q(\mu, \tau) = q_\mu(\mu) q_\tau(\tau)$  - MFVB factor.

$$\begin{aligned} 1. \quad \log q_\mu^*(\mu) &= \mathbb{E}_\tau [\log p(\mathcal{D} | \mu, \tau) + \log p(\mu | \tau)] + \text{const} = \\ &= - \frac{\mathbb{E}[\tau]}{2} \left\{ \lambda_0 (\mu - \mu_0)^2 + \sum_{i=1}^n (x_i - \mu)^2 \right\} + \text{const} \end{aligned}$$

$$q_\mu^*(\mu) = \mathcal{N}(\mu | \mu_n, \lambda_n^{-1})$$

$$\mu_n = \frac{\lambda_0 \mu_0 + n \bar{x}}{\lambda_0 + n}, \quad \lambda_n = (\lambda_0 + n) \mathbb{E}[\tau]$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$E_{q^*_\tau}[\tau]$$

2. Fix  $q_\mu(\mu)$  - find  $q^*_\tau(\tau)$

$$\begin{aligned} \log q^*_\tau(\tau) = & (a_0 - 1) \log \tau - b_0 \tau + \frac{n}{2} \log \tau - \\ & - \frac{\tau}{2} E_\mu \left[ \sum_{i=1}^n (x_i - \mu)^2 + \lambda_0 (\mu - \mu_0)^2 \right] \end{aligned}$$

$$q^*_\tau(\tau) = \Gamma(\tau | a_n, b_n)$$

$$a_n = a_0 + n/2$$

$$b_n = b_0 + \frac{1}{2} E_\mu \left[ \sum_{i=1}^n (x_i - \mu)^2 + \lambda_0 (\mu - \mu_0)^2 \right]$$

Alternative approach, SVI

$$(\vec{A}) = \mathbb{E}[\vec{A} | \vec{\mu}]$$

$$q(\mathcal{D}) = p(\mathcal{D} | \mu, \Sigma)$$

$$KL(q | p(\Theta | \mathcal{D})) \rightarrow \min_{\mu, \Sigma}$$

... more on this later

---

## Expectation propagation

VV:  $KL(q | p) \rightarrow \min$

EP:  $KL(p | q) \rightarrow \min_q, q \in Q - \text{exp. family}$

$$q_{\Theta}(\vec{z}) = h(\vec{z}) \exp(\vec{\Theta}^T T(\vec{z}) + A(\vec{\Theta}))$$

$$KL(p | q) = \int p(\vec{z}) \ln \frac{p(\vec{z})}{q_{\Theta}(\vec{z})} d\vec{z} =$$

$$= \text{const} - \int p(\vec{z}) [\ln h(\vec{z}) + A(\vec{\Theta}) + \vec{\Theta}^T T(\vec{z})] d\vec{z} =$$

$$= \text{const} - A(\vec{\theta}) - \vec{\theta}^T \mathbb{E}_p T(\vec{z})$$

$$- \nabla_{\theta} A(\vec{\theta}) = \mathbb{E}_p T(\vec{z}) \approx \frac{1}{n} \sum_{i=1}^n T(\vec{z}_i), \quad \vec{z}_i \sim p(\vec{z})$$

$$\mathbb{E}_q [T(\vec{z})] = \mathbb{E}_p T(\vec{z}) - \text{method of moments}$$

$$q(\vec{z}) \sim \mathcal{N}(\vec{z} | \vec{\mu}, \Sigma),$$

$\vec{\mu}, \Sigma$  - from  $p$



