

# Tutorial 8

## Theory review → what I learnt from the Tutorial

### 1. Variational Bayesian Methods

parameters  
latent variables  
unobserved var.

data  
observed var.

→ aim: provide approximation of posterior of unobserved variables  
→ derive elbo → lower bound of marginal likelihood.

ELBO → evidence lower bound



$p(z|D)$  - not a "close form"  
→ we don't know mean, variance  
→ we cannot sample

In general:

→ directed graphical models with observed variables  $X_{R^n}$  and latent  $Z_{K^0}$   
→ we have access to joint distribution  $p(x, z)$ , we observe  $X$  as  $D$  (data)  
e.g.  $X$  - images,  $Z$  - camera angle, lighting

task:

→ Use "inference" given  $X$ , estimate latent var.  $Z$   
→ look @ a couple of picture and estimate  $Z$ .

Using Bayes' rule

$$p(z|x=D) = \frac{\overbrace{p(x=D|z)}^{\text{likelihood of data}} \cdot \overbrace{p(z)}^{\text{prior}}}{p(x=D)}$$

$p(x=D)$   
Marginal → problem is we often do not know the normalization const. given by  $p(x=D)$

$$p(x) = \int_{z_0} \dots \int_{z_{n-1}} p(x, z) dz_0 \dots dz_{n-1}$$

→ this is intractable → incomputable!  
→ complex models, e.g. images where we want to find a close form for  $10^6$  dimensions (a lot) is computable infeasible (in general)

Solution:

Using surrogate → Gaussian, we can capture most of input info  
→  $q(z) \approx p(z|x=D)$  as "good" as possible  
→ "variational" optimize for a function

Metric

→ KL-divergence

→ distance between 2 distributions:  
identical: 0  
different:  $\infty$

Optimization problem

$$q^*(z) = \underset{q(z) \in \mathcal{Q}}{\text{argmin}} (KL(q(z) || p(z|D)))$$

simple distrib  
e.g. gauss

$$KL(q(z) \| p(z|D)) = \mathbb{E}_{z \sim q(z)} \left[ \log \frac{q(z)}{p(z|D)} \right] = \int_{z_0} \dots \int_{z_{D-1}} q(z) \log \frac{q(z)}{p(z|D)} dz_0 \dots dz_{D-1} \rightarrow \text{problem, we don't have the posterior, but we have joint } p(z|D)$$

$$\approx \int_z q(z) \log \left( \frac{q(z)}{\frac{p(z,D)}{p(D)}} \right) dz = \int_z q(z) \log \left( \frac{q(z)p(D)}{p(z,D)} \right) dz = \underbrace{\int_z q(z) \log \left( \frac{q(z)}{p(z,D)} \right) dz}_{\text{expectation}}$$

$$+ \underbrace{\int_z q(z) \log p(D) dz}_{\text{expectation}} = \mathbb{E}_{z \sim q(z)} \log \frac{q(z)}{p(z,D)} + \mathbb{E}_{z \sim q(z)} \log p(D) = \overset{\text{swap denominator}}{=} - \underbrace{\mathbb{E}_q \log \frac{p(z,D)}{q(z)}}_{\mathcal{L}(q)} + \mathbb{E}_q \log p(D)$$

$\uparrow$  we have the access (GM)       $\uparrow$  const quantity

$$\underbrace{KL}_{\geq 0} = \underbrace{-\mathcal{L}(q)}_{\leq 0} + \underbrace{\log p(D)}_{\text{marginal}}$$

evidence

→ something  $\leq 0$        $\log(0:1) \leq 0$   
→ fixed quantity

$\mathcal{L}(q) \leq \log p(D)$   
 (lower-bound of the evidence)  
 → ELBO

**ELBO**  $\mathcal{L}(q) = -\mathbb{E}_q \log \frac{p(z,D)}{q(z)}$

$$\mathcal{L}(q) = \log p(D) \quad \text{i.f.f.} \quad KL(q(z) \| p(z|D)) = 0$$

i.e.  $q(z)$  is the same as  $p(z|D)$

**Maximizing ELBO  $\Rightarrow$  Minimizing KL**

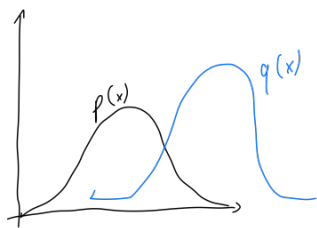
We can form an optimization problem  $\rightarrow$

**$\log p(x) = D_{KL}(q(z) \| p(z|D)) + \mathcal{L}(q)$**

this part is fixed! Even if we do not know the  $\theta$  of  $p$

This should be before but  $\rightarrow$  intuition for KL divergence, distance between distributions

What is the distance, how far are they from each other



What is the distance, how far are they from each other

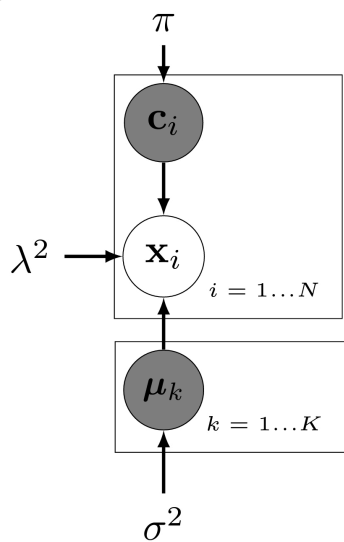
$$D_{KL}(p \| q) = \mathbb{E}_{x \sim p(x)} \left[ \log \left( \frac{p(x)}{q(x)} \right) \right] = \begin{cases} \sum_x p(x) \log \left( \frac{p(x)}{q(x)} \right) \rightarrow \text{discrete} \\ \int p(x) \log \left( \frac{p(x)}{q(x)} \right) \rightarrow \text{cont.} \end{cases}$$

How far description of one RV is away from another?

Actual tutorial starts below!



# ELBO (Assignment 1)



$\begin{cases} k - \text{components} \\ N - \text{sample size} \end{cases}$

$$\mathcal{L}(x|m, s^2, \phi) = \underbrace{\sum_{k=1}^K E_q[\log p(\mu_k)]}_{\text{I}} + \underbrace{\sum_{i=1}^N E_q[\log p(x_i)]}_{\text{II}} + \underbrace{\sum_{i=1}^N E_q[\log p(x_i|c_i, \mu)]}_{\text{III}} - \underbrace{\sum_{i=1}^N E_q[\log p(\mu_k)]}_{\text{IV}} - \underbrace{\sum_{i=1}^N E_q[\log p(c_i)]}_{\text{V}}$$

Compute closed form of ELBO

I.  $\mu_k \sim \mathcal{N}(\alpha, \Sigma, I)$

$$\sum_{k=1}^K E_q[\log p(\mu_k)] = \sum E_q[\log(\mathcal{N}(\alpha, \Sigma, I))] = \sum E_q\left[\log\left(\frac{1}{\sqrt{2\pi^p \Sigma}} e^{-\frac{1}{2}(\mu_k - \alpha)^T \Sigma^{-1} (\mu_k - \alpha)}\right)\right]$$

$$= \sum E_q\left[\log\left(\frac{1}{\sqrt{2\pi^p \Sigma}}\right) - \frac{1}{2}(\mu_k - \alpha)^T \Sigma^{-1} (\mu_k - \alpha)\right] = \sum \log\left(\frac{1}{\sqrt{2\pi^p \Sigma}}\right) - \frac{1}{2\Sigma} E[(\mu_k - \alpha)(\mu_k - \alpha)^T]$$

$\Sigma p(d-1)$

Since we calculate the expectation  $E_q$  then [according to "Var. Bay. Model Selection for mix. distr."]

$$\begin{cases} E[\mu] = m \\ E[\mu^2] = S^2 + m^2 \\ E[\text{const.}] = \text{const.} \end{cases}$$

thus,

$$= K \log\left(\frac{1}{\sqrt{2\pi^p \Sigma}}\right) - \frac{1}{2\Sigma} \sum \left( S_k^2 + m_k^T m_k - m_k^T \alpha - \alpha^T m_k + \alpha^T \alpha \right)$$

where "p" is the dim. of the data for multivar. gauss.

II.  $c_i \sim \text{Categorical}\left(\frac{1}{K}, \dots, \frac{1}{K}\right) = \log \prod_{k=1}^K \left(\frac{1}{K}\right)^{c_{i,k}} = \sum_{k=1}^K \log\left(\frac{1}{K}\right)^{c_{i,k}} = \sum_{k=1}^K -c_{i,k} \log k$

$$\sum_{i=1}^N E_q\left[\log \sum_{k=1}^K -c_{i,k} \log(k)\right] = \sum_{i=1}^N \sum_{k=1}^K \underbrace{-\phi_{i,k}}_{\substack{E[\phi_{i,k}] = \phi_k \\ \sum p(d-1 \rightarrow -1}} \log(k) = \sum_{k=1}^K -\log(k) = \underbrace{-\log(k) - \log(k) - \dots - \log(k)}_{N\text{-times}}$$

$$= -N \cdot \log k$$

\* closing the fact that  $p(x_i|c_i, \mu) = \prod_{k=1}^K p(x_i|\mu_k)^{c_{i,k}}$

III.  $x_i|c_i, \mu \sim \mathcal{N}(\alpha^T \mu, \lambda^2, I) = \prod_{k=1}^K p(x_i|\mu_k)^{c_{i,k}}$

$$\sum_{i=1}^N E_q[\log p(x_i|c_i, \mu)] = \sum_i E_q[\log p(x_i|\alpha, \mu)] = \sum_i E_q\left[\log \prod_{k=1}^K p(x_i|\mu_k)^{c_{i,k}}\right] = \sum_i E_q\left[\sum_{k=1}^K c_{i,k} \log p(x_i|\mu_k)\right] = \sum_i \sum_{k=1}^K \phi_{i,k} E_q[\log p(x_i|\mu_k)] = \sum_i \sum_{k=1}^K \phi_{i,k} E_q[\log p(x_i|\mu_k)]$$

$$\sum_{i=1}^N \sum_{k=1}^K \phi_{i,k} E_q\left[\log\left(\frac{1}{\sqrt{2\pi^p \lambda^2}} e^{-\frac{1}{2}(x_i - \mu)^T (\lambda^2 I)^{-1} (x_i - \mu)}\right)\right] =$$

$$= \sum_{i=1}^N \sum_{k=1}^K \phi_{i,k} \log\left(\frac{1}{\sqrt{2\pi^p \lambda^2}}\right) - \frac{1}{2\lambda^2} E_q\left[(x_i - \mu)^T (x_i - \mu)\right] \Rightarrow \text{similarly here as in I}$$

$$\sum_{i=1}^N \sum_{k=1}^K \phi_{i,k} \left( \log \left( \frac{1}{\sqrt{2\pi^2 \lambda^2}} \right) - \frac{1}{2\lambda^2} (x_i^T x_i - x_i^T m_k - m_k^T x_i + m_k^T m_k + S_k^2) \right)$$

$$\text{IV} \quad \sum_{k=1}^K E[\log q(\mu_k)]$$

$$q(\mu_k) \sim \mathcal{N}(\mu_k | m_k, S_k^{-1}) = \frac{1}{\sqrt{2\pi^2 S_k^2}} \cdot e^{-\frac{1}{2S_k^2} (\mu_k - m_k)^T (\mu_k - m_k)}$$

$$= \sum_{k=1}^K E_q \left[ \log \left( \frac{1}{\sqrt{2\pi^2 S_k^2}} \right) - \frac{1}{2S_k^2} (\mu_k - m_k)^T (\mu_k - m_k) \right] = \sum_{k=1}^K \log \frac{1}{\sqrt{2\pi^2 S_k^2}} - \frac{1}{2S_k^2} S_k^2 = \sum_{k=1}^K \left( \log \frac{1}{\sqrt{2\pi^2 S_k^2}} - 1 \right)$$

Expectation  $\rightarrow$

$$m_k^T m_k - m_k^T m_k - m_k^T m_k + m_k^T m_k + S_k^2 = S_k^2$$

$$\text{V} \quad \sum_{i=1}^N E[\log q(c_i)]$$

$$q(c_i) \sim \text{cat}(\phi_i) = \prod \phi_{i,k}^{c_{i,k}}$$

$$\sum_{i=1}^N E \log \prod \phi_{i,k}^{c_{i,k}} = \sum_{i=1}^N \sum_{k=1}^K E[c_{i,k} \log \phi_{i,k}] = \sum_{i=1}^N \sum_{k=1}^K \phi_{i,k} \cdot \log \phi_{i,k}$$

## Assignment 2

show that  $\phi_{i,k} \propto \exp \left\{ \frac{x_i^T E[\mu_k]}{\lambda^2} - \frac{E[\mu_k^T \mu_k]}{2\lambda^2} \right\}$

In order to show variational update, we have to calculate derivative of  $E[\mu_k]$  w.r.t.  $\phi_{i,k}$ . As suggested in 1<sup>st</sup> assig. and using the fact that  $\frac{\partial(f+g)}{\partial x} = \frac{\partial f}{\partial x} + \frac{\partial g}{\partial x}$  I will split the assignment into 5 terms

$$\text{I} \quad \frac{\partial}{\partial \phi_i} \left[ K \log \frac{1}{\sqrt{2\pi^2 \lambda^2}} - \frac{1}{2\lambda^2} \sum_{k=1}^K (S_k^2 + m_k^T m_k - m_k^T x_i - x_i^T m_k - x_i^T x_i) \right] = 0 \rightarrow \text{derivative w.r.t. } \phi_{i,k} = 0$$

$\frac{\partial \text{const}}{\partial \phi} = 0$        $\frac{\partial \text{const}}{\partial \phi} = 0$

$$\text{II} \quad \frac{\partial}{\partial \phi} -N \cdot \log K = 0$$

$$\text{III} \quad \frac{\partial}{\partial \phi_{i,k}} \sum_{i=1}^N \sum_{k=1}^K \phi_{i,k} \left( \log \left( \frac{1}{\sqrt{2\pi^2 \lambda^2}} \right) - \frac{1}{2\lambda^2} x_i^T x_i + x_i^T m_k - m_k^T x_i + m_k^T m_k + S_k^2 \right) =$$

$$= \sum_{i=1}^N \sum_{k=1}^K \log \left( \frac{1}{\sqrt{2\pi^2} \lambda^2} \right) - \frac{1}{2\lambda^2} (x_i^T x_i + x_i^T m_k - m_k^T x_i + m_k^T m_k + S_k^2)$$

$$\frac{\partial}{\partial \phi_{i,k}} \left( \log \frac{1}{\sqrt{2\pi^2} \lambda^2} - 1 \right) = 0$$

$$\frac{\partial}{\partial \phi_k} \sum_{i=1}^N \sum_{k=1}^K \phi_{i,k} \log \phi_{i,k} = \sum_{i=1}^N \sum_{k=1}^K \log \phi_{i,k} + \phi_{i,k} \cdot \frac{1}{\phi_{i,k}} = NK + \sum_{i=1}^N \sum_{k=1}^K \log \phi_{i,k}$$

Finally

$$\frac{\partial}{\partial \phi} \mathcal{L} = \sum_{i=1}^N \sum_{k=1}^K \log \left( \frac{1}{\sqrt{2\pi^2} \lambda^2} \right) - \frac{1}{2\lambda^2} (x_i^T x_i + x_i^T m_k - m_k^T x_i + m_k^T m_k + S_k^2) + \log \phi_{i,k} + 1$$

if we compare the term within summation to 0, then

$$\log \left( \frac{1}{\sqrt{2\pi^2} \lambda^2} \right) - \frac{1}{2\lambda^2} (x_i^T x_i + x_i^T m_k - m_k^T x_i + m_k^T m_k + S_k^2) + \log \phi_{i,k} + 1 = 0 \quad / \exp()$$

$$\phi_{i,k} = \frac{1}{\sqrt{2\pi^2} \lambda^2} \underbrace{\exp \left( \frac{1}{2\lambda^2} (x_i^T x_i + x_i^T m_k - m_k^T x_i + m_k^T m_k + S_k^2) \right)}_{\text{const.}} + \frac{1}{\text{const.}}$$

Since we do not care about exact value of  $\phi_{i,k}$  but the proportionality, let's focus on term  $*$ . If we look back to the assignment 1, we can see that this term

$$\text{was } \underline{E_q(x - \mu)^T (x - \mu)} = x^T E[\mu_i] + x_i^T x_i - E[\mu_i^T] x_i + E[\mu_i^T \mu_i]$$

if we assume that

$$x_i^T E[\mu_i] = E[\mu_i^T] x_i, \text{ then}$$

$$\exp \left( \frac{E[\mu_i^T \mu_i]}{2\lambda^2} - \frac{x_i^T E[\mu_i]}{\lambda^2} + \underbrace{\frac{x_i^T x_i}{2\lambda^2}}_{\text{const.}} \right), \text{ thus } \phi_{i,k} \propto \exp \left( \frac{E[\mu_i^T \mu_i]}{2\lambda^2} - \frac{x_i^T E[\mu_i]}{\lambda^2} \right)$$

# Assignment 3

$$= -\frac{\mu_k^T \mu_k}{2\sigma^2} + \sum_{i=1}^N \phi_{i,k} \left( -\frac{1}{2\lambda^2} (x_i - \mu_k)^T (x_i - \mu_k) \right) + \text{const}$$

## Assignment 3

Complete the square to find the parameters of the optimal Gaussian,  $\mathcal{N}(\mu_k, s_k^2 \mathbf{I})$ . Those parameters will be used for variational updates of the posterior of the mixture component means.

find value of  $\mu_k$  &  $s_k^2$

$$\mu_k \sim \mathcal{N}(\mu_k, s_k^2 \mathbf{I}) = \frac{1}{\sqrt{2\pi s_k^2}} \cdot e^{-\frac{1}{2} (\mu_k - \mu_k)^T s_k^{-2} (\mu_k - \mu_k)}$$

since we are comparing exp to exp we should focus on that term

$$-\frac{1}{2s_k^2} (\mu_k - \mu_k)^T (\mu_k - \mu_k) = -\frac{1}{2s_k^2} \left( \underbrace{\mu_k^T \mu_k - \mu_k^T \mu_k - \mu_k^T \mu_k + \mu_k^T \mu_k}_{\substack{\text{as in ass. 2} \\ \text{this term are equal}}} \right) = -\underbrace{\mu_k^T \mu_k \left( \frac{1}{2s_k^2} \right)}_I + \underbrace{\mu_k^T \left( \frac{\mu_k}{s_k^2} \right)}_{II}$$

$$\phi(\mu_k) \propto -\frac{\mu_k^T \mu_k}{2\sigma^2} + \sum_{i=1}^N \phi_{i,k} \left( -\frac{1}{2\lambda^2} (x_i - \mu_k)^T (x_i - \mu_k) \right) + \text{const} = -\frac{\mu_k^T \mu_k}{2\sigma^2} + \sum_i \left( \underbrace{\phi_{i,k}}_{\substack{\text{const.} \\ \text{we can skip}}} \right) \left( \underbrace{x_i^T x_i}_{\text{const.}} - \underbrace{x_i^T \mu_k}_{\text{we can skip}} - \underbrace{\mu_k^T x_i}_{\text{we can skip}} + \underbrace{\mu_k^T \mu_k}_{\text{const.}} \right) =$$

does not matter

$$= -\frac{\mu_k^T \mu_k}{2\sigma^2} + \sum_i \frac{\phi_{i,k} \mu_k^T \mu_k + \phi_{i,k} \cdot x_i \mu_k^T + \phi_{i,k} \mu_k^T x_i}{2\lambda^2} = \underbrace{\mu_k^T \mu_k \left( \frac{1}{2\sigma^2} + \frac{\sum_i \phi_{i,k}}{2\lambda^2} \right)}_{I^*} + \underbrace{\mu_k^T \left( \frac{\sum_i \phi_{i,k} \cdot x_i}{\lambda^2} \right)}_{II^*}$$

Now, by comparing  $I^*$  and  $II^*$  with  $I$  and  $II$

$$\begin{cases} \frac{1}{2s_k^2} = \frac{1}{\sigma^2} + \sum_i \frac{\phi_{i,k}}{2\lambda^2} \Rightarrow s_k = 2 \left( \frac{1}{\sigma^2} + \sum_i \frac{\phi_{i,k}}{2\lambda^2} \right)^{-1/2} \\ \frac{\mu_k}{s_k^2} = \frac{\sum_i \phi_{i,k} \cdot x_i}{\lambda^2} \Rightarrow \mu_k = \underbrace{s_k^2}_{\substack{\text{calc.} \\ \text{above}}} \frac{\sum_i \phi_{i,k} \cdot x_i}{\lambda^2} \end{cases}$$

$$s_k = \frac{1}{2 \left( \frac{1}{\sigma^2} + \sum_i \frac{\phi_{i,k}}{2\lambda^2} \right)^{1/2}}$$