

Variational Bayesian Learning

Shinichi nakajima (TU Berlin)

nakajima@tu-berlin.de

<http://sites.google.com/site/shinnkj23/>

Matrix Factorization Model

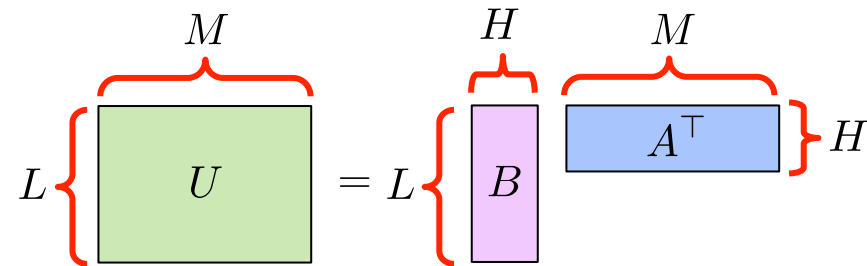
Ex 1. Collaborative Filtering

$$V = \begin{matrix} & \begin{matrix} \text{Movie 1} \\ \text{Movie 2} \\ \text{Movie 3} \\ \dots \\ \text{Movie M} \end{matrix} & \\ \begin{pmatrix} 3 & * & * & 5 & * & * & \dots & * \\ * & * & 2 & * & * & 5 & \dots & 5 \\ * & * & * & * & * & * & \dots & * \\ 3 & 1 & 4 & * & 4 & * & \dots & * \\ * & 3 & 3 & * & * & 5 & \dots & * \\ & & & & \vdots & & & \\ * & * & 5 & 5 & * & 5 & \dots & * \end{pmatrix} & \begin{matrix} \text{User 1} \\ \text{User 2} \\ \text{User 3} \\ \vdots \\ \text{User L} \end{matrix} \end{matrix}$$

Predict user preferences,
 which are not observed.

One can predict * by fitting the data
 with small degrees of freedom.

Low-rank Approximation



Approximate $V \in \mathbb{R}^{L \times M}$ with a low-rank matrix

$$V = U + \mathcal{E}$$

which can be naturally expressed as the product of two matrices:

$$U = BA^\top$$

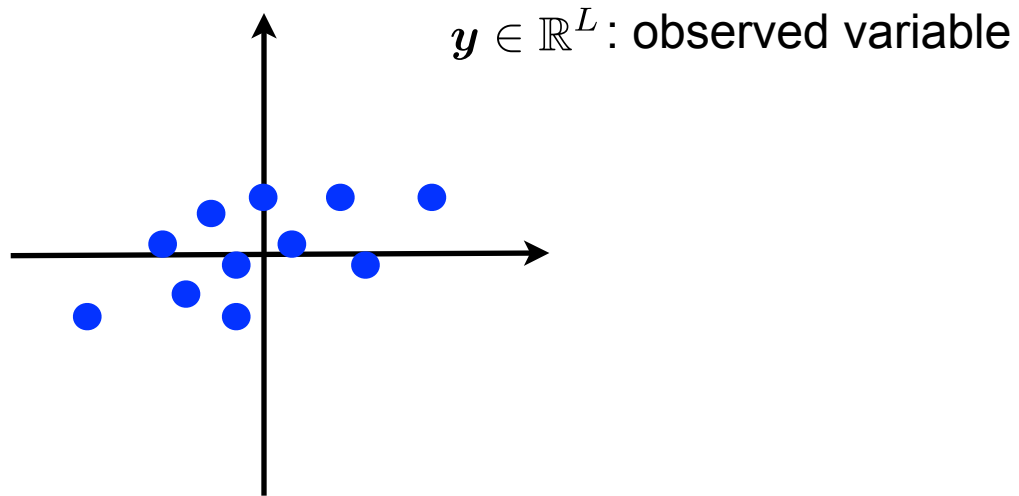
Here, $A \in \mathbb{R}^{M \times H}$, $B \in \mathbb{R}^{L \times H}$

Assuming Gaussian noise, the likelihood is given as

$$p(V|A, B) \propto \exp\left(-\frac{\|V - BA^\top\|_{\text{Fro}}^2}{2\sigma^2}\right) \quad \|V\|_{\text{Fro}}^2 = \sum_{l,m} V_{l,m}^2$$

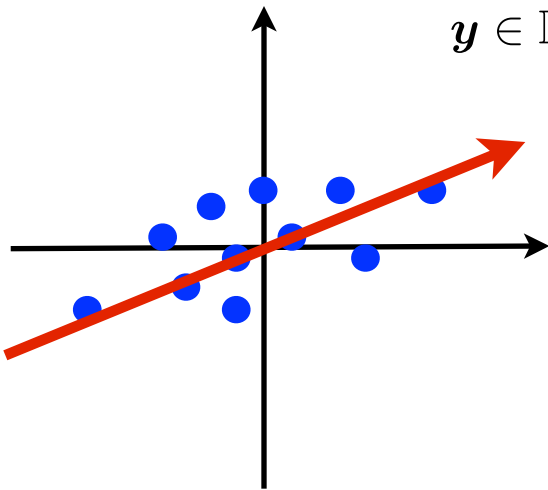
※Missing entries should be ignored.

Ex. 2: Probabilistic PCA

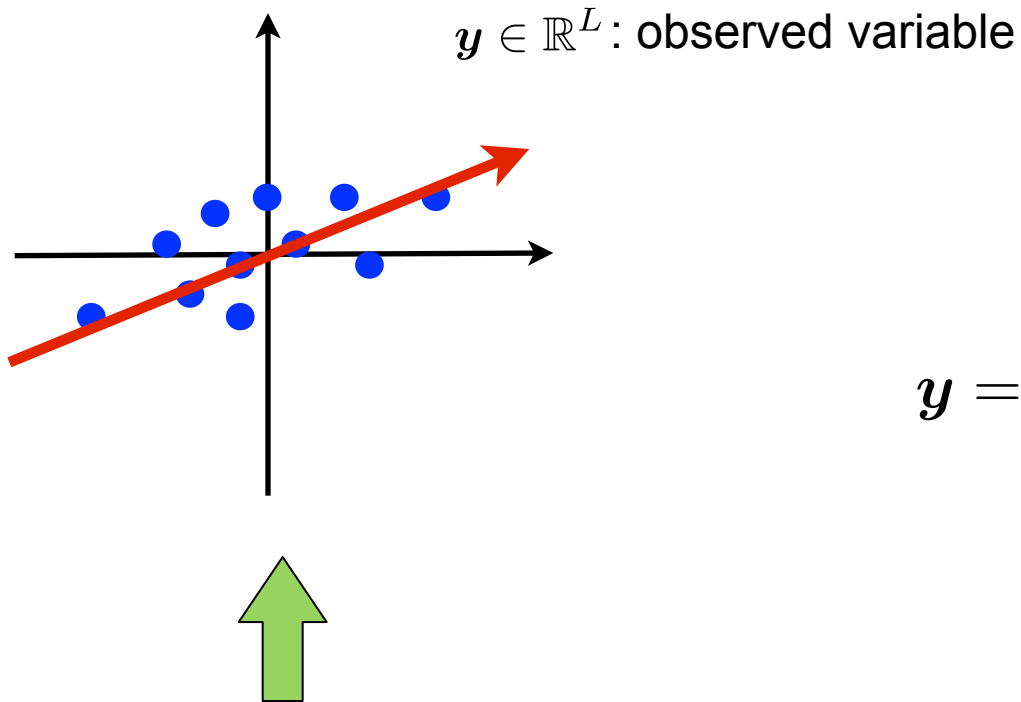


Ex. 2: Probabilistic PCA

$y \in \mathbb{R}^L$: observed variable



Ex. 2: Probabilistic PCA



$$L \left\{ \begin{array}{c} \color{green}{\mathbf{y}} \end{array} \right\} = L \left\{ \begin{array}{c} \color{purple}{\mathbf{B}} \\ \color{blue}{\mathbf{x}} \end{array} \right\}^H$$

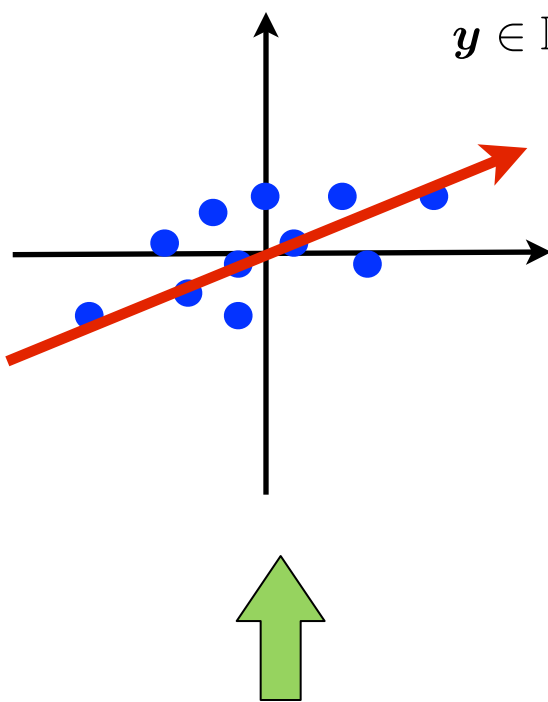
$$\mathbf{y} = \mathbf{B}\mathbf{x} + \boldsymbol{\varepsilon}$$

$\mathbf{B} \in \mathbb{R}^{L \times H}$: projection matrix 



$\mathbf{x} \in \mathbb{R}^H$: hidden variables

Ex. 2: Probabilistic PCA



$\mathbf{y} \in \mathbb{R}^L$: observed variable

$$L \left\{ \begin{array}{c} \mathbf{y} \end{array} \right\} = L \left\{ \begin{array}{c} \overbrace{B}^H \\ \mathbf{x} \end{array} \right\} H$$

For M samples,

$$Y = BX + \mathcal{E}$$

$B \in \mathbb{R}^{L \times H}$: projection matrix

$$Y = (\mathbf{y}_1, \dots, \mathbf{y}_M) \in \mathbb{R}^{L \times M}$$

$$X = (\mathbf{x}_1, \dots, \mathbf{x}_M) \in \mathbb{R}^{H \times M}$$

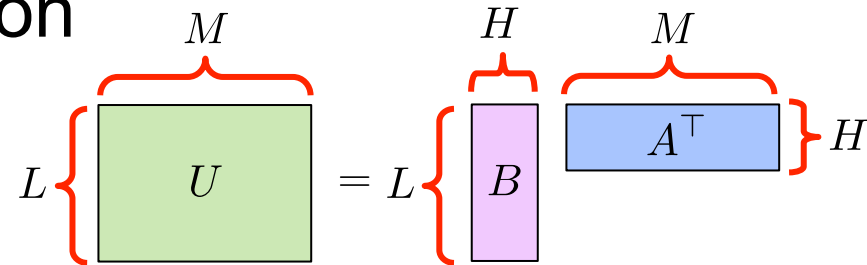
Since X is also unknown, we should estimate both B and X .



$\mathbf{x} \in \mathbb{R}^H$: hidden variables

$Y \rightarrow V, X \rightarrow A^\top$ gives ...

Bayesian Matrix Factorization



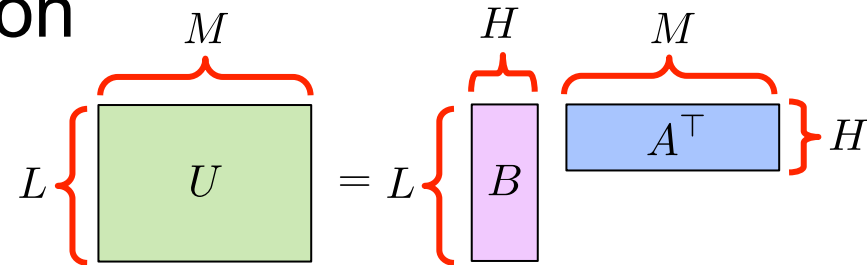
Likelihood: $p(V|A, B) \propto \exp\left(-\frac{\|V - BA^T\|_{\text{Fro}}^2}{2\sigma^2}\right)$

Observation:

$$V \in \mathbb{R}^{L \times M}$$

Variational Bayesian Learning

Bayesian Matrix Factorization



Likelihood: $p(V|A, B) \propto \exp\left(-\frac{\|V - BA^T\|_{\text{Fro}}^2}{2\sigma^2}\right)$

Observation:
 $V \in \mathbb{R}^{L \times M}$

No conjugate prior because of the quartic term:

$$p(A, B|V) \propto \exp\left(-\frac{\|V - BA^T\|_{\text{Fro}}^2}{2\sigma^2}\right) \cdot p(A, B)$$

But this is conditionally conjugate:

$$p(A|B, V) \propto \exp\left(-\frac{\|V - BA^T\|_{\text{Fro}}^2}{2\sigma^2}\right) \cdot p(A)$$

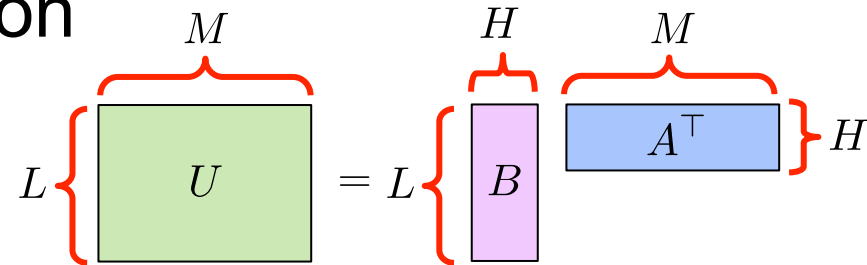
$$p(B|A, V) \propto \exp\left(-\frac{\|V - BA^T\|_{\text{Fro}}^2}{2\sigma^2}\right) \cdot p(B)$$

Independence makes
 Bayesian learning tractable!



$$r(A, B) = r_A(A)r_B(B)$$

Bayesian Matrix Factorization



Likelihood: $p(V|A, B) \propto \exp\left(-\frac{\|V - BA^T\|_{\text{Fro}}^2}{2\sigma^2}\right)$

Prior: $p(A) \propto \exp\left(-\frac{\text{tr}(AC_A^{-1}A^T)}{2}\right)$

$p(B) \propto \exp\left(-\frac{\text{tr}(BC_B^{-1}B^T)}{2}\right)$

Observation:

$$V \in \mathbb{R}^{L \times M}$$

Parameters:

$$A \in \mathbb{R}^{M \times H}$$

$$B \in \mathbb{R}^{L \times H}$$

Hyperparameters:

$$C_A = \text{diag}(c_{a_1}^2, \dots, c_{a_H}^2)$$

$$C_B = \text{diag}(c_{b_1}^2, \dots, c_{b_H}^2)$$

Free energy minimization

Trial distribution: $r(A, B)$

$$\begin{aligned} \text{Free energy: } F(r) &= \left\langle \log \frac{r(A, B)}{p(V|A, B)p(A)p(B)} \right\rangle_{r(A, B)} \\ &= \underbrace{\text{KL}(r(A, B) \| p(A, B|V))}_{\text{KL divergence to Bayes posterior}} - \underbrace{\log p(V)}_{\text{marginal likelihood (not depend on r)}} \end{aligned}$$

KL divergence to Bayes posterior

marginal likelihood
(not depend on r)



$$p(A, B|V) = \underset{r}{\operatorname{argmin}} F(r)$$

Minimizing free energy should give Bayes posterior, but intractable.

Minimize free energy under some constraints!

$$\min_r F(r)$$

$$\text{s.t. } r(A, B) \in \mathcal{G}$$

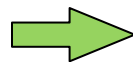
Free energy minimization

Trial distribution: $r(A, B)$

$$\begin{aligned} \text{Free energy: } F(r) &= \left\langle \log \frac{r(A, B)}{p(V|A, B)p(A)p(B)} \right\rangle_{r(A, B)} \\ &= \underbrace{\text{KL}(r(A, B) \| p(A, B|V))}_{\text{KL divergence to Bayes posterior}} - \underbrace{\log p(V)}_{\text{marginal likelihood (not depend on r)}} \end{aligned}$$

KL divergence to Bayes posterior

marginal likelihood
(not depend on r)



$$p(A, B|V) = \underset{r}{\operatorname{argmin}} F(r)$$

Minimizing free energy should give Bayes posterior, but intractable.

Minimize free energy under some constraints!

$$\min_r F(r)$$

$$\text{s.t. } r(A, B) = r_A(A)r_B(B)$$

Variational method

Free energy

$$F(r) = \left\langle \log \frac{r(A, B)}{p(V|A, B)p(A)p(B)} \right\rangle_{r(A, B)}$$

$$= \int r_A(A)r_B(B) \log \frac{r_A(A)r_B(B)}{p(V|A, B)p(A)p(B)} dA dB$$

$$\forall A \in \mathbb{R}^{M \times H},$$

$$0 = \delta I = \frac{\partial F}{\partial r_A} = \int r_B(B) \left(\log \frac{r_A(A)r_B(B)}{p(V|A, B)p(A)p(B)} + 1 \right) dB$$

$$= \left\langle \log \frac{r_A(A)r_B(B)}{p(V|A, B)p(A)p(B)} + 1 \right\rangle_{r_B(B)}$$

$$= \underbrace{\log r_A(A)} + \underbrace{\langle \log r_B(B) - \log (p(V|A, B)p(A)p(B)) \rangle_{r_B(B)}} + 1$$

$$= \log r_A(A) - \log p(A) - \langle \log p(V|A, B) \rangle_{r_B(B)} + \text{const.}$$

$$= \log \frac{r_A(A)}{p(A) \exp \left(\langle \log p(V|A, B) \rangle_{r_B(B)} \right)} + \text{const.}$$

Variational method

$$0 = \log \frac{r_A(A)}{p(A) \exp \left(\langle \log p(V|A, B) \rangle_{r_B(B)} \right)} + \text{const.}$$

This must hold for **any** A !

$$r_A(A) \propto p(A) \exp \left(\langle \log p(V|A, B) \rangle_{r_B(B)} \right)$$

Similarly, considering $\frac{\partial F}{\partial r_B}$ gives

$$r_B(B) \propto p(B) \exp \left(\langle \log p(V|A, B) \rangle_{r_A(A)} \right)$$

Variational method

Local minimizer of $\min_r F(r)$
s.t. $r(A, B) = r_A(A)r_B(B)$

satisfies $r_A(A) \propto p(A) \exp \left(\langle \log p(V|A, B) \rangle_{r_B(B)} \right)$
 $r_B(B) \propto p(B) \exp \left(\langle \log p(V|A, B) \rangle_{r_A(A)} \right)$

Variational method

Local minimizer of $\min_r F(r)$
 s.t. $r(A, B) = r_A(A)r_B(B)$

$$p(V|A, B) \propto \exp\left(-\frac{\|V - BA^\top\|_{\text{Fro}}^2}{2\sigma^2}\right)$$

$$p(A) \propto \exp\left(-\frac{\text{tr}(AC_A^{-1}A^\top)}{2}\right)$$

$$p(B) \propto \exp\left(-\frac{\text{tr}(BC_B^{-1}B^\top)}{2}\right)$$

satisfies $r_A(A) \propto p(A) \exp\left(\underbrace{\langle \log p(V|A, B) \rangle_{r_B(B)}}\right)$
 $r_B(B) \propto p(B) \exp\left(\langle \log p(V|A, B) \rangle_{r_A(A)}\right)$

Focusing on dependence on A, we have

$$\begin{aligned} \langle \log p(V|A, B) \rangle_{r(B)} &= -\frac{1}{2\sigma^2} \langle \|V - BA^\top\|_{\text{Fro}}^2 \rangle_{r(B)} \\ &= -\frac{1}{2\sigma^2} \langle \text{tr}(-2V^\top BA^\top + AB^\top BA^\top) \rangle_{r(B)} + \text{const.} \\ &= -\frac{1}{2\sigma^2} \text{tr}\left(-2V^\top \langle B \rangle_{r_B(B)} A^\top + A \langle B^\top B \rangle_{r_B(B)} A^\top\right) + \text{const.} \end{aligned}$$

$$r_A(A) \propto \exp\left(-\frac{1}{2\sigma^2} \text{tr}\left(-2V^\top \langle B \rangle_{r_B(B)} A^\top + A \langle B^\top B \rangle_{r_B(B)} A^\top + \sigma^2 AC_A^{-1}A^\top\right)\right)$$

$r_A(A)$ is Gaussian!

Variational method

$$r_A(A) \propto \exp \left(-\frac{\text{tr} \left((A - \hat{A}) \hat{\Sigma}_A^{-1} (A - \hat{A}) \right)}{2} \right)$$

$$\hat{A} = V^\top \langle B \rangle_{r_B(B)} \frac{\hat{\Sigma}_A}{\sigma^2} \quad \hat{\Sigma}_A = \sigma^2 \left(\langle B^\top B \rangle_{r_B(B)} + \sigma^2 C_A^{-1} \right)^{-1}$$

Similarly we have

$$r_B(B) \propto \exp \left(-\frac{\text{tr} \left((B - \hat{B}) \hat{\Sigma}_B^{-1} (B - \hat{B}) \right)}{2} \right)$$

$$\hat{B} = V^\top \langle A \rangle_{r_A(A)} \frac{\hat{\Sigma}_B}{\sigma^2} \quad \hat{\Sigma}_B = \sigma^2 \left(\langle A^\top A \rangle_{r_A(A)} + \sigma^2 B_B^{-1} \right)^{-1}$$

Since we now know $r_A(A)$ and $r_B(B)$ are Gaussian, the moments are given as

$$\langle A \rangle_{r_A(A)} = \hat{A} \quad \langle A^\top A \rangle_{r_A(A)} = \hat{A}^\top \hat{A} + M \hat{\Sigma}_A$$

$$\langle B \rangle_{r_B(B)} = \hat{B} \quad \langle B^\top B \rangle_{r_B(B)} = \hat{B}^\top \hat{B} + L \hat{\Sigma}_B$$

VB posterior

$$r_A(A) \propto \exp \left(-\frac{\text{tr} \left((A - \hat{A}) \hat{\Sigma}_A^{-1} (A - \hat{A}) \right)}{2} \right)$$

$$r_B(B) \propto \exp \left(-\frac{\text{tr} \left((B - \hat{B}) \hat{\Sigma}_B^{-1} (B - \hat{B}) \right)}{2} \right)$$

$$\hat{A} = V^\top \hat{B} \frac{\hat{\Sigma}_A}{\sigma^2} \quad \hat{\Sigma}_A = \sigma^2 \left(\hat{B}^\top \hat{B} + L \hat{\Sigma}_B + \sigma^2 C_A^{-1} \right)^{-1}$$

$$\hat{B} = V \hat{A} \frac{\hat{\Sigma}_B}{\sigma^2} \quad \hat{\Sigma}_B = \sigma^2 \left(\hat{A}^\top \hat{A} + M \hat{\Sigma}_A + \sigma^2 C_B^{-1} \right)^{-1}$$

Iterative algorithm

Starting from initial values for $\hat{A}, \hat{\Sigma}_A, \hat{B}, \hat{\Sigma}_B$, they are updated by equations above.

Hyperparameter Estimation

Likelihood: $p(V|A, B) \propto \exp\left(-\frac{\|V - BA^\top\|_{\text{Fro}}^2}{2\sigma^2}\right)$

Prior: $p(A) \propto \exp\left(-\frac{\text{tr}(AC_A^{-1}A^\top)}{2}\right)$

$$p(B) \propto \exp\left(-\frac{\text{tr}(BC_B^{-1}B^\top)}{2}\right)$$

(σ^2, C_A, C_B) can also be estimated by minimizing free energy

$$F = \underbrace{\text{KL}(r(A, B) \| p(A, B|V, \sigma^2, C_A, C_B))}_{\text{minimize KL divergence to Bayes posterior}} - \underbrace{\log p(V|\sigma^2, C_A, C_B)}_{\text{make marginal likelihood larger}}$$

minimize KL divergence to Bayes posterior, and make marginal likelihood larger.

Hyperparameter Estimation

F can be explicitly written as a function of $(A, B, \Sigma_A, \Sigma_B, C_A, C_B, \sigma^2)$.

(will be a homework)

Differentiating F, we have update rules for hyperparameters:

$$c_{a_h}^2 = \|\hat{\mathbf{a}}_h\|^2 / M + (\Sigma_A)_{hh}$$

$$c_{b_h}^2 = \|\hat{\mathbf{b}}_h\|^2 / L + (\Sigma_B)_{hh}$$

$$\sigma^2 = \frac{\|V\|_{\text{Fro}}^2 - \text{tr}(2V^\top \hat{B} \hat{A}^\top) + \text{tr} \left((\hat{A}^\top \hat{A} + M \Sigma_A) (\hat{B}^\top \hat{B} + L \Sigma_B) \right)}{LM}$$

Mixture of Gaussians

Mixture of Gaussians

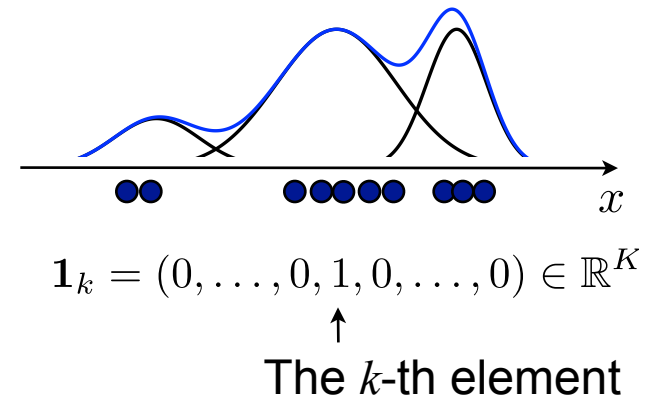
Observed var.: $\mathbf{x} \in \mathbb{R}^M$

Latent var.: $\mathbf{z} \in \{\mathbf{1}_k\}_{k=1}^K$

parameters: $\{a_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k; 0 \leq a_k \leq 1, \sum_{k=1}^K a_k = 1\}$

likelihood: $p(\{\mathbf{x}^{(n)}, \mathbf{z}^{(n)}\}_{n=1}^N | \{a_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K) = \prod_{n=1}^N \prod_{k=1}^K \left(a_k \mathcal{N}_M(\mathbf{x}^{(n)}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right)^{z_k^{(n)}}$

Prior: $p(\{a_k\}_{k=1}^K | \phi) \propto \prod_{k=1}^K a_k^{\phi-1}$ $p(\{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K) \propto \prod_{k=1}^K \mathcal{N}_M(\boldsymbol{\mu}_k; \mathbf{d}_0, c_\mu^2 \boldsymbol{\Sigma}) \mathcal{W}(\boldsymbol{\Sigma}^{-1}; \rho, \Lambda)$



Mixture of Gaussians (known covariance)

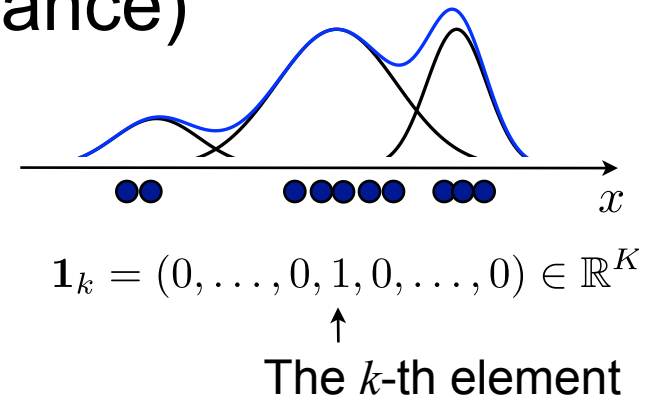
Observed var.: $\mathbf{x} \in \mathbb{R}^M$

Latent var.: $\mathbf{z} \in \{\mathbf{1}_k\}_{k=1}^K$

parameters: $\{a_k, \boldsymbol{\mu}_k, \Sigma_k; 0 \leq a_k \leq 1, \sum_{k=1}^K a_k = 1\}$

likelihood: $p(\{\mathbf{x}^{(n)}, \mathbf{z}^{(n)}\}_{n=1}^N | \{a_k, \boldsymbol{\mu}_k\}_{k=1}^K) = \prod_{n=1}^N \prod_{k=1}^K \left(a_k \mathcal{N}_M(\mathbf{x}^{(n)}; \boldsymbol{\mu}_k, \mathbf{I}_M) \right)^{z_k^{(n)}}$

Prior: $p(\{a_k\}_{k=1}^K | \phi) \propto \prod_{k=1}^K a_k^{\phi-1}$ $p(\{\boldsymbol{\mu}_k\}_{k=1}^K) \propto \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k; 0, \sigma_0^2 \mathbf{I}_M)$



Mixture of Gaussians (known covariance)

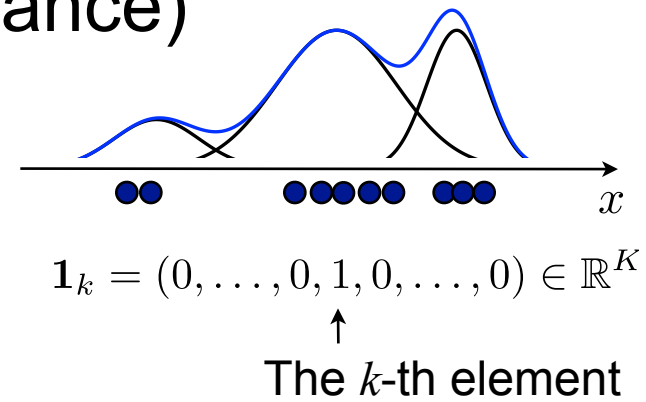
Observed var.: $\mathbf{x} \in \mathbb{R}^M$

Latent var.: $\mathbf{z} \in \{\mathbf{1}_k\}_{k=1}^K$

parameters: $\{a_k, \boldsymbol{\mu}_k, \Sigma_k; 0 \leq a_k \leq 1, \sum_{k=1}^K a_k = 1\}$

likelihood: $p(\{\mathbf{x}^{(n)}, \mathbf{z}^{(n)}\}_{n=1}^N | \{a_k, \boldsymbol{\mu}_k\}_{k=1}^K) = \prod_{n=1}^N \prod_{k=1}^K \left(a_k \mathcal{N}_M(\mathbf{x}^{(n)}; \boldsymbol{\mu}_k, \mathbf{I}_M) \right)^{z_k^{(n)}}$

Prior: $p(\{a_k\}_{k=1}^K | \phi) \propto \prod_{k=1}^K a_k^{\phi-1}$ $p(\{\boldsymbol{\mu}_k\}_{k=1}^K) \propto \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k; 0, \sigma_0^2 \mathbf{I}_M)$



Independence constraint:

$$r(\mathcal{H}, \mathbf{w}) = r_{\mathcal{H}}(\mathcal{H}) r_w(\mathbf{w})$$

where $\mathcal{H} = \{\mathbf{z}^{(n)}\}_{n=1}^N$,
 $\mathbf{w} = \{a_k, \boldsymbol{\mu}_k\}_{k=1}^K$

Variational Bayesian Learning

$$\begin{aligned} & \min F(r), \\ \text{s.t. } & r(\mathcal{H}, \boldsymbol{\omega}) = r_{\mathcal{H}}(\mathcal{H})r_{\boldsymbol{\omega}}(\boldsymbol{\omega}) \end{aligned}$$

$$\begin{aligned} \text{where } F(r) &= \left\langle \log \frac{r_{\mathcal{H}}(\mathcal{H})r_{\boldsymbol{\omega}}(\boldsymbol{\omega})}{p(\mathcal{D}, \mathcal{H}|\boldsymbol{\omega})p(\boldsymbol{\omega})} \right\rangle_{r_{\mathcal{H}}(\mathcal{H})r_{\boldsymbol{\omega}}(\boldsymbol{\omega})} \\ &= \sum_{\mathcal{H}} \int r_{\mathcal{H}}(\mathcal{H})r_{\boldsymbol{\omega}}(\boldsymbol{\omega}) \log \frac{r_{\mathcal{H}}(\mathcal{H})r_{\boldsymbol{\omega}}(\boldsymbol{\omega})}{p(\mathcal{D}, \mathcal{H}|\boldsymbol{\omega})p(\boldsymbol{\omega})} d\boldsymbol{\omega} \end{aligned}$$

Applying variational method gives

$$r_{\mathcal{H}}(\mathcal{H}) \propto \exp \langle \log p(\mathcal{D}, \mathcal{H}|\boldsymbol{\omega}) \rangle_{r_{\boldsymbol{\omega}}(\boldsymbol{\omega})}$$

$$r_{\boldsymbol{\omega}}(\boldsymbol{\omega}) \propto p(\boldsymbol{\omega}) \exp \langle \log p(\mathcal{D}, \mathcal{H}|\boldsymbol{\omega}) \rangle_{r_{\mathcal{H}}(\mathcal{H})}$$

Variational Bayesian Learning

Thus,

$$r_z(\{\mathbf{z}^{(n)}\}_{n=1}^N) \propto \prod_{n=1}^N \left(\prod_{k=1}^K (\bar{z}_k^{(n)})^{z_k^{(n)}} \right) \quad \text{Multinomial}$$

where
$$\bar{z}_k^{(n)} = \exp \left(\langle \log \alpha_k \rangle_{r_{\alpha, \mu}(\alpha, \{\mu_k\}_{k=1}^K)} - \frac{1}{2} \left\langle \|\mathbf{x}^{(n)} - \mu_k\|^2 \right\rangle_{r_{\alpha, \mu}(\alpha, \{\mu_k\}_{k=1}^K)} \right)$$

$$r_{\alpha, \mu}(\alpha, \{\mu_k\}_{k=1}^K) \propto \prod_{k=1}^K \alpha_k^{\bar{N}_k + \phi - 1} \exp \left(- \frac{(\bar{N}_k + \sigma_0^{-2}) \left\| \mu_k - \frac{\bar{N}_k \bar{\mathbf{x}}_k}{\bar{N}_k + \sigma_0^{-2}} \right\|^2}{2} \right)$$

Dirichlet x Gauss (independent!)

where
$$\bar{N}_k = \sum_{n=1}^N \langle z_k^{(n)} \rangle_{r_z(\{\mathbf{z}^{(n)}\}_{n=1}^N)}$$

$$\bar{\mathbf{x}}_k = \frac{1}{\bar{N}_k} \sum_{n=1}^N \mathbf{x}^{(n)} \langle z_k^{(n)} \rangle_{r_z(\{\mathbf{z}^{(n)}\}_{n=1}^N)}$$

Variational Bayesian Learning

$$\Psi(x) = \frac{\partial \log \Gamma(x)}{\partial x}$$

Since we now know form of posteriors,

$$r_z(\{\mathbf{z}^{(n)}\}_{n=1}^N) = \prod_{n=1}^N \text{Multinomial}_{K,1}(\mathbf{z}^{(n)}; \hat{\mathbf{z}}^{(n)})$$

$$r_\alpha(\boldsymbol{\alpha}) = \text{Dirichlet}(\boldsymbol{\alpha}; \hat{\boldsymbol{\alpha}})$$

$$r_\mu(\{\boldsymbol{\mu}_k\}_{k=1}^K) = \prod_{k=1}^K \text{Norm}_M(\boldsymbol{\mu}_k; \hat{\boldsymbol{\mu}}_k, \hat{\sigma}_k^2 \mathbf{I}_M)$$

we can compute the expectations

$$\langle z_k^{(n)} \rangle_{r_z(\{\mathbf{z}^{(n)}\}_{n=1}^N)} = \hat{z}_k^{(n)}$$

$$\begin{aligned} \langle \log \alpha_k \rangle_{r_{\alpha, \mu}(\boldsymbol{\alpha}, \{\boldsymbol{\mu}_k\}_{k=1}^K)} &= \langle \log \alpha_k \rangle_{r_\alpha(\boldsymbol{\alpha})} \\ &= \Psi(\hat{\alpha}_k) - \Psi\left(\sum_{k'=1}^K \hat{\alpha}_{k'}\right) \end{aligned}$$

$$\begin{aligned} \langle \|\mathbf{x}^{(n)} - \boldsymbol{\mu}_k\|^2 \rangle_{r_{\alpha, \mu}(\boldsymbol{\alpha}, \{\boldsymbol{\mu}_k\}_{k=1}^K)} &= \langle \|\mathbf{x}^{(n)} - \hat{\boldsymbol{\mu}}_k + (\hat{\boldsymbol{\mu}}_k - \boldsymbol{\mu}_k)\|^2 \rangle_{r(\boldsymbol{\mu}_k)} \\ &= \|\mathbf{x}^{(n)} - \hat{\boldsymbol{\mu}}_k\|^2 + M\hat{\sigma}_k^2 \end{aligned}$$

Variational Bayesian Learning

$$\Psi(x) = \frac{\partial \log \Gamma(x)}{\partial x}$$

and the update rules:

$$\hat{z}_k^{(n)} = \frac{\bar{z}_k^{(n)}}{\sum_{k'=1}^K \bar{z}_{k'}^{(n)}}$$

$$\hat{\alpha}_k = \bar{N}_k + \phi$$

$$\hat{\boldsymbol{\mu}}_k = \frac{\bar{N}_k \bar{\mathbf{x}}_k}{\bar{N}_k + \sigma_0^{-2}}$$

$$\hat{\sigma}_k^2 = \frac{1}{\bar{N}_k + \sigma_0^{-2}}$$

where

$$\bar{z}_k^{(n)} = \exp \left(\Psi(\hat{\alpha}_k) - \frac{1}{2} \left\| \mathbf{x}^{(n)} - \hat{\boldsymbol{\mu}}_k \right\|^2 + M \hat{\sigma}_k^2 + \text{const.} \right)$$

$$\bar{N}_k = \sum_{n=1}^N \hat{z}_k^{(n)}$$

$$\bar{\mathbf{x}}_k = \frac{1}{\bar{N}_k} \sum_{n=1}^N \mathbf{x}^{(n)} \hat{z}_k^{(n)}$$

Variational Bayesian Learning

$$\Psi(x) = \frac{\partial \log \Gamma(x)}{\partial x}$$

and the update rules:

$$\hat{z}_k^{(n)} = \frac{\bar{z}_k^{(n)}}{\sum_{k'=1}^K \bar{z}_{k'}^{(n)}}$$

$$\hat{\alpha}_k = \bar{N}_k + \phi$$

$$\hat{\mu}_k = \frac{\bar{N}_k \bar{x}_k}{\bar{N}_k + \sigma_0^{-2}}$$

$$\hat{\sigma}_k^2 = \frac{1}{\bar{N}_k + \sigma_0^{-2}}$$

Iterative algorithm

Iterate these equations until convergence.

where

$$\bar{z}_k^{(n)} = \exp \left(\Psi(\hat{\alpha}_k) - \frac{1}{2} \left\| \mathbf{x}^{(n)} - \hat{\mu}_k \right\|^2 + M \hat{\sigma}_k^2 + \text{const.} \right)$$

$$\bar{N}_k = \sum_{n=1}^N \hat{z}_k^{(n)}$$

$$\bar{x}_k = \frac{1}{\bar{N}_k} \sum_{n=1}^N \mathbf{x}^{(n)} \hat{z}_k^{(n)}$$

Properties of Variational Bayes

- ✿ Applicable to many models with conditional conjugacy (typically, models built up with well-known distributions).
- ✿ (Generally) faster than sampling methods, and slower than MAP.
- ✿ Provides model selection procedure, and tends to give sparse solution.
- ✿ (Sometimes crude) approximation. Some theoretical guarantee was obtained for special cases.

