

EM Algorithm:

$$X = \{x_i\}$$

$$Y = \{y_i\} \text{. latent}$$

$$\text{Note } q(Y) = \prod_i q_i(y_i)$$

$$q_i(y_i) = P(y_i | x_i, \theta)$$

Log-likelihood :

$$\begin{aligned} l(\theta) &= \log P(X|\theta) = \log \int P(X, Y|\theta) dY \\ &= \log \int \frac{P(X, Y|\theta)}{q(Y)} \cdot q(Y) dY \\ &\geq \int \log \left(\frac{P(X, Y|\theta)}{q(Y)} \right) \cdot q(Y) dY \end{aligned} \quad \begin{matrix} \rightarrow \text{Jensen's Inequality.} \\ \equiv \mathcal{F}(q, \theta) \quad (\text{Appendix}) \end{matrix}$$

Definition : Free Energy

$$\mathcal{F}(q(Y), \theta) = \int \log \left(\frac{P(X, Y|\theta)}{q(Y)} \right) \cdot q(Y) dY = \left\langle \log \frac{P(X, Y|\theta)}{q(Y)} \right\rangle_q$$

$$\text{Hence : } \log P(X|\theta) \geq \mathcal{F}(q(Y), \theta) = \left\langle \log \frac{P(X, Y|\theta)}{q(Y)} \right\rangle_q$$

$$\textcircled{1} \quad \mathcal{F}(q(Y)|\theta) = \left\langle \log \frac{P(X, Y|\theta)}{q(Y)} \right\rangle_q + \underline{\text{H}(q)} \quad \text{The entropy of } q(Y)$$

$$\text{proof: } \mathcal{F}(q(Y)|\theta) = \left\langle \log \frac{P(X, Y|\theta)}{q(Y)} \right\rangle_q = \left\langle \log P(X, Y|\theta) \right\rangle - \left\langle \log q(Y) \right\rangle_q \\ = \text{H}(q(Y)) \\ = \left\langle \log \frac{1}{q(Y)} \right\rangle_q$$

$$\textcircled{2} \quad \mathcal{F}(q(Y)|\theta) = \log P(X|\theta) - \underset{\parallel}{KL} \left[q(Y) \parallel P(Y|X, \theta) \right]$$

$$\text{proof: } \mathcal{F}(q(Y)|\theta) = \left\langle \log \frac{P(X, Y|\theta)}{q(Y)} \right\rangle_q = \int \log \left(\frac{P(Y|X, \theta) \cdot P(X|\theta)}{q(Y)} \right) q(Y) dY \\ = \int \log(P(X|\theta)) \cdot q(Y) \cdot dY + \int \log \left(\frac{P(Y|X, \theta)}{q(Y)} \right) q(Y) dY \\ = \log P(X|\theta) - \underset{\parallel}{KL} \left[q(Y) \parallel P(Y|X, \theta) \right] \\ \geq \log P(X|\theta) \quad \text{as } \underset{\parallel}{KL} \left[q(Y) \parallel P(Y|X, \theta) \right] \geq 0.$$

The lower bound for $\mathcal{F}(q(Y)|\theta)$ is $l(\theta)$ (See Appendix)

EM Steps :

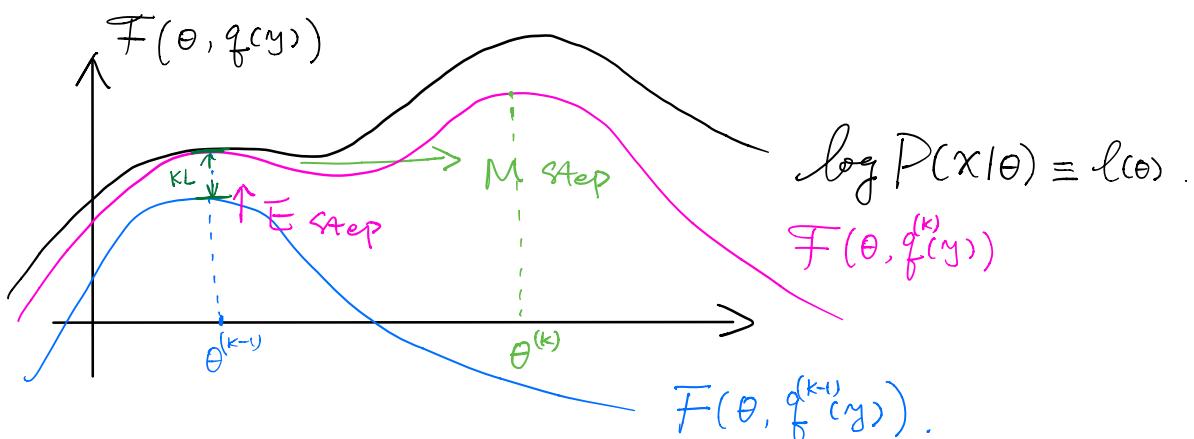
$$\begin{aligned} \text{E-step: } f^{(k)}(\mathbf{y}) &:= \underset{f(\mathbf{y})}{\text{argmax}} \mathcal{F}(q^{(k-1)}(\mathbf{y}), \theta^{(k-1)}) \\ &= l(\theta^{(k-1)}) - \underset{f(\mathbf{y})}{\text{argmin}} \underbrace{\text{KL}\left[q^{(k)}(\mathbf{y}) \parallel P(\mathbf{y}|\mathbf{x}, \theta^{(k-1)})\right]}_{=} \\ &= l(\theta^{(k-1)}) \quad = 0 \text{ by choosing: } q^{(k)}(\mathbf{y}) = P(\mathbf{y}|\mathbf{x}, \theta). \end{aligned}$$

"maximise $\mathcal{F}(q(\mathbf{y}), \theta)$ with θ fixed w.r.t. latent variables \mathbf{y} "

"After E step, $\mathcal{F}(q, \theta) = l(\theta) \Leftrightarrow$ maximum of free energy is maximum of likelihood".

$$\begin{aligned} \text{M-step: } \theta^{(k)} &:= \underset{\theta}{\text{argmax}} \mathcal{F}(q^{(k)}(\mathbf{y}), \theta^{(k-1)}) \quad \begin{array}{l} \text{Ignore, as } H \text{ does not depend} \\ \text{on } \theta \end{array} \\ &= \underset{\theta}{\text{argmax}} \langle \log P(\mathbf{y}, \mathbf{x}|\theta) \rangle_{q^{(k)}} + H(q^{(k)}(\mathbf{y})) \end{aligned}$$

"maximise $\mathcal{F}(q(\mathbf{y}), \theta)$ with latent \mathbf{y} fixed w.r.t. θ "



M Step moves to a larger point at $\mathcal{F}(\theta, q^{(k)}(\mathbf{y}))$

* E and M steps never decreases the log likelihood.

$$\begin{array}{c} l(\theta^{(k-1)}) = \mathcal{F}(q^{(k)}, \theta^{(k-1)}) \leq \mathcal{F}(q^{(k)}, \theta^{(k)}) \leq l(\theta^{(k)}) \\ \downarrow \quad \downarrow \quad \downarrow \\ \text{E-step} \quad \text{M-step} \quad \text{Jensen.} \end{array}$$

Mixture of Gaussian

Data: $\mathcal{X} = \{x_1, \dots, x_n\}$, Latent process: $s_i \stackrel{\text{iid}}{\sim} \text{Disc}(\pi)$.

Component distributions:

$$x_i | s_i = m \sim \mathcal{N}(\mu_m, \Sigma_m)$$

Marginal distribution:

$$P(x_i) = \sum_{m=1}^k \pi_m P_m(x_i; \theta_m).$$

Log-likelihood:

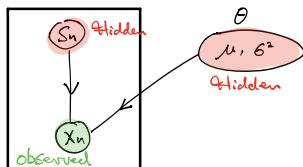
$$\ell(\{\mu_m\}, \{\Sigma_m\}, \pi) = \sum_{i=1}^n \log \sum_{m=1}^k \frac{\pi_m}{\sqrt{2\pi \Sigma_m}} \exp \left\{ -\frac{1}{2} (x_i - \mu_m)^T \Sigma_m^{-1} (x_i - \mu_m) \right\}$$

1. Univariate Case. (Single Gaussian).

Density of a data point x :

$$P(x|\theta) = \sum_{m=1}^k \underbrace{P(x|s_i=m, \theta)}_{\text{Sn: which component generated observation } x_i} \cdot \underbrace{P(s_i=m|\theta)}_{\pi_m} \propto \underbrace{\frac{1}{\sqrt{2\pi \Sigma_m}}}_{\text{Sn}}$$

Sn: which component generated observation x_i



E-step: computes the posterior for s_i given the current parameters:

$$q(s_i) = \underbrace{P(s_i|x_i, \theta)}_{\text{Set this to minimise KL.}} \propto P(x_i|s_i, \theta) \cdot P(s_i|\theta)$$

Set this to minimise KL.

Define Responsibility: $\pi_m \equiv q(s_i=m)$

(which component best explains the data)

Here $\pi_m = q(s_i=m) \propto \frac{\pi_m}{\sqrt{2\pi \Sigma_m}} \exp \left\{ -\frac{1}{2\pi \Sigma_m} (x_i - \mu_m)^2 \right\}$ And check $\sum_m \pi_m = 1$ (normalization)

M Step :

$$\begin{aligned}\Sigma &= \langle \log P(x, s | \theta) \rangle_{q(s)} \\ &= \sum f(s) \log \{P(x|s, \theta) \cdot P(s|\theta)\} \quad \text{Discrete} \\ &= \sum_{i,m} \tau_{im} \left\{ \log \pi_m - \log \sigma_m - \frac{1}{2\sigma_m^2} (x_i - \mu_m)^2 \right\}.\end{aligned}$$

Optimum is found by setting partial derivative of Σ to 0.

$$\begin{aligned}\frac{\partial \Sigma}{\partial \mu_m} &= \sum_i \tau_{im} \frac{(x_i - \mu_m)}{\sigma_m^2} = 0 \\ \frac{\partial \Sigma}{\partial \sigma_m} &= \sum_i \tau_{im} \left(-\frac{1}{\sigma_m} + \frac{(x_i - \mu_m)^2}{\sigma_m^3} \right) = 0 \\ \frac{\partial \Sigma}{\partial \pi_m} &= \sum_i \tau_{im} \frac{1}{\pi_m}, \quad \frac{\partial \Sigma}{\partial \pi_m} + \lambda = 0\end{aligned}$$

Then update
the parameters

$$\begin{aligned}\Rightarrow \mu_m &= \frac{\sum_i \tau_{im} x_i}{\sum_i \tau_{im}} \quad \text{Weighted mean of data} \\ \Rightarrow \sigma_m^2 &= \frac{\sum_i \tau_{im} (x_i - \mu_m)^2}{\sum_i \tau_{im}} \quad \text{Weighted covariance of data} \\ \Rightarrow \pi_m &= \frac{1}{n} \sum_i \tau_{im}\end{aligned}$$

Note that τ_{im} explains how good the cluster m explains the data. If cluster m has good explanation, the corresponding mean or covariance has higher contribution to μ_m and σ_m^2 .

Mixture of Gaussians (Bishop Book)

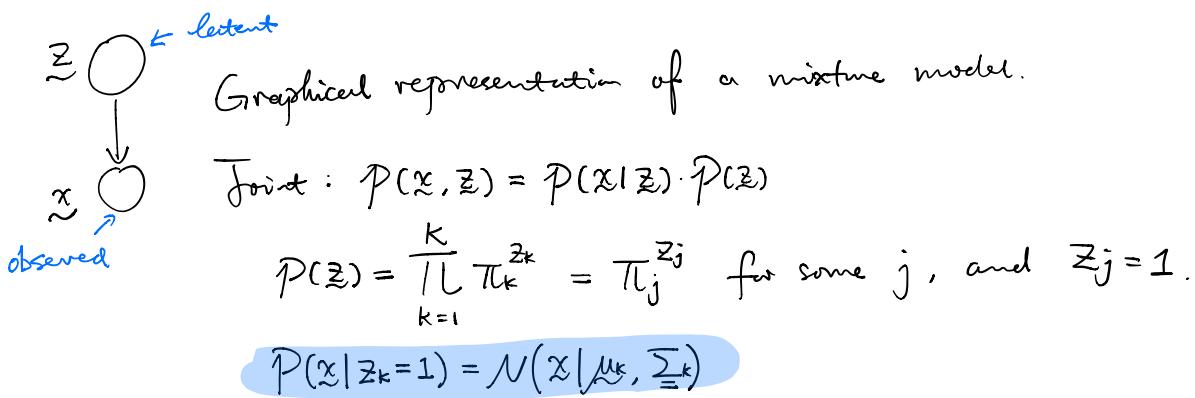
The Gaussian mixture can be written as a superposition of Gaussians in this form:

$$P(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$$

Let $\mathbf{z} \in \mathbb{R}^{Kx1}$, satisfying $z_k \in \{0, 1\}$, and $\sum_k z_k = 1$.

The marginal distribution over \mathbf{z} is specified in terms of π_k

$$\text{s.t. } P(z_k=1) = \pi_k \quad \text{and } \pi_k \in [0, 1], \quad \sum_{k=1}^K \pi_k = 1$$



$$\text{Hence } P(x|z) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$$

$$P(x) = \sum_z P(x, z) = \sum_z P(x|z) \cdot P(z) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$$

Also: $\underline{P(z_k)} \equiv P(z_k=1 | x)$

Probability that it belongs to cluster k
given data x

$$\textcircled{1} \text{ posterior probability} = \frac{P(x|z_k=1) \cdot P(z_k=1)}{P(x)}$$

once we observed x

$$\textcircled{2} \text{ Responsibility that component } k \text{ takes for explaining the observation } x = \frac{P(x|z_k=1) \cdot P(z_k=1)}{\sum_{j=1}^K P(x|z_j=1) \cdot P(z_j=1)}$$

for explaining the observation x .

$$= \frac{N(x | \mu_k, \Sigma_k) \cdot \pi_k}{\sum_{j=1}^K \pi_j N(x | \mu_j, \Sigma_j)}$$

→ regarded as prior probability of $z_k=1$

Maximum likelihood: $\mathbb{D} = \{x_1, \dots, x_n\}$. (CAN solve in closed form)

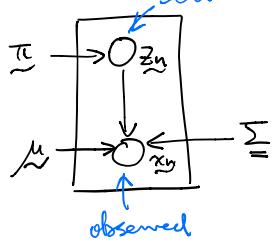
$$\underline{x} = \begin{pmatrix} \cdots & \underline{x}_1^T & \cdots \\ \cdots & \underline{x}_2^T & \cdots \\ \vdots & & \\ \cdots & \underline{x}_n^T & \cdots \end{pmatrix} \in \mathbb{R}^{n \times k}, \quad \underline{z} = \begin{pmatrix} \cdots & \underline{z}_1^T & \cdots \\ \cdots & \underline{z}_2^T & \cdots \\ \vdots & & \\ \cdots & \underline{z}_n^T & \cdots \end{pmatrix} \in \mathbb{R}^{n \times k}$$

\uparrow observed \uparrow latent

$$P(\underline{x} | \pi, \mu, \Sigma) = \prod_{n=1}^N \left(\sum_{k=1}^K \pi_k N(x_n | \underline{\mu}_k, \underline{\Sigma}_k) \right)$$

$$\log P(\underline{x} | \pi, \mu, \Sigma) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K \pi_k N(x_n | \underline{\mu}_k, \underline{\Sigma}_k) \right\}$$

\uparrow latent.



Try normal maximise likelihood: (Doesn't provide a closed form)
Solution for this

$$* \frac{\partial}{\partial \underline{\mu}_k} \log P(\underline{x} | \pi, \mu, \Sigma) = 0 \Rightarrow 0 = - \sum_{n=1}^N \frac{\pi_k N(x_n | \underline{\mu}_k, \underline{\Sigma}_k)}{\sum_j \pi_j N(x_n | \underline{\mu}_j, \underline{\Sigma}_j)} \sum_{n=1}^N (x_n - \underline{\mu}_k)$$

$$\Rightarrow \underline{\mu}_k = \frac{\sum_{n=1}^N \gamma(Z_{nk}) \cdot x_n}{\sum_{n=1}^N \gamma(Z_{nk})} = \frac{1}{N_k} \sum_{n=1}^N \gamma(Z_{nk}) \cdot x_n$$

$$* \frac{\partial}{\partial \underline{\Sigma}_k} \log P(\underline{x} | \pi, \mu, \Sigma) = 0 \quad \text{w.r.t. } \underline{\Sigma}_k$$

$$\Rightarrow \underline{\Sigma}_k = \frac{\sum_{n=1}^N \gamma(Z_{nk}) (x_n - \underline{\mu}_k)(x_n - \underline{\mu}_k)^T}{\sum_{n=1}^N \gamma(Z_{nk})} = \frac{1}{N_k} \sum_{n=1}^N \gamma(Z_{nk}) (x_n - \underline{\mu}_k)(x_n - \underline{\mu}_k)^T$$

$$* \log P(\underline{x} | \pi, \mu, \Sigma) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \quad \text{w.r.t. } \pi_k, \text{ by lagrange multiplier}$$

$$\Rightarrow \pi_k = \frac{\sum_{n=1}^N \gamma(Z_{nk})}{N} = \frac{N_k}{N} \quad \begin{array}{l} \text{The mixing coefficient for the } k^{\text{th}} \text{ component} \\ \text{is given by average responsibility which that} \\ \text{component takes for explaining the data points.} \end{array}$$

EM for Gaussian Mixtures.

Given a Gaussian Mixture model. Goal: maximise the log likelihood w.r.t. parameters

- * mean, covariances of the components μ_k, Σ_k
- * mixing coefficients π_k

EM Steps:

1. Initialize μ_k, Σ_k and π_k , the first log likelihood

2. E step: evaluate the responsibilities:

$$\gamma(z_{nk}) = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)}$$

3. M step: Re-estimate parameters by current responsibilities.

$$\mu_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \cdot x_n,$$

$$\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \cdot (x_n - \mu_k^{new})(x_n - \mu_k^{new})^T$$

$$\pi_k^{new} = \frac{N_k}{N}.$$

$$\text{where } N_k = \sum_{n=1}^N \gamma(z_{nk})$$

4. Evaluate new log-likelihood:

$$\log P(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K \pi_k^{new} N(x_n | \mu_k^{new}, \Sigma_k^{new}) \right\} \quad \dots \quad \text{?}$$

EM for Gaussian Mixture with Latent Variable Known ^{observed}

Consider maximising log likelihood of COMPLETE dataset $\{\underline{x}, \underline{z}\}$.

$$P(\underline{x}, \underline{z} | \mu, \Sigma, \pi) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} N(x_n | \mu_k, \Sigma_k)^{z_{nk}}$$

where z_{nk} is the k^{th} component of \underline{z}_n .

$$\underbrace{\log P(\underline{x}, \underline{z} | \mu, \Sigma, \pi)}_{\text{joint}} = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left\{ \log \pi_k + \log N(x_n | \mu_k, \Sigma_k) \right\}$$

Compare with the  equation, we can find closed form solution here, the same as single Gaussian case.

CAN Find closed form solution

$$P(\underline{z} | \underline{x}, \mu, \Sigma, \pi) \propto \prod_{n=1}^N \prod_{k=1}^K \left\{ \pi_k N(x_n | \mu_k, \Sigma_k) \right\}^{z_{nk}}$$

$E(z_{nk}) = \tau(z_{nk})$ is the responsibility

EM steps :

1. Initialise $\mu^{\text{old}}, \Sigma^{\text{old}}, \pi^{\text{old}}$.

2. E step: Find the responsibility with $\mu^{\text{old}}, \Sigma^{\text{old}}, \pi^{\text{old}}$.

3. M step: update parameters

$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \tau(z_{nk}) \cdot \underline{x}_n ,$$

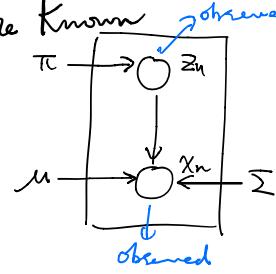
$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \tau(z_{nk}) \cdot (\underline{x}_n - \mu_k^{\text{new}}) (\underline{x}_n - \mu_k^{\text{new}})^T$$

$$\pi_k^{\text{new}} = \frac{N_k}{N}.$$

$$\text{where } N_k = \sum_{n=1}^N \tau(z_{nk})$$

4. Evaluate new log-likelihood:

$$\log P(\underline{x}, \underline{z} | \mu^{\text{new}}, \Sigma^{\text{new}}, \pi^{\text{new}}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left\{ \log \pi_k^{\text{new}} + \log N(x_n | \mu_k^{\text{new}}, \Sigma_k^{\text{new}}) \right\}$$



Gaussian Mixtures and K-means.

K-means: hard assignment of data points to clusters

EM : soft assignment based on posterior probabilities.

K-means:

$$* \text{ Objective function : } J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

represents the sum of squares of the distances of each data point to its assigned vector μ_k .

$$* r_{nk} = \begin{cases} 1 & \text{if } k = \arg\min_j \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise.} \end{cases}$$

$$* \frac{\partial J}{\partial \mu_k} = 0 \Leftrightarrow \underline{\mu}_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$

similarly as $r(z_{nk})$.

but this posterior probability is certain

* Connect Gaussian Mixture with EM with K-means.

Consider a Gaussian Mixture model where $\Sigma = \varepsilon I = \begin{pmatrix} \varepsilon & & \\ & \ddots & \\ & & \varepsilon \end{pmatrix}$

$$P(x|\underline{\mu}, \Sigma) = \frac{1}{\sqrt{2\pi\varepsilon}} \cdot \exp \left\{ -\frac{1}{2\varepsilon} \|x - \underline{\mu}\|^2 \right\}, \text{ with } \varepsilon \text{ fixed.}$$

$$r(z_{nk}) = \frac{\pi_k N(x_n | \underline{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \underline{\mu}_j, \Sigma_j)} = \frac{\pi_k \exp \left\{ -\frac{1}{2\varepsilon} \|x_n - \underline{\mu}_k\|^2 \right\}}{\sum_{j=1}^K \pi_j \exp \left\{ -\frac{1}{2\varepsilon} \|x_n - \underline{\mu}_j\|^2 \right\}}.$$

and $\lim_{\varepsilon \rightarrow 0} r(z_{nk}) = r_{nk}$ from the K-means

let $\varepsilon \rightarrow 0$, consider the log-likelihood:

$$\mathbb{E}_{\Sigma} \left\{ \log P(x, z | \underline{\mu}, \Sigma, \Pi) \right\} \rightarrow -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \underline{\mu}_k\|^2 + \text{constant.}$$

Hence maximizing the expected complete-data log-likelihood is equivalent to minimizing J for K-means.

Factor Analysis

$$\tilde{x} \sim \mathcal{N}(0, \Delta \Delta^T + \Psi), \quad \Theta = \{\Delta, \Psi\}$$

E-step: for each \tilde{x}_n , compute posterior distribution of hidden factors given the observed data.

Set this in E-step

$$\uparrow q(\tilde{y}_n) = P(\tilde{y}_n | \tilde{x}_n, \theta) = \frac{P(\tilde{y}_n, \tilde{x}_n | \theta)}{P(\tilde{x}_n | \theta)}$$

Calculate $P(\tilde{y}_n, \tilde{x}_n | \theta)$ as a function of \tilde{y}_n and \tilde{x}_n fixed.

$$P(\tilde{y}_n, \tilde{x}_n) = P(\tilde{x}_n | \tilde{y}_n) \cdot P(\tilde{y}_n)$$

$$\propto \exp \left\{ -\frac{1}{2} \left[\tilde{y}_n^\top \tilde{y}_n + (\tilde{x}_n - \Delta \tilde{y}_n)^\top \Psi^{-1} (\tilde{x}_n - \Delta \tilde{y}_n) \right] \right\}$$

$$\propto \exp \left\{ -\frac{1}{2} \left[\tilde{y}_n^\top (\mathbb{I} + \Delta^\top \Psi^{-1} \Delta) \tilde{y}_n - 2 \tilde{y}_n^\top \Delta^\top \Psi^{-1} \tilde{x}_n \right] \right\}$$

$$\propto \exp \left\{ -\frac{1}{2} \left[\tilde{y}_n^\top \sum_{\Delta}^{-1} \tilde{y}_n - 2 \tilde{y}_n^\top \sum_{\Delta}^{-1} \tilde{x}_n + \tilde{x}_n^\top \sum_{\Delta}^{-1} \tilde{x}_n \right] \right\}$$

$$\text{Hence } \sum_{\Delta} = (\mathbb{I} + \Delta^\top \Psi^{-1} \Delta)^{-1}, \quad \mu_n = \sum_{\Delta} \Delta^\top \Psi^{-1} \tilde{x}_n$$

μ_n is a linear function of \tilde{x}_n

\sum_{Δ} does not depend \tilde{x}_n .

M-step: find $\theta^{(t+1)}$ by maximizing

$$\theta^{(t+1)} = \underset{\theta}{\operatorname{argmax}} \mathcal{F}(q(\tilde{y}), \theta) = \underset{\theta}{\operatorname{argmax}} \langle \log P(X, Y | \theta) \rangle + \text{H}$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_n \langle \log P(y_n | \theta) + \log P(x_n | y_n, \theta) \rangle_{q_n(y_n)} + \text{H}$$

$$\log P(y_n | \theta) + \log P(x_n | y_n, \theta)$$

$$= C - \frac{1}{2} \tilde{y}_n^\top \tilde{y}_n - \frac{1}{2} \log |\Psi| - \frac{1}{2} (\tilde{x}_n - \Delta \tilde{y}_n)^\top \Psi^{-1} (\tilde{x}_n - \Delta \tilde{y}_n)$$

$$= C' - \frac{1}{2} \log |\Psi| - \frac{1}{2} \left[\tilde{x}_n^\top \Psi^{-1} \tilde{x}_n - 2 \tilde{x}_n^\top \Psi^{-1} \Delta \tilde{y}_n + \tilde{y}_n^\top \Delta^\top \Psi^{-1} \Delta \tilde{y}_n \right]$$

$$= C' - \frac{1}{2} \log |\Psi| - \frac{1}{2} \left[\tilde{x}_n^\top \Psi^{-1} \tilde{x}_n - 2 \tilde{x}_n^\top \Psi^{-1} \Delta \tilde{y}_n + \text{Tr}(\Delta^\top \Psi^{-1} \Delta \tilde{y}_n \tilde{y}_n^\top) \right]$$

mean

Variance

Take expectation w.r.t. $q_n(y_n)$ Sufficient statistics of multivariate

$$\langle \log P(y_n|\theta) + \log P(x_n|y_n, \theta) \rangle_{q_n(y_n)}$$

$$= C' - \frac{1}{2} \log |\Psi| - \frac{1}{2} \left[\underline{x}_n^T \underline{\Psi}^{-1} \underline{x}_n - 2 \underline{x}_n^T \underline{\Psi}^{-1} \underline{\Delta} \underline{\mu}_n + \text{Tr}(\underline{\Delta} \underline{\Psi}^{-1} \underline{\Delta} (\underline{\mu}_n \underline{\mu}_n^T + \underline{\Sigma})) \right]$$

$$\bar{F} = \langle \log P(y_n|\theta) + \log P(x_n|y_n, \theta) \rangle_{q_n(y_n)} + A$$

$$= C' - \frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \left[\underline{x}_n^T \underline{\Psi}^{-1} \underline{x}_n - 2 \underline{x}_n^T \underline{\Psi}^{-1} \underline{\Delta} \underline{\mu}_n + \text{Tr}(\underline{\Delta} \underline{\Psi}^{-1} \underline{\Delta} (\underline{\mu}_n \underline{\mu}_n^T + \underline{\Sigma})) \right] \in \mathbb{R}$$

Partial derivatives : Note that $\frac{\partial}{\partial B} \text{Tr}[AB] = A^T$, $\frac{\partial}{\partial A} \log |A| = A^{-T}$

$$\frac{\partial \bar{F}}{\partial \underline{\Delta}} = \underline{\Psi}^{-1} \sum_n \underline{x}_n \underline{\mu}_n^T - \underline{\Psi}^{-1} \underline{\Delta} \left(N \underline{\Sigma} + \sum_n \underline{\mu}_n \underline{\mu}_n^T \right) = 0.$$

$$\Rightarrow \hat{\underline{\Delta}} = \left(\sum_n \underline{x}_n \underline{\mu}_n^T \right) \cdot \left(N \underline{\Sigma} + \sum_n \underline{\mu}_n \underline{\mu}_n^T \right)^{-1}$$

$$\frac{\partial \bar{F}}{\partial \underline{\Psi}^{-1}} = \frac{N}{2} \underline{\Psi} - \frac{1}{2} \sum_n \left[\underline{x}_n \underline{x}_n^T - \underline{\Delta} \underline{\mu}_n \underline{x}_n^T - \underline{x}_n \underline{\mu}_n^T \underline{\Delta}^T + \underline{\Delta} (\underline{\mu}_n \underline{\mu}_n^T + \underline{\Sigma}) \underline{\Delta}^T \right] = 0$$

$$\Rightarrow \hat{\underline{\Psi}} = \frac{1}{N} \sum_n \left[\underline{x}_n \underline{x}_n^T - \underline{\Delta} \underline{\mu}_n \underline{x}_n^T - \underline{x}_n \underline{\mu}_n^T \underline{\Delta}^T + \underline{\Delta} (\underline{\mu}_n \underline{\mu}_n^T + \underline{\Sigma}) \underline{\Delta}^T \right] \\ = \underbrace{\underline{\Delta} \underline{\Sigma} \underline{\Delta}^T}_{D \times D} + \frac{1}{N} \cdot \sum_n (\underline{x}_n - \underline{\Delta} \underline{\mu}_n) \cdot (\underline{x}_n - \underline{\Delta} \underline{\mu}_n)^T$$

Note that if $\underline{\Sigma} \rightarrow 0$, $\hat{\underline{\Psi}}$ becomes the likelihood of ML linear regression.

Mixture of FA :

$$P(\underline{x}|\theta) = \sum_k \pi_k N(\underline{\mu}_k, \underline{\Delta}_k \underline{\Delta}_k^T + \underline{\Psi}), \quad \left\{ \begin{array}{l} \pi_k : \text{mixing proportion for FA.} \\ \underline{\mu}_k : \text{centre} \\ \theta = \{\{\pi_k, \underline{\mu}_k, \underline{\Delta}_k\}_{k=1,\dots,K}, \underline{\Psi}\}. \end{array} \right.$$

Details see Hinton + Zouhri's Paper.

EM for Exponential Families

For joint over $\underline{z} = (\underline{y}, \underline{x})$ with exponential form:

$$P(\underline{z}|\theta) = f(\underline{z}) \exp\{\theta^T \cdot T(\underline{z})\} / Z(\theta).$$

The free energy dependence on θ is given by:

$$\begin{aligned} F(f, \theta) &= \langle \log P(\underline{y}, \underline{x}|\theta) \rangle_f + H(f) \\ &= \int f(\underline{y}) \log \underbrace{P(\underline{y}, \underline{x}|\theta)}_{Z(\theta)} d\underline{y} + H(f) \\ &= \int f(\underline{y}) \cdot \log \left[f(\underline{z}) \exp\{\theta^T \cdot T(\underline{z})\} / Z(\theta) \right] d\underline{y} + H(f) \\ &= \int f(\underline{y}) \cdot \left[\theta^T \cdot T(\underline{z}) - \log Z(\theta) \right] d\underline{y} + \text{constant} \\ &= \theta^T \langle T(\underline{z}) \rangle_{f(\underline{y})} - \log Z(\theta) + \text{constant}. \end{aligned}$$

E-step: Compute the expected sufficient statistics under f :

$$\text{i.e. } \frac{\partial}{\partial \theta} \log Z(\theta) = \langle T(\underline{z}) | \theta \rangle$$

$$\text{M-step: } \frac{\partial F}{\partial \theta} = \langle T(\underline{z}) \rangle_{f(\underline{y})} - \frac{\partial}{\partial \theta} \log Z(\theta)$$

EM for exponential families Mixtures

note $s_i = m \Leftrightarrow \underline{s}_i = (0, 0, \dots, 0, \underbrace{1}_{m^{\text{th position}}}, 0, \dots, 0)$

Let M component distributions' parameters into $\underline{\Theta} = [\theta_m]$.

$$\log P(\underline{x}, \underline{s}) = \sum_i \left\{ (\log \pi_i)^T \underline{s}_i + \underline{s}_i^T \underline{\Theta}^T T(x_i) - \underline{s}_i^T \log Z(\underline{\Theta}) \right\} + \text{constant}.$$

where $\log Z(\underline{\Theta})$ collects log-normalisers M components into an M -element vector.

E-step : Expected sufficient statistics

$$* \sum_i \langle S_{ij} \rangle_f \quad \text{responsibilities}$$

$$* \sum_i T(x_i) \langle S_i^T \rangle_f \quad \text{responsibility-weighted sufficient stats.}$$

M-step : maximisation of expected log-joint

$$\pi^{(k+1)} \propto \sum_i \langle S_{ij} \rangle_f.$$

$$\langle T(x) | \theta_m^{(k+1)} \rangle = \frac{\sum_i T(x_i) \langle [S_i]_m \rangle_f}{\sum_i \langle [S_i]_m \rangle_f}$$

EM for MAP

$$P(z|\theta) = \frac{f(z) \cdot \exp\{\theta^T T(z)\}}{Z(\theta)}, \quad P(\theta) = \frac{F(\nu, \tau) \cdot \exp\{\theta^T \tau\}}{Z(\theta)}$$

$$\begin{aligned} F_{MAP}(q, \theta) &= \int q(y) \log P(y, z, \theta) dy + H(q) \leq \log P(z|\theta) + \log P(\theta). \\ &= \int q(y) \left\{ \theta^T \left(\sum_i T(z_i) + \tau \right) - (\nu + \nu) \log Z(\theta) \right\} dy + \text{const} \\ &= \theta^T \left(\langle T(z) \rangle_{q(y)} + \tau \right) - (\nu + \nu) \log Z(\theta) + \text{const}. \end{aligned}$$

Hence expected sufficient stats in E step are unchanged.

Appendix :

* Positive Definiteness of KL.

Claim : $KL[q(x) \| p(x)] \geq 0$ with equality iff $\forall x : p(x) = q(x)$.

$$\text{Proof: } KL(q \| p) = \sum_i q_i \log \frac{q_i}{p_i}$$

$$E \stackrel{\text{def}}{=} KL(q \| p) + \lambda(1 - \sum_i q_i) = \sum_i q_i \log \frac{q_i}{p_i} + \lambda(1 - \sum_i q_i)$$

$$\frac{\partial E}{\partial q_i} = \log q_i - \log p_i + 1 - \lambda = 0 \Rightarrow q_i = p_i \exp(\lambda - 1).$$

$$\frac{\partial E}{\partial \lambda} = 1 - \sum_i q_i = 0 \Rightarrow \sum_i q_i = 1.$$

Check curvature (Hessian) :

$$\frac{\partial^2 E}{\partial q_i \partial q_j} = \frac{1}{q_i} > 0; \quad \frac{\partial^2 E}{\partial q_i \partial q_j} = 0.$$

Hence unique stationary point $q_i = p_i$ is a minimum. \square

* Fixed Points of EM are Stationary Points in ℓ

Let a fixed point of EM with θ^* , then:

$$\frac{\partial}{\partial \theta} \left\langle \log P(Y, X | \theta) \right\rangle_{P(Y|X, \theta^*)} \Big|_{\theta^*} = 0$$

$$\ell(\theta) = \log P(X | \theta) = \left\langle \log P(X | \theta) \right\rangle_{P(Y|X, \theta^*)}$$

$$= \left\langle \log \frac{P(Y, X | \theta)}{P(Y | X, \theta)} \right\rangle_{P(Y|X, \theta^*)}.$$

$$= \left\langle \log P(Y, X | \theta) \right\rangle_{P(Y|X, \theta^*)} - \left\langle \log P(Y | X, \theta) \right\rangle_{P(Y|X, \theta^*)}$$

$$\text{Hence } \frac{d}{d\theta} \ell(\theta) = \frac{d}{d\theta} \left\langle \log P(Y, X | \theta) \right\rangle_{P(Y|X, \theta^*)} - \underbrace{\frac{d}{d\theta} \left\langle \log P(Y | X, \theta) \right\rangle_{P(Y|X, \theta^*)}}_{= 0 \text{ if derivative exists.}}$$

$$\Rightarrow \frac{d}{d\theta} \ell(\theta) \Big|_{\theta^*} = \frac{d}{d\theta} \left\langle \log P(Y, X | \theta) \right\rangle_{P(Y|X, \theta^*)} \Big|_{\theta^*} = 0. \quad (\min \text{ of } KL)$$

This term is not necessarily KL. But we can add some extra terms to make it KL.

$$\begin{aligned}
 & \frac{d}{d\theta} \text{KL}\left[P(\underline{y}|\underline{x}, \theta^*) \| P(\underline{y}|\underline{x}, \theta)\right] \\
 &= \frac{d}{d\theta} \left\{ -\langle \log P(\underline{y}|\underline{x}, \theta) \rangle_{P(\underline{y}|\underline{x}, \theta^*)} + \langle \log P(\underline{y}|\underline{x}, \theta^*) \rangle_{P(\underline{y}|\underline{x}, \theta^*)} \right\}.
 \end{aligned}$$