

## Expectation-Maximization

notes from notes from Stanford Machine Learning,  
CS229\_Lecture\_Notes.pdf

**Factor Analysis model** is an example of EM for missing data, but m is << number of features n

$\{x^{(1)}, \dots, x^{(m)}\}$  training set

Observed, incomplete-data of length n w/ d number of features

$$x^{(i)} \in \mathbb{R}^n$$

$$z \sim \mathcal{N}(0, I)$$

$$\epsilon \sim \mathcal{N}(0, \Psi)$$

$$x = \mu + \Lambda z + \epsilon.$$

$\epsilon$  and  $z$  are independent.

$\mu$ ,  $\Lambda$ , and  $\epsilon$  are the model parameters.  
 $k$  is a hyper-parameter

$$\begin{bmatrix} z \\ x \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \vec{0} \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda\Lambda^T + \Psi \end{bmatrix}\right). \quad (3)$$

Hence, we also see that the marginal distribution of  $x$  is given by  $x \sim \mathcal{N}(\mu, \Lambda\Lambda^T + \Psi)$ . Thus, given a training set  $\{x^{(i)}; i = 1, \dots, m\}$ , we can write down the log likelihood of the parameters:

$$\underline{\ell(\mu, \Lambda, \Psi)} = \log \prod_{i=1}^m \frac{1}{(2\pi)^{n/2} |\Lambda\Lambda^T + \Psi|^{1/2}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu)^T (\Lambda\Lambda^T + \Psi)^{-1} (x^{(i)} - \mu)\right).$$

E-step is easy.

we find that  $z^{(i)}|x^{(i)}; \mu, \Lambda, \Psi \sim \mathcal{N}(\mu_{z^{(i)}|x^{(i)}}, \Sigma_{z^{(i)}|x^{(i)}})$ , where

$$\begin{aligned} \mu_{z^{(i)}|x^{(i)}} &= \Lambda^T (\Lambda\Lambda^T + \Psi)^{-1} (x^{(i)} - \mu), \\ \Sigma_{z^{(i)}|x^{(i)}} &= I - \Lambda^T (\Lambda\Lambda^T + \Psi)^{-1} \Lambda. \end{aligned}$$

$$Q_i(z^{(i)}) = \frac{1}{(2\pi)^{k/2} |\Sigma_{z^{(i)}|x^{(i)}}|^{1/2}} \exp\left(-\frac{1}{2}(z^{(i)} - \mu_{z^{(i)}|x^{(i)}})^T \Sigma_{z^{(i)}|x^{(i)}}^{-1} (z^{(i)} - \mu_{z^{(i)}|x^{(i)}})\right).$$

Let's now work out the M-step. Here, we need to maximize

$$\sum_{i=1}^m \int_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \mu, \Lambda, \Psi)}{Q_i(z^{(i)})} dz^{(i)} \quad (4)$$

with respect to the parameters  $\mu, \Lambda, \Psi$ . We will work out only the optimiza-

$$\Lambda = \left( \sum_{i=1}^m (x^{(i)} - \mu) \mu_{z^{(i)}|x^{(i)}}^T \right) \left( \sum_{i=1}^m \mu_{z^{(i)}|x^{(i)}} \mu_{z^{(i)}|x^{(i)}}^T + \Sigma_{z^{(i)}|x^{(i)}} \right)^{-1}. \quad (8)$$

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}. \quad \text{this can be calculated just once and needs not be further updated as the algorithm is run}$$

and setting  $\Psi_{ii} = \Phi_{ii}$  (i.e., letting  $\Psi$  be the diagonal matrix containing only the diagonal entries of  $\Phi$ ).

$$\Phi = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} - x^{(i)} \mu_{z^{(i)}|x^{(i)}}^T \Lambda^T - \Lambda \mu_{z^{(i)}|x^{(i)}} x^{(i)T} + \Lambda (\mu_{z^{(i)}|x^{(i)}} \mu_{z^{(i)}|x^{(i)}}^T + \Sigma_{z^{(i)}|x^{(i)}}) \Lambda^T,$$