

# Introduction to Machine Learning. CSCI-UA 9473, Lecture 9.

Augustin Cosse

Ecole Normale Supérieure, DMA & NYU  
Fondation Sciences Mathématiques de Paris.



2018

# Statistical intuition and Latent Variable models

- ▶ Many models discussed so far can in fact be considered as particular instances of latent variable models.
- ▶ The general factorization for a latent variable model of the form  $\mathbf{z}_i \rightarrow \mathbf{x}_i$  is  $p(\mathbf{x}_i|\mathbf{z}_i)p(\mathbf{z}_i)$ . Recall that here  $p(\mathbf{x}_i|\mathbf{z}_i)$  denotes the likelihood and the prior  $p(\mathbf{z}_i)$  indicates the distribution of the latent state

$p(\mathbf{x}_i \mathbf{z}_i)$	$p(\mathbf{z}_i)$	Model
prod. Gauss.	prod. Gauss.	Factor Analysis/proba. PCA
MVN	cat.	Mixture of Gaussians
prod. cat	cat.	Mixture of Multinomials
prod. Gauss.	prod. Laplace	proba. ICA/sparse coding
prod. cat.	prod. Gauss.	multinomial PCA
prod. cat.	Dirichlet	Latent Dirichlet alloc.

# Statistical intuition and Latent Variable models

- ▶ The simplest latent variable models assume that the **data** was **generated** from a distribution governed by a **single latent state**  $z_i$ . In a Bayesian approach, we let  $p(z_i)$  to denote the prior for this latent state.
- ▶ We further let  $p(\mathbf{x}_i|z_i = k)$  to denote the (assumed) distribution of the prototypes given the latent state  $z_i = k$ .
- ▶ A popular choice for the prior is the **categorical distribution**.

# Statistical intuition and Latent Variable models

- ▶ The **categorical distribution** is usually used when we represent measurements as belonging to 1 of  $K$  possible exclusive classes.
- ▶ Clusters are usually represented through **dummy encodings**.
- ▶ The general form of a dummy encoding is  $\mathbf{x} = [\mathbb{I}(x = 1), \dots, \mathbb{I}(x = K)]$  where  $\mathbb{I}(x = \alpha) = 1$  if  $x$  is in the cluster  $C_\alpha$ .
- ▶ Then if we assume that the class probabilities are independent, and label those probabilities as  $\theta_j$ , the **categorical distribution** reads as a particular case of the multinomial

$$\text{Cat}(\mathbf{x}|\mathbf{1}, \theta) = \prod_{j=1}^K \theta_j^{\mathbb{I}(x_j=1)}$$

## Statistical intuition and Latent Variable models

- ▶ More Generally, we will use the notation  $\pi_k$  to denote the class probabilities (probabilities to belong to the class  $C_k$ ).
- ▶ With this notation, we thus have  $p(z_i) = \text{Cat}(\boldsymbol{\pi})$  and in particular, and  $p(z_i = k) = \pi_k$
- ▶ Given those probabilities, the probability to observe a prototype  $\mathbf{x}_i$  from the dataset is given by

$$p(\mathbf{x}_i | \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}_i | \boldsymbol{\theta})$$

Here  $p_k(\mathbf{x}_i | \boldsymbol{\theta})$  denote the  $k^{\text{th}}$  **base distribution** (i.e distribution to observe the prototype  $\mathbf{x}_i$  knowing it belongs to the  $k^{\text{th}}$  cluster.). The general model above is known as a **Mixture Model**.

# Statistical intuition and Latent Variable models

- ▶ The two most popular models are the **Mixture of Gaussians**

$$p(\mathbf{x}_i | \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- ▶ and the **mixture of Multinoullis**,

$$p(\mathbf{x}_i, z_i = k | \boldsymbol{\theta}) = \prod_{j=1}^D \text{Ber}(x_{ij} | \mu_{jk}) = \prod_{j=1}^D \mu_{jk}^{x_{ij}} (1 - \mu_{jk})^{1-x_{ij}}$$

- ▶ The Mixture of Multinoullis is **useful** when using **dummy encodings**, i.e. when  $\mathbf{x}_i = (0, 1, 0, \dots, 0)$  and we use  $\mu_{jk}$  to denote the probability that the  $j^{\text{th}}$  bit is 1 in sequences from cluster  $k$ .

# The EM algorithm

- ▶ When we want to fit a probabilistic model to the data, we usually minimize the negative log-likelihood (we look for the parameters that make it maximally likely to observe the given data)
- ▶ In the case of a GMM, the log likelihood reads as

$$\log p(\mathbf{x}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \log \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

- ▶ Minimizing this expression directly is hard because the log cannot be pushed inside the sum

# The EM algorithm

- ▶ The general form of a probability distribution from the exponential family reads as

$$p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp(\boldsymbol{\theta}^T \phi(\mathbf{x}, \boldsymbol{\theta}))$$

- ▶ And the parameters we want to find also appear in the normalizing constant. I.e for a distribution from the exponential family, the **observed log-likelihood** reads as

$$\begin{aligned} \ell(\boldsymbol{\theta}) &= \sum_i \log \sum_{\mathbf{z}_i} p(\mathbf{x}_i, \mathbf{z}_i|\boldsymbol{\theta}) \\ &= \sum_i \log \left[ \sum_{\mathbf{z}_i} e^{\boldsymbol{\theta}^T \phi(\mathbf{z}_i, \mathbf{x}_i)} \right] - N \log Z(\boldsymbol{\theta}) \end{aligned}$$

- ▶ The log sum exp is convex and the normalizing constant  $Z(\boldsymbol{\theta})$  is convex as well. However, the difference of two convex function



# The EM algorithm

- ▶ An easier approach would be to optimize the **complete data log-likelihood**

$$\ell_c(\boldsymbol{\theta}) = \sum_i \log p(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta}) = \boldsymbol{\theta}^T \left( \sum_i \boldsymbol{\phi}(\mathbf{x}_i, \mathbf{z}_i) \right) - NZ(\boldsymbol{\theta})$$

- ▶ In the exponential family the normalizing constant is convex so that the whole function is concave and can be optimized efficiently

# The EM algorithm

- ▶ For the reasons listed above, we would clearly prefer to work with the complete data log-likelihood

$$\ell_c(\theta) = \sum_{i=1}^N \log p(\mathbf{x}_i, \mathbf{z}_i | \theta)$$

- ▶ The problem is that we don't have access to the joint probability  $p(\mathbf{x}_i, \mathbf{z}_i | \theta)$
- ▶ To get round the difficulty, and to estimate the parameters of the mixture together with the latent states, the EM algorithm works on the expected complete data log-likelihood

$$Q(\theta, \theta^{t-1}) = \mathbb{E} \{ \ell_c(\theta) | \mathcal{D}, \theta^{t-1} \}$$

Here  $Q(\theta, \theta^{t-1})$  is called the **auxiliary function**.

# The EM algorithm

- ▶ For the reasons listed above, **we would** clearly **prefer** to work with the **complete** data **log-likelihood**

$$\ell_c(\theta) = \sum_{i=1}^N \log p(\mathbf{x}_i, \mathbf{z}_i | \theta)$$

- ▶ To get round the difficulty, and to estimate the parameters of the mixture together with the latent states, the EM algorithm works on the **expected** complete data log-likelihood

$$Q(\theta, \theta^{t-1}) = \mathbb{E} \{ \ell_c(\theta) | \mathcal{D}, \theta^{t-1} \}$$

Here  $Q(\theta, \theta^{t-1})$  is called the **auxilliary function**.

- ▶ The *E*-step computes the expression of  $Q$  (or the terms needed to express  $Q$ ). The *M*-step optimizes  $Q$  with respect to  $\theta$ .

# The Auxilliary function

- ▶ The auxilliary function reads as

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{t-1}) &= \mathbb{E} \left\{ \sum_i \log p(\mathbf{x}_i, z_i | \boldsymbol{\theta}) \right\} \\ &= \sum_i \mathbb{E} \left\{ \log \left[ \prod_{k=1}^K (\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k))^{\mathbb{I}(z_i=k)} \right] \right\} \\ &= \sum_i \sum_k \mathbb{E} \{ \mathbb{I}(z_i = k) \} \log(\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k)) \\ &= \sum_i \sum_k p(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}^{t-1}) \log(\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}_k)) \\ &= \sum_i \sum_k r_{ik} \log \pi_k + \sum_i \sum_k r_{ik} \log p(\mathbf{x}_i | \boldsymbol{\theta}_k) \end{aligned}$$

- ▶ It is fully determined from the **responsibilities**  $r_{ik}$  of the cluster  $k$  in the realization of the prototype  $\mathbf{x}_i$  as well as the likelihoods  $p(\mathbf{x}_i | \boldsymbol{\theta}_k)$  which follow from gaussianity and  $\theta$

# The EM algorithm for GMMs

- ▶ Given the likelihoods, the **E-step** updates the **responsibilities** as

$$r_{ik} = \frac{\pi_k p(\mathbf{x}_i | \boldsymbol{\theta}^{t-1})}{\sum_{k'} \pi_{k'} p(\mathbf{x}_i | \boldsymbol{\theta}_{k'}^{t-1})}$$

- ▶ The **M-step** then optimizes  $Q$  with respect to the class probabilities  $\pi_k$  and the parameters  $\boldsymbol{\theta}_k$
- ▶ The estimates for  $\pi_k$  are given by

$$\pi_k = \frac{1}{N} \sum_i r_{ik} = \frac{r_k}{N}$$

- ▶ The estimates for  $\boldsymbol{\theta} = (\mu_k, \sigma_k)$  are obtained by substituting the Normal distribution for  $p(\mathbf{x}_i, \boldsymbol{\theta}_k)$  in the log likelihood and minimizing

## The EM algorithm for GMMs (E-step)

- ▶ Substituting the normal distributions for the  $p(\mathbf{x}_i|\theta_k)$ , we get

$$\begin{aligned}\ell(\mathbf{m}_k, \mathbf{\Sigma}_k) &= \sum_k \sum_i r_{ik} \log p(\mathbf{x}_i|\theta_k) \\ &= -\frac{1}{2} \sum_i r_{ik} \left[ \log |\mathbf{\Sigma}_k| + (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \mathbf{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \right]\end{aligned}$$

- ▶ Setting the derivatives to zero, we get

$$\begin{aligned}\boldsymbol{\mu}_k &= \frac{\sum_i r_{ik} \mathbf{x}_i}{r_k} \\ \boldsymbol{\sigma}_k &= \frac{\sum_i r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T}{r_k} = \frac{\sum_i r_{ik} \mathbf{x}_i \mathbf{x}_i^T}{r_k} - \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T\end{aligned}$$

## Relation to K-means

- ▶ Remember the K means algorithm? **K means** can be considered as a **particular instance** of the EM algorithm. If we assume that  $\Sigma_k = \sigma^2 \mathbf{I}_D$  and  $\pi_k = \frac{1}{K}$  is fixed, we **only update** the **centers of the clusters**
- ▶ Instead of the previous responsibilities, one can assume that the probability of a prototype belonging to a cluster is either 1 or 0.
- ▶ We can choose the only possible cluster for a prototype to be the one that maximizes the likelihood

$$z_i^* = \underset{k}{\operatorname{argmax}} p(z_i = k | \mathbf{x}_i, \theta)$$

- ▶ We then set the probability that  $\mathbf{x}_i$  belongs to this cluster to 1. We can do this because we assumed  $\pi_k = 1/K$  fixed and  $\Sigma_k = \sigma^2 \mathbf{I}$

- ▶ Under the earlier hypotheses, maximizing the posterior to find the most likely assignment reduces to the minimization

$$z_i^* = \underset{k}{\operatorname{argmin}} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|$$

- ▶ And the M-step updates the centers as

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i|z_i=k} \mathbf{x}_i$$



# K-means vs EM

- ▶ K means
  - ▶ + Better running time
  - ▶ + More interesting for high dimensional data
  - ▶ + Interpretation is easier
  - ▶ - Assumes clusters are spherical (see Mouse dataset). So does not work well with complex shape.
  - ▶ - The "Hard assignment" approach might lead to misclassification
- ▶ EM Clustering
  - ▶ + Works usually better when there is some uncertainty regarding the assignment
  - ▶ + Does not assume any predefined geometry for the clusters
  - ▶ - Uses more information than  $K$ -means so more difficult to implement in high dimension.
  - ▶ - More difficult to interpret

# Statistical intuition for FA, PCA and ICA

- ▶ Now that we have introduced GMMs and the notion of latent variable model, we are ready to discuss the [statistical intuition](#) for [Factor analysis](#), [PCA](#) and [ICA models](#).
- ▶ [Gaussian mixture models](#) are [very general](#) in that every observation is assumed to have been generated from one of  $k$  independent clusters with their respective mean and covariance.
- ▶ An [alternative](#) is to [view](#) the [distribution of prototypes](#) as something [smoother](#) and to replace the hard assignment (i.e. each of the prototype belongs (exclusively) to one of the  $k$  clusters) by the assumption that the prototypes are organized according to a set of gaussian distributions concentrated around a single point.

# Statistical intuition for FA, PCA and ICA

- ▶ Such model then relies on a first **continuous prior** for the latent variables  $z_i$ ,

$$p(z_i) = \mathcal{N}(z_i | \mu_0, \Sigma_0)$$

An then model each of the **prototypes distributions** as **gaussian distributions** centered around a mean defined from this continuously varying latent variable  $z_i$

$$p(\mathbf{x}_i | z_i, \theta) = \mathcal{N}(\mathbf{W}z_i + \mu, \Psi)$$

In this case,  $\theta = (\mathbf{W}, \mu, \Psi)$

- ▶ **Factor analysis** can thus be understood as a **GMM** with **constraints** on the **mean** and **covariance** but **continuous latent variables**

# Statistical intuition for FA, PCA and ICA

$$p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

- ▶ Here  $\mathbf{W}$  is known as the **factor loading matrix** and  $\boldsymbol{\Psi}$  is the covariance matrix.
- ▶ In practice, the **covariance matrix**  $\boldsymbol{\Psi}$  is **usually** taken to be **diagonal** and we therefore turn to the  $\mathbf{z}_i$  and their connection through the latent distribution  $\mathcal{N}(\mathbf{z}_i | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$
- ▶ When the **covariance** is taken to be **spherical**, i.e.  $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$ , we get the **probabilistic PCA** model as we will see.

## Low rank covariance

- ▶ The model  $p(\mathbf{x}_i|\boldsymbol{\theta})$  given by the combination of the prior  $p(\mathbf{z}_i)$  and the likelihood  $p(\mathbf{x}_i|\mathbf{z}_i, \boldsymbol{\theta})$  is known as a **linear gaussian system**.
- ▶ For a general Gaussian system

$$\begin{cases} p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x) \\ p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \boldsymbol{\Sigma}_y) \end{cases}$$

The distribution of  $\mathbf{y}$  (a.k.a normalizing constant) is defined as

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu}_x + \mathbf{b}, \boldsymbol{\Sigma}_y + \mathbf{A}\boldsymbol{\Sigma}_x\mathbf{A}^T)$$

- ▶ Using this to derive the distribution of prototypes in FA, we get the distribution

$$\begin{aligned} p(\mathbf{x}_i|\boldsymbol{\theta}) &= \int \mathcal{N}(\mathbf{x}_i|\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi})\mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) d\mathbf{z}_i \\ &= \mathcal{N}(\mathbf{x}_i|\mathbf{W}\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\boldsymbol{\Sigma}_0\mathbf{W}^T) \end{aligned}$$

## FA as a low rank model for the covariance

- ▶ The **mean** and **covariance** of the factor analysis model thus read as

$$\begin{aligned}\mathbb{E}\{\mathbf{x}_i\} &= \mathbf{W}\boldsymbol{\mu}_0 + \boldsymbol{\mu}, & \text{cov}\{\mathbf{x}_i\} &= \mathbf{W}\mathbb{E}\{\mathbf{z}\mathbf{z}^T\}\mathbf{W}^T + \boldsymbol{\Psi} \\ & & &= \mathbf{W}\boldsymbol{\Sigma}_0\mathbf{W}^T + \boldsymbol{\Psi}\end{aligned}$$

- ▶ An **additional interpretation** of the factor analysis model can be obtained by noting that one can always write  $\boldsymbol{\mu}' = \boldsymbol{\mu} + \mathbf{W}\boldsymbol{\mu}_0$  and take  $\boldsymbol{\mu}'_0 = 0$ . We can also always take  $\boldsymbol{\Sigma}_0 = \mathbf{I}$  as we can always write the model by introducing the factor  $\tilde{\mathbf{W}} = \mathbf{W}\boldsymbol{\Sigma}_0^{-1/2}$ ,

$$\text{cov}(\mathbf{x}_i) = \mathbf{W}\boldsymbol{\Sigma}_0\mathbf{W}^T + \boldsymbol{\Psi} = \tilde{\mathbf{W}}\tilde{\mathbf{W}}^T + \boldsymbol{\Psi} \quad (1)$$

- ▶ Factor analysis can thus be understood as a **Gaussian model** with a **low rank covariance** matrix !

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \boldsymbol{\Psi})$$

# Unidentifiability

- ▶ Learning a factor analysis model is **ill posed** in general (i.e we say that the parameters of the FA model are **unidentifiable**)
- ▶ Replacing the factor  $\mathbf{W}$  by any other matrix of the form  $\mathbf{WR}$  where  $\mathbf{R}$  is an **orthogonal matrix**  $\mathbf{RR}^T = \mathbf{I}$ , we get

$$\text{cov}[\mathbf{x}] = \mathbf{WRR}^T\mathbf{W}^T + \Psi = \mathbf{WW}^T + \Psi$$

- ▶ There exists a couple of approaches to **reduce the number of dof**
  - ▶ Force  $\mathbf{W}$  to be **orthonormal** (this is the approach followed by PCA)
  - ▶ Force  $\mathbf{W}$  to be **lower triangular** together with  $W_{ii} > 0$  for all  $i$
  - ▶ **Sparsity prior** on the  $\mathbf{W}$  in the form of  $\ell_1$  regularization (this is known as sparse factor analysis)
  - ▶ Select the **rotation matrix**  $\mathbf{R}$  that leads to easier interpretations (e.g. enforce **sparsity**)
  - ▶ Use **non gaussian priors** on the latent variables  $\mathbf{z}_i$  (ICA)

# Mixture of Factor Analyzers

- ▶ So far we assumed that the means  $\boldsymbol{\mu} + \mathbf{W}\mathbf{z}_i$  all live in the same affine subspace.
- ▶ An alternative if we want to capture the low dimensional nature of the data locally (and keep a small number of latent variables) is to introduce multiple subspaces,  $\{\boldsymbol{\mu}_k, \mathbf{W}_k\}_{k=1}^K$
- ▶ The model, which is known as mixture of factor analyzers (MFA) then read as

$$\begin{aligned}p(\mathbf{x}_i | \mathbf{z}_i, q_i = k, \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k + \mathbf{W}_k \mathbf{z}_i, \boldsymbol{\Psi}) \\p(\mathbf{z}_i | \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{z}_i | \mathbf{0}, \mathbf{I}) \\p(q_i | \boldsymbol{\theta}) &= \text{Cat}(q_i | \boldsymbol{\theta})\end{aligned}$$

where we introduced the latent variables  $q_i$  which indicates the local subspace to be used and use the Categorical distribution to encode the corresponding distributions of those variables.



# Fitting Mixtures of FA and the EM algorithm

- ▶ When learning a FA model, we learn the parameters of the **posterior**  $p(\mathbf{x}_i|\mathbf{z}_i, \theta)$  and of the **prior**, or **latent distribution**  $p(\mathbf{z}_i)$ . Once we have those parameters, we usually want to see whether we can discover something meaningful on the data based the latent variables  $\mathbf{z}_i$
- ▶ One can then **analyze** the **shape** of  $p(\mathbf{z}_i|\mathbf{m}_i, \Sigma_i)$  (as we deal with Gaussian distribution it is possible to compute a closed form expression for this distribution)

$$\begin{aligned}p(\mathbf{z}_i|\mathbf{x}_i, \theta) &= \mathcal{N}(\mathbf{z}_i|\mathbf{m}_i, \Sigma_i) \\ \Sigma_i &= (\Sigma_0^{-1} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} \\ \mathbf{m}_i &= \Sigma_i (\mathbf{W}^T \Psi^{-1} (\mathbf{x}_i - \mu) + \Sigma_0^{-1} \mu_0)\end{aligned}$$

# Fitting Mixtures of FA and the EM algorithm

- ▶ The simplest way to fit an FA model is to use the EM algorithm
- ▶ Applying the exact same steps as for the GMM, we first estimate the responsibilities of each pair (cluster, prototype),  $(c, i)$  by using Bayes rule (E-step)

$$r_{i,c} = p(q_i = c | \mathbf{x}_i, \boldsymbol{\theta}) \propto \pi_c \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_c, \mathbf{W}_c \mathbf{W}_c^T + \boldsymbol{\Psi})$$

- ▶ In the M-Step, we then update the parameters using the parametrization of the posteriors  $p(\mathbf{z}_i | \mathbf{x}_i, q_i = c, \boldsymbol{\theta})$  (we assume  $\boldsymbol{\mu}_0 = \mathbf{0}$  and  $\boldsymbol{\Sigma}_0 = \mathbf{I}$ ),

$$\begin{aligned} p(\mathbf{z}_i | \mathbf{x}_i, q_i = c, \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{z}_i | \mathbf{m}_{ic}, \boldsymbol{\Sigma}_{ic}) \\ \boldsymbol{\Sigma}_{ic} &= (\mathbf{I}_L + \mathbf{W}_c^T \boldsymbol{\Psi}_c^{-1} \mathbf{W}_c)^{-1} \\ \mathbf{m}_{ic} &= \boldsymbol{\Sigma}_{ic} (\mathbf{W}_c^T \boldsymbol{\Psi}_c^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_c)) \end{aligned}$$

- ▶ The **M-step** is then completed by **estimating** the parameters  $\hat{\mathbf{W}}_c$ ,  $\hat{\Psi}$  and  $\hat{\pi}_c$  as

$$\hat{\mathbf{W}}'_c = \left[ \sum_i r_{ic} \mathbf{x}_i \mathbf{b}_{ic}^T \right] \left[ \sum_i r_{ic} \mathbf{C}_{ic} \right]^{-1}$$

$$\hat{\Psi} = \frac{1}{N} \text{diag} \left\{ \sum_{ic} r_{ic} \left( \mathbf{x}_i - \hat{\mathbf{W}}'_c \mathbf{b}_{ic} \right) \mathbf{x}_i^T \right\}$$

$$\hat{\pi}_c = \frac{1}{N} \sum_{i=1}^N r_{ic}$$

where we defined  $\mathbf{W}'_c$ ,  $\mathbf{b}_{ic}$  and  $\mathbf{C}_{ic}$  as  $\mathbf{W}'_c = [\mathbf{W}_c, \boldsymbol{\mu}_c]$

$$\mathbf{b}_{ic} = \mathbb{E} \left\{ \mathbf{z}' | \mathbf{x}_i, q_i = c \right\} = [\mathbf{m}_{ic}; 1], \quad \mathbf{z}' = (\mathbf{z}, 1),$$

$$\mathbf{m}_{ic} = \mathbb{E} \left\{ \mathbf{z}' (\mathbf{z}')^T | \mathbf{x}_i, q_i = c \right\},$$

$$= \begin{pmatrix} \mathbb{E} \left\{ \mathbf{z} \mathbf{z}^T | \mathbf{x}_i, q_i = c \right\} & \mathbb{E} \left\{ \mathbf{z} | \mathbf{x}_i, q_i = c \right\} \\ \mathbb{E} \left\{ \mathbf{z} | \mathbf{x}_i, q_i = c \right\} & 1 \end{pmatrix}$$

# Probabilistic PCA and classical PCA

- ▶ Constraint the covariance matrix of a FA model by requiring  $\Psi = \sigma^2 \mathbf{I}$  and  $\mathbf{W}$  to be orthonormal leads to the classical PCA model when  $\sigma^2 \rightarrow \infty$ .
- ▶ When we only require  $\sigma^2 > 0$ , the model is known as probabilistic PCA.
- ▶ The connections between PCA and probabilistic PCA as well as their respective (statistical) interpretation is given by writing down the data log likelihood  $\log p(\mathbf{X} | \mathbf{W}, \sigma^2)$

## From Probabilistic to classical PCA

### Probabilistic PCA (see Tipping, Bishop '99)

We consider a factor analysis model with  $\Psi = \sigma^2 \mathbf{I}$ . The data (or observed) log-likelihood is given by

$$\log p(\mathbf{X} | \mathbf{W}, \sigma^2) = -\frac{N}{2} \ln |\mathbf{C}| - \frac{1}{2} \sum_{i=1}^N \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i$$

where  $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$  and  $\mathbf{S} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T = \left(\frac{1}{N}\right) \mathbf{X}\mathbf{X}^T$  (again we assumed that the  $\mathbf{x}_i$  have been centered). The maxima of the log-likelihood are defined as

$$\hat{\mathbf{W}} = \mathbf{V}(\mathbf{\Lambda} - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}$$

Where  $\mathbf{R}$  is an arbitrary  $L \times L$  orthogonal matrix,  $\mathbf{V}$  is the  $D \times L$  matrix whose columns are the first  $L$  eigenvectors of  $\mathbf{S}$  and  $\mathbf{\Lambda}$  is the corresponding diagonal matrix of eigenvalues. Without loss of generality we can set  $\mathbf{R} = \mathbf{I}$ .

# Independent Component Analysis

- ▶ Just as PCA, ICA can be expressed a special instance of a Factor Analysis model. Recall that in FA we were expression the parameters as a linear function in the latent variables

$$\mathbf{x}_t = \mathbf{W}\mathbf{z}_t + \varepsilon_t$$

- ▶  $\mathbf{W}$  is thus called the **mixing matrix** and  $\varepsilon_t$  is viewed as some Gaussian noise  $\varepsilon_t \sim \mathcal{N}(0, \Psi)$
- ▶ In PCA we assumed that the source were independent and distributed following a Gaussian distribution,

$$p(\mathbf{z}_t) = \prod_{j=1}^L \mathcal{N}(z_{t_j} | 0, 1)$$

- ▶ In ICA, we relax the Gaussian assumption and let the source distributions be any non Gaussian distribution

$$p(\mathbf{z}_i) = \prod_{j=1}^L p_j(\mathbf{z}_{t_j})$$

# Independent Component Analysis as MLE

- ▶ Just as before, we can write the log-likelihood for ICA. Here we assume that the data has been **centered** and **whitened** (which can be done by a first application of PCA)
- ▶ The covariance reads as  $\mathbb{E} \{ \mathbf{x} \mathbf{x}^T \} = \mathbf{W} \mathbb{E} \{ \mathbf{z} \mathbf{z}^T \} \mathbf{W}^T$
- ▶ Using the whitening assumption,  $\mathbb{E} \{ \mathbf{z} \mathbf{z}^T \} = \mathbf{I}$  as well as the fact that the data is centered,  $\mathbb{E} \{ \mathbf{x} \mathbf{x}^T \}$ , we see that the matrix  $\mathbf{W}$  **must be orthogonal**
- ▶ Now using a change of variables, we can write the sample posterior  $p(\mathbf{x} | \mathbf{z}, \mathbf{W})$  as

$$\begin{aligned} p(\mathbf{x} | \mathbf{W}, \mathbf{z}) &= p_{\mathbf{x}}(\mathbf{W} \mathbf{z}) \\ &= p_{\mathbf{z}}(\mathbf{z}) |\det(\mathbf{W}^{-1})| \\ &= p_{\mathbf{z}}(\mathbf{W}^{-1} \mathbf{x}) |\det(\mathbf{W}^{-1})| \end{aligned}$$

# Independent Component Analysis as MLE

- ▶ From the posterior  $p(\mathbf{x}|\mathbf{W}, \mathbf{z})$ , the **data log-likelihood** for a set of  $T$  samples follows as

$$\frac{1}{T} \log p(\mathcal{D}|\mathbf{V}) = \log |\det(\mathbf{V})| + \frac{1}{T} \sum_{j=1}^L \sum_{t=1}^T \log p_j(\mathbf{v}_j^T \mathbf{x}_t)$$

- ▶ Using **orthogonality of the rows** of  $\mathbf{V}$ ,  $\mathbf{v}_j$ , and replacing the sum over the data with a population average, we get the reduced formulation for the negative LL

$$NLL(\mathbf{V}) = \sum_{j=1}^L \mathbb{E} \{ G_j(z_j) \}$$

where  $z_j = \mathbf{v}_j^T \mathbf{x}$  and  $G_j(z) = -\log p_j(z)$ .

- ▶ We then **minimize the NLL** under the constraints that the rows  $\mathbf{v}_j$  are **orthogonal** and have **unit norm** (which follows from the whitening assumption and  $\mathbb{E} \{ \mathbf{v}_j^T \mathbf{x} \mathbf{x}^T \mathbf{v}_j \} = \|\mathbf{v}_j\|^2 = \mathbb{E} \{ z_j^2 \}$ )



# Fast ICA on the NLL (I)

- ▶ **Whitening** and **centering** are essentially used to **reduce** the **computational complexity** as they reduce the number of parameters from  $n^2$  to  $n(n-1)/2$
- ▶ **Whitening** can be obtained with a **first application of PCA** from where one can then apply **any Fast ICA** algorithm relying on the orthogonality of the matrix  $\mathbf{W}$ .
- ▶ There exists several algorithms to perform ICA. Here we focus on an algorithm that can be used to minimize the NLL.
- ▶ Fast ICA on the NLL can be considered a particular instance of a **Newton method**.

## Fast ICA on the NLL (II)

- ▶ For the negative log-likelihood derived earlier, if we let  $g = \frac{d}{dz} G(z)$ , in the constrained framework, the contributions from each independent component to objective **function**, **gradient** and **Hessian** can respectively read as

$$\begin{aligned}f(\mathbf{v}) &= \mathbb{E} \left\{ G(\mathbf{v}^T \mathbf{x}) \right\} + \lambda(1 - \mathbf{v}^T \mathbf{v}) \\ \nabla f(\mathbf{v}) &= \mathbb{E} \left\{ \mathbf{x} g(\mathbf{v}^T \mathbf{x}) \right\} - \beta \mathbf{v} \\ \mathbf{H}(\mathbf{v}) &= \mathbb{E} \left\{ \mathbf{x} \mathbf{x}^T g'(\mathbf{v}^T \mathbf{x}) \right\} - \beta \mathbf{I}\end{aligned}$$

$\beta = 2\lambda$  is a Lagrange multiplier.

- ▶ If we make the **approximation**,  
 $\mathbb{E} \left\{ \mathbf{x} \mathbf{x}^T g'(\mathbf{v}^T \mathbf{x}) \right\} \approx \mathbb{E} \left\{ \mathbf{x} \mathbf{x}^T \right\} \mathbb{E} \left\{ g'(\mathbf{v}^T \mathbf{x}) \right\} = \mathbb{E} \left\{ g'(\mathbf{v}^T \mathbf{x}) \right\}$ ,  
The Hessian is easy to invert and we get the **Newton step**

$$\mathbf{v} \leftarrow \mathbf{v} - \frac{\mathbb{E} \left[ \mathbf{x} g(\mathbf{v}^T \mathbf{x}) \right] - \beta \mathbf{v}}{\mathbb{E} \left[ g'(\mathbf{v}^T \mathbf{x}) \right] - \beta}$$

## Fast ICA on the NLL (III)

- ▶ After the Newton step has been applied, we simply **project** the resulting vector  $\mathbf{v}$  onto the subspace **orthogonal to the other independent components** and normalize it.
- ▶ As the objective is **non convex**, there are **multiple local minimas**

## Possible distributions

- ▶ As we have seen, Gaussian priors won't work well for ICA so **what distributions** can we use instead?
- ▶ There are several possible distributions one can use besides the Gaussian distribution:
  - ▶ **Super-Gaussian** distributions (e.g. **Laplace** distribution). Super Gaussian distributions are distributions with a big spike at the mean and heavy tails. Generally speaking we say that a distribution is Super Gaussian when its **kurtosis**,  $\text{kurt}(z) = \mu^4/\sigma^4 - 3$  is positive,  $\text{kurt}(z) > 0$ . Here  $\mu_k = \mathbb{E} \{ (X - \mathbb{E}(X))^k \}$
  - ▶ **Sub-Gaussian** distributions. (e.g. **uniform distribution**). Subgaussian distributions have negative kurtosis.
  - ▶ **Skewed distributions** (e.g. **Gamma** distribution). A distribution can be different from the Gaussian distribution by being asymmetric. We define the skewness (measure of asymmetry) of a distribution as  $\text{skew}(z) = \mu^3/\sigma^3$ .