

# Bayesian Models

Sandro Cumani

sandro.cumani@polito.it

Politecnico di Torino

The frequentist approach to inference we have considered in previous lectures consists in estimating “optimal” model parameters by defining a suitable objective function and computing model parameters that optimize the objective function

A possible criterion consists in computing Maximum Likelihood estimates, i.e. estimating parameters that maximize the likelihood of our data (observed features in the generative framework, observed class labels in the discriminative framework)

ML point estimates, however, do not consider alternative, possible explanations of our data that may be less likely, but still properly model the distribution of our data — they neglect, as all point estimates, that our parameter estimates are **uncertain**

# Bayesian models

The Bayesian framework allows taking into account our uncertainty over the model parameters

Since model parameters are not known, they are treated as uncertain quantities, that can be described in terms of **probability distributions**

Model parameters are thus represented in terms of Random Variables, and inference is based on probabilistic reasoning

In the following we consider again a classification problem (although our discussion can be adapted to density estimation problems)

We also focus on generative approaches

We want to predict the class label for an unlabeled sample  $x_t$

We have at our disposal a labeled training set of samples

$$\mathcal{D} = \{(\mathbf{x}_1, c_1) \dots (\mathbf{x}_n, c_n)\}$$

or, in a more compact way:

$$\mathcal{D} = (\mathcal{X}, \mathcal{C})$$

For a generative model, we have seen that we can compute the class posterior distribution for the test samples from the joint distribution of observed features and class labels

Assuming that features are described by model  $\mathcal{M}$  (i.e. a Gaussian model, a GMM model, ...), then inference is based on computing a **predictive distribution**

$$f_{X_t, C_t | \mathcal{D}, \mathcal{M}}(\mathbf{x}_t, c | \mathcal{D})$$

i.e., the joint distribution for features and classes, **given the knowledge we have** (the training set) and **the model we chose** ( $\mathcal{M}$ )

The predictive distribution can be expressed as a ratio between the joint distribution of the observed samples (training samples and test samples) and labels (observed training labels, hypothesized test label  $c$ ) and the joint distribution for the observed training data (Bayes rule):

$$f_{\mathbf{X}_t, C_t | \mathcal{D}, \mathcal{M}}(\mathbf{x}_t, c | \mathcal{D}) = \frac{f_{\mathbf{X}_t, C_t, \mathcal{D} | \mathcal{M}}(\mathbf{x}_t, c, \mathcal{D})}{f_{\mathcal{D} | \mathcal{M}}(\mathcal{D})}$$

Our problem thus consists in modeling the **joint distributions** of sets of feature vectors and corresponding labels

# Bayesian models

Model  $\mathcal{M}$  describes the joint distribution for a set of samples and classes in terms of a parametric distribution family

It allows defining conditional distributions of the form

$$f_{X_1 \dots X_m, C_1 \dots C_m | \Theta}(\mathbf{x}_1, \dots, \mathbf{x}_m, c_1, \dots, c_m | \theta)$$

where  $X_1 \dots X_m$  is a (generic) set of  $m$  R.V.s representing  $m$  features vectors, and  $C_1 \dots C_m$  are the R.V.s representing the corresponding labels

The model provides conditional distribution given a value  $\theta$  for the parameters

$\Theta$  is the R.V. that describes the model parameters

# Bayesian models

In the following we restrict the analysis to models that assume that observations are **independent and identically distributed** (i.i.d.) given the model parameters:

$$[(X_i, C_i) \perp\!\!\!\perp (X_j, C_j)] | (\Theta = \theta)$$

and

$$(X_i, C_i) | (\Theta = \theta) \sim (X_j, C_j) | (\Theta = \theta) \sim (X, C) | (\Theta = \theta)$$

thus

$$f_{X_1 \dots X_m, C_1 \dots C_m | \Theta}(\mathbf{x}_1, \dots, \mathbf{x}_m, c_1, \dots, c_m | \theta) = \prod_{i=1}^m f_{X, C | \Theta}(\mathbf{x}_i, c_i | \theta)$$



# Bayesian models

Given a **prior distribution**  $f_{\Theta}(\theta)$  for the model parameters  $\Theta$ , we can then compute the joint distribution of a set of samples  $\mathcal{D}_S = [\mathcal{X}_S, \mathcal{C}_S]$ , with  $\mathcal{X}_S = [x_1 \dots x_m]$ ,  $\mathcal{C}_S = [c_1 \dots c_m]$ :

$$\begin{aligned} f_{\mathcal{D}_S, \Theta}(\mathcal{D}_S, \theta) &= f_{X_S, C_S | \Theta}(\mathcal{X}_S, \mathcal{C}_S | \theta) f_{\Theta}(\theta) \\ &= \left[ \prod_{i=1}^m f_{X_i, C_i | \Theta}(x_i, c_i | \theta) \right] f_{\Theta}(\theta) \end{aligned}$$

The joint distribution consists of a product of sample **likelihoods** and parameters **prior**

The distribution for the data R.V.  $\mathcal{D}_S$  is obtained by **marginalizing** with respect to the model parameters:

$$\begin{aligned} f_{\mathcal{D}_S}(\mathcal{D}_S) &= \int f_{\mathcal{D}_S, \Theta}(\mathcal{D}_S, \theta) d\theta \\ &= \int \left[ \prod_{i=1}^m f_{X_i, C_i | \Theta}(x_i, c_i | \theta) \right] f_{\Theta}(\theta) d\theta \end{aligned}$$

We can apply these expressions to compute our predictive distribution (we drop the explicit conditioning on  $\mathcal{M}$  for readability):

$$f_{X_t, C_t | \mathcal{D}}(\mathbf{x}_t, c | \mathcal{D}) = \frac{f_{X_t, C_t, \mathcal{D}}(\mathbf{x}_t, c, \mathcal{D})}{f_{\mathcal{D}}(\mathcal{D})}$$

with

$$f_{X_t, C_t, \mathcal{D}}(\mathbf{x}_t, c, \mathcal{D}) = \int f_{X_t, C_t | \Theta}(\mathbf{x}_t, c_t | \theta) f_{\mathcal{D} | \Theta}(\mathcal{D} | \theta) f_{\Theta}(\theta) d\theta$$

$$f_{\mathcal{D}}(\mathcal{D}) = \int f_{\mathcal{D} | \Theta}(\mathcal{D} | \theta) f_{\Theta}(\theta) d\theta$$

The factorizations derive from our independence assumptions

We can also apply Bayes rule to write the predictive distribution in an equivalent form:

$$\begin{aligned}f_{X_t, C_t | \mathcal{D}}(\mathbf{x}_t, c | \mathcal{D}) &= \int f_{X_t, C_t, \Theta | \mathcal{D}}(\mathbf{x}_t, c, \theta | \mathcal{D}) d\theta \\&= \int f_{X_t, C_t | \Theta, \mathcal{D}}(\mathbf{x}_t, c | \theta, \mathcal{D}) f_{\Theta | \mathcal{D}}(\theta | \mathcal{D}) d\theta \\&= \int f_{X_t, C_t | \Theta}(\mathbf{x}_t, c | \theta) f_{\Theta | \mathcal{D}}(\theta | \mathcal{D}) d\theta\end{aligned}$$

The second step derives from the application of Bayes rule, whereas the third step derives from the assumption that samples are independent given the model parameters

The first term inside the integral is the likelihood

$$f_{X_t, C_t | \Theta}(\mathbf{x}_t, c | \theta)$$

If we assume that the prior class distribution is independent of the model parameters, we can also write

$$f_{X_t, C_t | \Theta}(\mathbf{x}_t, c | \theta) = f_{X_t | C_t, \Theta}(\mathbf{x}_t | c, \theta) P(C_t = c)$$

i.e., the product of class-conditional likelihoods and class priors.

The predictive distribution becomes

$$\int f_{X_t | C_t, \Theta}(\mathbf{x}_t | c, \theta) f_{\Theta | \mathcal{D}}(\theta | \mathcal{D}) d\theta \cdot P(C_t = c)$$

i.e. a product of the class prior and the class-conditional predictive distribution

$$f_{X_t | C_t, \mathcal{D}}(\mathbf{x}_t | c, \mathcal{D}) = \int f_{X_t | C_t, \Theta}(\mathbf{x}_t | c, \theta) f_{\Theta | \mathcal{D}}(\theta | \mathcal{D}) d\theta$$

# Bayesian models

The second term

$$f_{\Theta|\mathcal{D}}(\theta|\mathcal{D})$$

is the **posterior distribution** for the model parameters **given the training data** (and the chosen model  $\mathcal{M}$ )

It represents our belief over the possible model parameters, given the training dataset

Before we observe the training data, our belief corresponds to the prior distribution  $f_{\Theta}(\theta)$

Once we have observed data  $\mathbf{D} = \mathcal{D}$ , our belief is updated to  $f_{\Theta|\mathcal{D}}(\theta|\mathcal{D})$

# Bayesian models

The posterior distribution can be expressed in terms of conditional likelihoods and priors:

$$f_{\Theta|\mathcal{D}}(\theta|\mathcal{D}) = \frac{f_{\Theta,\mathcal{D}}(\theta,\mathcal{D})}{f_{\mathcal{D}}(\mathcal{D})} = \frac{f_{\mathcal{D}|\Theta}(\mathcal{D}|\theta)f_{\Theta}(\theta)}{\int f_{\mathcal{D}|\Theta}(\mathcal{D}|\theta)f_{\Theta}(\theta)d\theta} \\ \propto f_{\mathcal{D}|\Theta}(\mathcal{D}|\theta)f_{\Theta}(\theta)$$

The posterior distribution is **proportional** to the product of the **likelihood** term

$$f_{\mathcal{D}|\Theta}(\mathcal{D}|\theta)$$

and the **prior** distribution

$$f_{\Theta}(\theta)$$

The proportionality factor is the value that ensures that the posterior distribution integrates to 1.

# Bayesian models

In the Bayesian framework “learning” consists in updating our beliefs as we collect more evidence

For example, if we have computed the posterior distribution given a set  $\mathcal{D}_1$  of training data, and we then collect additional data  $\mathcal{D}_2$ , we can refine our beliefs as

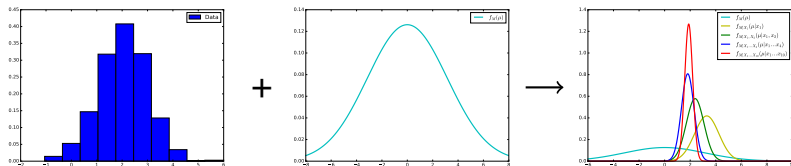
$$\begin{aligned} f_{\Theta|\mathcal{D}_1, \mathcal{D}_2}(\theta|\mathcal{D}_1, \mathcal{D}_2) &\propto f_{\mathcal{D}_1|\Theta}(\mathcal{D}_1|\theta)f_{\mathcal{D}_2|\Theta}(\mathcal{D}_2|\theta)f_{\Theta}(\theta) \\ &\propto f_{\mathcal{D}_2|\Theta}(\mathcal{D}_2|\theta)f_{\Theta|\mathcal{D}_1}(\theta|\mathcal{D}_1) \end{aligned}$$

i.e., the posterior given the two datasets is proportional to the product of **likelihood computed over the new data** and the **posterior distribution given the initial dataset**

The posterior distribution given dataset  $\mathcal{D}_1 = \mathcal{D}_1$  acts as a **prior** for the remaining data, and collects all the information we have extracted from the initial dataset

# Bayesian models

From prior to posterior distributions — modeling the mean of a Gaussian distribution:



We can observe that, as we increase the data size, the posterior distribution tends to concentrate more and more around a single value

This corresponds to a decrease in our uncertainty over the possible model parameters



# Bayesian models

In theory, once we have chosen a prior distribution, we have all the elements to compute the predictive distribution and thus to perform inference

However, in practice computing predictive distributions requires solving integrals, which may prove a difficult task

Furthermore, the results will depend on our choice of the prior, thus we need to decide what kind of prior is appropriate for our task

The first issue can be addressed by using **conjugate prior distributions**

We say that a prior distribution is **conjugate** to a likelihood function if the corresponding posterior distribution belongs to the same family as the prior distribution, i.e. they have the same **functional form**

In this case, since we know the posterior distribution up to a normalization factor, the normalization factor can be recovered from the distribution family

For example, we can consider a Gaussian likelihood for a (un-labeled) dataset  $\mathcal{X}$ , parametrized by the distribution mean (we assume that the covariance matrix is known and fixed)

Let the log-likelihood (a function of the mean  $\mu$ ) be

$$\begin{aligned}\ell(\mu) &= \log \mathcal{L}(\mu) = \log f_{\mathcal{X}|\mathbf{M}}(\mathcal{X}|\mu) \\ &= \sum_i \frac{1}{2} \log |\Lambda| - \frac{1}{2} (x_i - \mu)^T \Lambda (x_i - \mu) + \xi_1 \\ &= \mu^T \Lambda F - \frac{n}{2} \mu^T \Lambda \mu + \xi_2\end{aligned}$$

where  $F = \sum_i x_i$ ,  $n$  is the number of samples in  $\mathcal{X}$  and  $\xi_2$  collects all terms that do not depend on  $\mu$

# Bayesian models

The posterior distribution for  $\mu$  is proportional to the product of the prior distribution and the likelihood.

Let the prior distribution belong to the Gaussian family

$$f_M(\mu) = \mathcal{N}(\mu|m, P^{-1})$$

where  $m$  and  $P$  are the mean and precision matrix of the prior

The functional form of the prior distribution is

$$\log f_M(\mu) = \mu^T P m - \frac{1}{2} \mu^T P \mu + \xi$$

A log-density defined, up to a constant term, by the above expression corresponds to a normal density  $\mathcal{N}(\mu|m, P^{-1})$

The constant term  $\xi$  corresponds to the logarithm of the normalization constant

The posterior distribution for  $\mathbf{M}$  corresponds to

$$f_{\mathbf{M}|\mathcal{X}}(\boldsymbol{\mu}|\mathcal{X}) \propto \mathcal{L}(\boldsymbol{\mu})\mathcal{N}(\boldsymbol{\mu}|\mathbf{m}, \mathbf{P}^{-1})$$

In log form

$$\begin{aligned}\log f_{\mathbf{M}|\mathcal{X}}(\boldsymbol{\mu}|\mathcal{X}) &= \ell(\boldsymbol{\mu}) - \frac{1}{2}\boldsymbol{\mu}^T \mathbf{P} \boldsymbol{\mu} + \boldsymbol{\mu}^T \mathbf{P} \mathbf{m} + \xi_3 \\ &= \boldsymbol{\mu}^T (\boldsymbol{\Lambda} \mathbf{F} + \mathbf{P} \mathbf{m}) - \frac{1}{2}\boldsymbol{\mu}^T (n\boldsymbol{\Lambda} + \mathbf{P}) \boldsymbol{\mu} + \xi_4\end{aligned}$$

The expression has the functional form of a multivariate Gaussian distribution:

$$f_{\mathbf{M}|\mathcal{X}}(\boldsymbol{\mu}|\mathcal{X}) = \mathcal{N}(\boldsymbol{\mu}|\mathbf{m}_x, \mathbf{P}_x^{-1}) \iff \log f_{\mathbf{M}|\mathcal{X}}(\boldsymbol{\mu}|\mathcal{X}) = \boldsymbol{\mu}^T \mathbf{P}_x \mathbf{m}_x - \frac{1}{2}\boldsymbol{\mu}^T \mathbf{P}_x \boldsymbol{\mu} + \xi$$

The posterior distribution for the mean is therefore again a Gaussian distribution  $\mathcal{N}(\mu|\mathbf{m}_x, \mathbf{P}_x^{-1})$

The parameters can be recovered by inspection:

$$\mathbf{P}_x = \mathbf{P} + n\mathbf{\Lambda}$$

$$\mathbf{m}_x = \mathbf{P}_x^{-1} (\mathbf{\Lambda}\mathbf{F} + \mathbf{P}\mathbf{m}) = (\mathbf{P} + n\mathbf{\Lambda})^{-1} (\mathbf{\Lambda}\mathbf{F} + \mathbf{P}\mathbf{m})$$

Using conjugate priors we can avoid solving difficult integrals — the posterior is in the same family as the prior, so we already know the values for the normalization constant

With conjugate priors we are also able to compute the integral required for the predictive distribution

We can observe that we can also express the predictive distribution in terms of posterior distributions using Bayes rule:

$$f_{X_t|C_t,D}(x_t|c, \mathcal{D}) = \frac{f_{X_t|C_t,\Theta}(x_t|c, \theta) f_{\Theta|D}(\theta|\mathcal{D})}{f_{\Theta|D,X_t,C_t}(\theta|\mathcal{D}, x_t, c)}$$

The expression is valid for any value of  $\theta$  for which the denominator is non-zero, and provides an alternative to solving the predictive distribution integral (in practice, we already solved the integration when computing the posterior density)

# Bayesian models

Let's see what we get if we use the ML solution for  $\theta$  to compute the previous expression:

$$f_{X_t|C_t, \mathcal{D}}(\mathbf{x}_t|c, \mathcal{D}) = f_{X_t|C_t, \Theta}(\mathbf{x}_t|c, \theta^*) \frac{f_{\Theta|\mathcal{D}}(\theta^*|\mathcal{D})}{f_{\Theta|\mathcal{D}, X_t, C_t}(\theta^*|\mathcal{D}, \mathbf{x}_t, c)}$$

We can contrast with the ML framework: our predictive distribution was

$$f_{X_t|C_t, \Theta}(\mathbf{x}_t|c, \theta^*)$$

For a binary classification problem we can see that the log-likelihood would be

$$\log \frac{f_{X_t|C_t, \mathcal{D}}(\mathbf{x}_t|\mathcal{H}_T, \mathcal{D})}{f_{X_t|C_t, \mathcal{D}}(\mathbf{x}_t|\mathcal{H}_F, \mathcal{D})} = \log \frac{f_{X_t|C_t, \Theta}(\mathbf{x}_t|\mathcal{H}_T, \theta^*)}{f_{X_t|C_t, \Theta}(\mathbf{x}_t|\mathcal{H}_F, \theta^*)} - \log \frac{f_{\Theta|\mathcal{D}, X_t, C_t}(\theta^*|\mathcal{D}, \mathbf{x}_t, \mathcal{H}_T)}{f_{\Theta|\mathcal{D}, X_t, C_t}(\theta^*|\mathcal{D}, \mathbf{x}_t, \mathcal{H}_F)}$$



# Bayesian models

The prior distribution encodes our prior beliefs over the model parameters, and can have a large impact over the model when we have limited amounts of data

Priors may be chosen to encode expert knowledge

In other cases, we may not have strong indications on preferred models

What kind of priors should we use?

If we have no reason to choose a strongly informative prior, we may resort to **non-informative** or **weakly informative** priors

- Flat prior  $f_{\Theta}(\theta) \propto \alpha$
- Jeffrey's priors
- ...

The flat prior, for example, encodes that all values of the parameters are, a priori, equally likely

We should note that a flat prior is often **improper**, in that it cannot be normalized, since its integral diverges (e.g. this happens if we use a flat prior over  $\mathbb{R}$ )

Improper priors may still, in many cases, provide proper posteriors, and thus may be still used in the Bayesian framework

In practice, it may be better to represent an improper prior as a limit of a proper prior. For example, we may represent an “almost” flat prior as

$$f_{\Theta}(\theta) = \mathcal{N}(0, \varepsilon^{-1} \mathbf{I})$$

and then let  $\varepsilon \rightarrow 0$

In this case, it's better to work with the proper prior and pass to the limit only at the end

Another option is to use **empirical** priors: we can represent the prior as a parametric distribution, and learn the parameters of the prior distribution **from the data**, for example maximizing the marginal likelihood

These kind of models are also known as **empirical Bayes** approaches

We will see an example shortly

The frequentist approach we have discussed in previous lectures can be interpreted as an approximation of the Bayesian approach

We can interpret frequentist estimators as approximations of posterior distribution

$$f_{\Theta|\mathcal{D}}(\theta|\mathcal{D}) \approx \delta_{\theta^*}(\theta)$$

where the posterior is approximated by a **delta** distribution centered around the estimated value  $\theta^*$

The predictive distribution becomes

$$f_{X_t|C_t,\mathcal{D}}(\mathbf{x}_t|c, \mathcal{D}) = \int f_{X_t|C_t,\Theta}(\mathbf{x}_t|c, \theta) \delta_{\theta^*}(\theta) d\theta = f_{X_t|C_t,\Theta}(\mathbf{x}_t|c, \theta^*)$$

Since the delta distribution is approximating our posterior, we may select as value for  $\theta^*$  the value that is **most likely** according to the posterior distribution

$$\begin{aligned}\theta_{MAP}^* &= \arg \max_{\theta} f_{\Theta|D}(\theta|D) \\ &= \arg \max_{\theta} f_{D|\Theta}(\mathcal{D}|\theta)f_{\Theta}(\theta) \\ &= \arg \max_{\theta} \mathcal{L}(\theta)f_{\Theta}(\theta)\end{aligned}$$

The estimator  $\theta_{MAP}^*$  is called **Maximum A Posteriori** (MAP) estimator, since its the most likely value according to the **posterior** distribution

The approximation will become more precise as the number of samples increases, since the variance of the posterior distribution will tend to decrease

If we choose a flat prior  $f_{\Theta}(\theta) \propto \alpha$ , then the MAP estimate becomes

$$\theta_{ML}^* = \arg \max_{\theta} \mathcal{L}(\theta)$$

i.e., the ML estimator that **maximizes the likelihood of the training data**

We can thus re-interpret ML estimators as an approximation of the Bayesian posterior assuming flat priors

The approximation is good if we have a large number of training samples

We can also re-interpret the regularized logistic regression model

$$\arg \min_{\mathbf{w}, b} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_i \log \left[ 1 + e^{-z_i(\mathbf{w}^T \mathbf{x}_i + b)} \right]$$

as a **MAP estimate**, where the prior for the model parameters  $\mathbf{w}$  is a Gaussian distribution

$$f_{\mathbf{w}}(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$



# Bayesian models

We now consider some examples that allow us to better understand how to perform inference using Bayesian models

We start considering again the multinomial model for text classification

The multinomial model assumes that tokens of each class can be modeled by multinomial distributions with parameters  $\pi_1 \dots \pi_K$ , where  $K$  is the number of classes

The log-likelihood of the model is

$$\ell(\pi_1, \dots, \pi_K) = \sum_{c=1}^K \sum_{j=1}^m N_{c,j} \log \pi_{c,j} + \xi$$

# Bayesian models

In the following, we assume that the model parameters are represented by i.i.d. R.V.s  $\Pi_1 \dots \Pi_K \sim \Pi$ , with prior distribution given by

$$f_{\Pi_i}(\pi) = f_{\Pi}(\pi)$$

We can show that a conjugate prior for  $\Pi_1 \dots \Pi_K$  is the Dirichlet prior

$$\Pi \sim \text{Dir}(\alpha) , \quad f_{\Pi}(\pi) = \frac{1}{B(\alpha)} \prod_{i=1}^m \pi_i^{\alpha_i - 1}$$

where  $B(\alpha)$  is the Beta function

$$B(\alpha) = \frac{\prod_{i=1}^m \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^m \alpha_i)}$$

$\Gamma$  is the Gamma function, a generalization of the factorial function to non-integer values — for example,  $\Gamma(n) = (n-1)!$  and  $\Gamma(x) = x\Gamma(x-1)$

$\alpha$  is a vector of  $m$  positive, real values that parametrize the Dirichlet density, and  $m$  is the number of unique tokens (size of the dictionary)

The prior log-density for  $\Pi_c$  therefore has the form

$$\log f_{\Pi_c|\alpha}(\pi_c|\alpha) = \sum_{i=1}^m (\alpha_i - 1) \log \pi_{c,i} + \xi$$

for some constant<sup>1</sup>  $\xi$

Since we assume that  $\Pi_1 \dots \Pi_K$  are i.i.d., the prior log-density is

$$\log f_{\Pi_1 \dots \Pi_K|\alpha}(\pi_1 \dots \pi_K|\alpha) = \sum_{c=1}^K \sum_{i=1}^m (\alpha_i - 1) \log \pi_{c,i} + \xi$$

---

<sup>1</sup>In the following we will denote all irrelevant constants with  $\xi$ , even if they refer to different actual quantities

The posterior distribution is proportional to the product of likelihood and prior, thus

$$\begin{aligned}\log f_{\Pi_1 \dots \Pi_K | \mathcal{D}, \alpha}(\pi_1 \dots \pi_K | \mathcal{D}, \alpha) \\&= \sum_{c=1}^K \sum_{i=1}^m N_{c,i} \log \pi_{c,i} + \sum_{c=1}^K \sum_{i=1}^m (\alpha_i - 1) \log \pi_{c,i} + \xi \\&= \sum_{c=1}^K \sum_{i=1}^m (N_{c,i} + \alpha_i - 1) \log \pi_{c,i} + \xi\end{aligned}$$

We can observe that the posterior distribution factorizes over the terms  $\Pi_c$  of the different classes:

$$\log f_{\Pi_1 \dots \Pi_K | \mathcal{D}, \alpha}(\pi_1 \dots \pi_K | \mathcal{D}, \alpha) = \sum_{c=1}^K \log f_{\Pi_c | \mathcal{D}_c, \alpha}(\pi_c | \mathcal{D}_c, \alpha)$$

where  $\mathcal{D}_c$  are the training samples for class  $c$

The terms  $f_{\Pi_c|\mathcal{D}_c,\alpha}(\pi_c|\mathcal{D}_c,\alpha)$  are the posterior distributions of the different model parameters, given the data of the corresponding class

$$f_{\Pi_c|\mathcal{D}_c,\alpha}(\pi_c|\mathcal{D}_c,\alpha) = \sum_{i=1}^m (N_{c,i} + \alpha_i - 1) \log \pi_{c,i} + \xi$$

We can recognize that this is the log-density of a Dirichlet distribution

$$\Pi_c|\mathcal{D}_c = \mathcal{D}_c, \alpha \sim \text{Dir}(\alpha + N_c)$$

where  $N_c = [N_{c,1} \dots N_{c,m}]$

We can interpret the parameters of the posterior distribution as counts of each token. The term  $\alpha$  represents “pseudo-counts” that we are adding to the counts

Before we consider the predictive distribution, we focus on the selection of the values  $\alpha$  of the prior distribution

From the form of the density of the Dirichlet distribution, we can observe that setting

$$\alpha_i = 1, \quad \forall i$$

corresponds to a uniform (flat) prior

$$f_{\Pi|\alpha=\mathbf{1}}(\pi|\alpha=\mathbf{1}) = \frac{1}{B(\mathbf{1})}$$

In this case, the prior is proper, since the support for  $\Pi$  is the unit-box

If we choose to use the flat prior, then the posterior distribution becomes

$$\boldsymbol{\Pi}_c | \mathcal{D}_c = \mathcal{D}_c, \boldsymbol{\alpha} = \mathbf{1} \sim \text{Dir}(N_c + \mathbf{1})$$

The MAP solution in this case corresponds to the ML solution

However, we can observe that the mean of the posterior distribution is a vector with components

$$\mathbb{E} [\boldsymbol{\Pi}_c | \mathcal{D}_c = \mathcal{D}_c, \boldsymbol{\alpha} = \mathbf{1}]_i = \frac{N_{c,i} + 1}{\sum_{j=1}^m N_{c,j} + 1}$$

If, rather than the MAP / ML estimates, we approximate the posterior with the **posterior mean**, then we retrieve the pseudo-count solution with pseudo-counts  $\varepsilon = 1.0$  of Laboratory 6

The same solution corresponds also to a MAP estimate, but with a non-flat prior  $\alpha_i = 1 + \varepsilon$

The predictive distribution can be obtained by integrating the product of the class-conditional likelihood for the test sample and the posterior probability of the parameters given the training samples.

Let  $\mathbf{y}_t$  be the vector of token occurrences for the test document. The predictive distribution is

$$f_{Y_t|C,D,\alpha}(\mathbf{y}_t|c, \mathcal{D}, \alpha) = \int f_{Y_t|\Pi_c}(\mathbf{y}_t|\pi) f_{\Pi_c|D_c,\alpha}(\pi|\mathcal{D}_c, \alpha) d\pi$$

The first term in the integral is a multinomial p.d.f., whereas the second term is the Dirichlet density of the posterior distribution for the parameters. We can show that

$$f_{Y_t|C,D,\alpha}(\mathbf{y}_t|c, \mathcal{D}, \alpha) = \frac{\left(\sum_{j=1}^m \mathbf{y}_{t,j}\right) B\left(\sum_{j=1}^m (N_{c,j} + \alpha_j), \sum_{j=1}^m \mathbf{y}_{t,j}\right)}{\prod_{j=1}^m \mathbf{y}_{t,j} B(N_{c,j} + \alpha_j, \mathbf{y}_{t,j})}$$



Compared to the ML solution, the Bayesian approach can provide better results. For the task of Laboratory 6 we have

Model	Accuracy
ML $\varepsilon = 10^{-6}$	52%
Mean point estimate (ML + pseudo-counts $\varepsilon = 1$ )	58%
Bayes, flat prior $\alpha_i = 1$	58%

The Bayesian approach achieves better solution than the ML solution (note: the ML solution is actually already using small pseudo-counts to avoid numerical issues)

The posterior mean provides a good approximation of the posterior distribution, and indeed allows us to obtain the same results as for the Bayesian model

The second example we analyze considers a tied covariance Gaussian model, but introduces a Bayesian modeling of the class means

As we have seen, the tied covariance model can be interpreted as an extension of LDA

The model we are going to consider can thus be interpreted as a probabilistic version of LDA, in the sense that class means will be modeled in a probabilistic way

As we shall see, Bayesian treatment of the class means will also allow us to extend the model to open-set classification

We consider a  $K$ -class problem, where the samples of each class are modeled by a Gaussian density

$$f_{X_i|C, \mu_1, \dots, \mu_K, \Lambda}(\mathbf{x}_i|c, \mu_1 \dots \mu_K, \Lambda) = \mathcal{N}(\mathbf{x}_i|\mu_c, \Lambda^{-1})$$

We also adopt a slightly different notation, and represent with  $\mathbf{x}_{c,i}$  the  $i$ -th training set sample of class  $c$

We have seen that the model allows interpreting a sample as the sum of the corresponding class mean and i.i.d., zero-mean, Gaussian distributed noise

$$\mathbf{x}_{c,i} = \mu_c + \varepsilon_{c,i}, \quad \varepsilon_{c,i} \sim \mathbf{E}_{c,i} \sim \mathcal{N}(\mathbf{0}, \Lambda^{-1})$$

$\mathbf{E}_{c,i}$  is the R.V. that represents the noise.  $\varepsilon_{c,i}$  is the noise realization

The frequentist approach we have followed consisted in performing ML estimation of the class means  $\mathbf{M}$  and of the noise covariance matrix  $\Sigma = \Lambda^{-1}$

We now consider a Bayesian treatment of the class means

Again, we resort to a conjugate prior to be able to compute closed form solutions for our integrals

We have seen that a Gaussian distribution is a conjugate prior for the mean of a Gaussian likelihood

We can use a parametric prior, and assume that the means are a priori distributed as

$$\mathbf{M}_1 \sim \mathbf{M}_2 \sim \dots \sim \mathbf{M}_K \sim \mathcal{N}(\mathbf{m}, \mathbf{B})$$

Our model is thus specified by the likelihood function

$$f_{X|C, \mathbf{M}_1 \dots \mathbf{M}_K}(\mathbf{x}|c, \boldsymbol{\mu}_1 \dots \boldsymbol{\mu}_K) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_c, \boldsymbol{\Lambda}^{-1})$$

and the prior for the class means

$$f_{\mathbf{M}_c|m, \mathbf{B}}(\boldsymbol{\mu}|m, \mathbf{B}) = \mathcal{N}(m, \mathbf{B})$$

The **complete data** likelihood for the training set is given by

$$f_{D, \mathbf{M}_1 \dots \mathbf{M}_K|m, \mathbf{B}, \boldsymbol{\Lambda}}(\mathcal{D}, \boldsymbol{\mu}_1 \dots \boldsymbol{\mu}_K|m, \mathbf{B}, \boldsymbol{\Lambda}) = \\ \prod_{c=1}^K \left[ \left( \prod_{i=1}^{n_c} \mathcal{N}(x_{c,i}|\boldsymbol{\mu}_c, \boldsymbol{\Lambda}^{-1}) \right) \mathcal{N}(\boldsymbol{\mu}_c|m, \mathbf{B}) \right]$$

We can see that the likelihood factorizes over classes, and thus the same applies to the posterior distribution for  $\mathbf{M}_1 \dots \mathbf{M}_K$

# Bayesian models

For each term  $M_i$ , the complete data likelihood is proportional to a product of a Gaussian likelihood and a Gaussian prior

We have already shown that this corresponds to a Gaussian posterior distribution:

$$f_{M_c|D,m,P,\Lambda}(\mu_c|\mathcal{D},m,P,\Lambda) = \mathcal{N}(\mu_c|m_c,P_c^{-1})$$

with

$$\begin{aligned}P_c &= P + n_c\Lambda \\m_c &= P_c^{-1} \left( \Lambda \sum_i x_{c,i} + Pm \right)\end{aligned}$$

The class-conditional predictive distribution can be computed in closed form as

$$\begin{aligned} f_{X_t|C_t,D,m,B,\Lambda}(\mathbf{x}_t|c, \mathcal{D}, m, B, \Lambda) \\ &= \int f_{X_t|C_t,M_c,\Lambda}(\mathbf{x}_t|c, \mu, \Lambda) f_{M_c|D,m,B,\Lambda}(\mu|\mathcal{D}, \mu, B, \Lambda) d\mu \\ &= \int \mathcal{N}(\mathbf{x}_t|\mu, \Lambda^{-1}) \mathcal{N}(\mu|m_c, P_c^{-1}) d\mu \\ &= \mathcal{N}(\mathbf{x}_t|m_c, \Lambda^{-1} + P_c^{-1}) \end{aligned}$$

# Bayesian models

The model we derived depends on the values of the likelihood parameter  $\Lambda$  and the prior parameters  $m, P$

A Bayesian treatment of  $\Lambda$  would make the problem much harder, therefore we adopt, in this example, a frequentist approach for the estimation of  $\Lambda$

We still have to specify the prior parameters  $m, P$

Rather than fixing the values of  $m$  and  $P$ , we follow an **empirical Bayes** approach and we try to estimate the model parameters **from the data** (an alternative would be to use improper priors, or use hyper-priors to model also the parameters  $m$  and  $P$  in a Bayesian way)

We consider again Maximum Likelihood estimates



We compute  $\Lambda$ ,  $m$  and  $P$  as the maximizer of the likelihood of the **observed data**, i.e. after we have marginalized over the **latent** factors  $M_c$

$$\begin{aligned}\mathcal{L}_D(m, P, \Lambda) &= f_{D|m, P, \Lambda}(\mathcal{D}|m, B, \Lambda) \\ &= \prod_{c=1}^K \int \prod_{i=1}^{n_c} \mathcal{N}(x_{c,i}|\mu, \Lambda^{-1}) \mathcal{N}(\mu|m, P^{-1}) d\mu\end{aligned}$$

Direct maximization of the likelihood is not straightforward. However, the complete data likelihood has a simple expression

$$\begin{aligned}\mathcal{L}_{D,M}(m, P, \Lambda) &= f_{D,M|m, P, \Lambda}(\mathcal{D}, \mu_1 \dots \mu_K|m, P, \Lambda) \\ &= \prod_{c=1}^K \left[ \left( \prod_{i=1}^{n_c} \mathcal{N}(x_{c,i}|\mu_c, \Lambda^{-1}) \right) \mathcal{N}(\mu_c|m, P^{-1}) \right]\end{aligned}$$

In log-form:

$$\begin{aligned}\ell_{D,M}(\mathbf{m}, \mathbf{P}, \Lambda) &= \log f_{D,M|\mathbf{m},\mathbf{P},\Lambda}(\mathcal{D}, \mu_1 \dots \mu_K | \mathbf{m}, \mathbf{P}, \Lambda) \\ &= \sum_{c=1}^K \left[ \left( \sum_{i=1}^{n_c} \log \mathcal{N}(\mathbf{x}_{c,i} | \mu_c, \Lambda^{-1}) \right) + \log \mathcal{N}(\mu_c | \mathbf{m}, \mathbf{P}^{-1}) \right]\end{aligned}$$

We can resort to the **EM algorithm** to perform iterative maximizations of the expected complete data log-likelihood

The EM latent variables are the class means  $\mathbf{M}_1 \dots \mathbf{M}_K$

We start from an initial estimate of the parameters — for example, we can use the data mean for  $\mathbf{m}$ , and the between and within-class covariance matrices for  $\mathbf{P}^{-1}$  and  $\Lambda^{-1}$ , respectively

In the E-step, we need to compute the **posterior distributions** for  $M_1 \dots M_K$ , **conditioned on the current estimate** of the parameters

$$\begin{aligned} f_{M_1 \dots M_K | D, m_t, P_t, \Lambda_t}(\mu_1 \dots \mu_K | \mathcal{D}, m_t, P_t, \Lambda_t) \\ = \prod_{c=1}^K f_{M_c | D_c, m_t, P_t, \Lambda_t}(\mu_c | \mathcal{D}_c, m_t, P_t, \Lambda_t) \end{aligned}$$

As we have seen, each term  $f_{M_c | D_c, m_t, P_t, \Lambda_t}(\mu_c | \mathcal{D}_c, m_t, P_t, \Lambda_t)$  corresponds to a Gaussian density

$$f_{M_c | D, m_t, P_t, \Lambda_t}(\mu_c | \mathcal{D}, m_t, P_t, \Lambda_t) = \mathcal{N}(\mu_c | m_{t,c}, P_{t,c}^{-1})$$

with

$$\begin{aligned} P_{t,c} &= P_t + n_c \Lambda_t \\ m_{t,c} &= P_{t,c}^{-1} \left( \Lambda_t \sum_i x_{c,i} + P_t m_t \right) \end{aligned}$$

The **EM auxiliary function** corresponds to the expectation of the complete-data log-likelihood with respect to the posterior distribution of the latent variables

$$\begin{aligned} \mathcal{Q}(\mathbf{m}, \mathbf{P}, \mathbf{\Lambda}, \mathbf{m}_t, \mathbf{P}_t, \mathbf{\Lambda}_t) \\ = \mathbb{E}_{\mathbf{M}_1 \dots \mathbf{M}_K | \mathcal{D} = \mathcal{D}, \mathbf{m}_t, \mathbf{P}_t, \mathbf{\Lambda}_t} [\log f_{\mathcal{D}, \mathbf{M}_1 \dots \mathbf{M}_K | \mathbf{m}, \mathbf{P}, \mathbf{\Lambda}}(\mathcal{D}, \boldsymbol{\mu}_1 \dots \boldsymbol{\mu}_K | \mathbf{m}, \mathbf{P}, \mathbf{\Lambda})] \end{aligned}$$

Since the likelihood and the posterior distribution factorize over classes, we can show that the auxiliary function can be expressed as

$$\begin{aligned} \mathcal{Q}(\mathbf{m}, \mathbf{P}, \mathbf{\Lambda}, \mathbf{m}_t, \mathbf{P}_t, \mathbf{\Lambda}_t) \\ = \sum_{c=1}^K \mathbb{E}_{\mathbf{M}_c | \mathcal{D}_c = \mathcal{D}_c, \mathbf{m}_t, \mathbf{P}_t, \mathbf{\Lambda}_t} [\log f_{\mathcal{D}_c, \mathbf{M}_c | \mathbf{m}, \mathbf{P}, \mathbf{\Lambda}}(\mathcal{D}_c, \boldsymbol{\mu}_c | \mathbf{m}, \mathbf{P}, \mathbf{\Lambda})] \\ = \sum_{c=1}^K \mathbb{E}_{\mathbf{M}_c | \mathcal{D}_c = \mathcal{D}_c, \mathbf{m}_t, \mathbf{P}_t, \mathbf{\Lambda}_t} \left[ \left( \sum_{i=1}^{n_c} \log \mathcal{N}(\mathbf{x}_{c,i} | \boldsymbol{\mu}_c, \mathbf{\Lambda}^{-1}) \right) + \log \mathcal{N}(\boldsymbol{\mu}_c | \mathbf{m}, \mathbf{P}^{-1}) \right] \end{aligned}$$

Ignoring all terms that do not depend on the parameters, and denoting the expectations simply as  $\mathbb{E}[\cdot]$ , we have

$$\begin{aligned} \mathcal{Q}(\mathbf{m}, \mathbf{P}, \mathbf{\Lambda}, \mathbf{m}_t, \mathbf{P}_t, \mathbf{\Lambda}_t) &= \frac{1}{2} \sum_{c=1}^K \mathbb{E} \left[ n_c |\mathbf{\Lambda}| - \text{Tr}(\mathbf{\Lambda} \mathbf{S}_c) + 2 \text{Tr}(\mathbf{\Lambda} \mathbf{F}_c \boldsymbol{\mu}_c^T) - \text{Tr}(n_c \mathbf{\Lambda} \boldsymbol{\mu}_c \boldsymbol{\mu}_c^T) \right. \\ &\quad \left. + \log |\mathbf{P}| - \text{Tr}(\mathbf{P} \boldsymbol{\mu}_c \boldsymbol{\mu}_c^T) + 2 \text{Tr}(\mathbf{m}^T \mathbf{P} \boldsymbol{\mu}_c) - \text{Tr}(\mathbf{P} \mathbf{m} \mathbf{m}^T) \right] \\ &= \frac{1}{2} n \log |\mathbf{\Lambda}| - \text{Tr}(\mathbf{S} \mathbf{\Lambda}) + 2 \text{Tr} \left( \mathbf{\Lambda} \sum_{c=1}^K \mathbf{F}_c \mathbb{E} [\boldsymbol{\mu}_c]^T \right) \\ &\quad - \text{Tr} \left( \mathbf{\Lambda} \sum_{c=1}^K n_c \mathbb{E} [\boldsymbol{\mu}_c \boldsymbol{\mu}_c^T] \right) + K \log |\mathbf{P}| - \text{Tr} \left( \mathbf{P} \sum_{c=1}^K \mathbb{E} [\boldsymbol{\mu}_c \boldsymbol{\mu}_c^T] \right) \\ &\quad + 2 \text{Tr} \left( \mathbf{m}^T \mathbf{P} \sum_{c=1}^K \mathbb{E} [\boldsymbol{\mu}_c] \right) - K \text{Tr} (\mathbf{P} \mathbf{m} \mathbf{m}^T) \end{aligned}$$

The terms  $\mathbf{S}$ ,  $\mathbf{S}_c$  and  $\mathbf{F}_c$  are the statistics

$$\mathbf{F}_c = \sum_{i=1}^{n_c} \mathbf{x}_{c,i}, \quad \mathbf{S}_c = \sum_{i=1}^{n_c} \mathbf{x}_{c,i} \mathbf{x}_{c,i}^T, \quad \mathbf{S} = \sum_{c=1}^K \mathbf{S}_c$$

Setting the derivatives w.r.t  $\mathbf{m}$ ,  $\mathbf{P}$ ,  $\mathbf{\Lambda}$  equal to zero we obtain

$$\mathbf{m}_{t+1} = \frac{1}{K} \sum_{c=1}^K \mathbb{E} [\boldsymbol{\mu}_c]$$

$$\mathbf{P}_{t+1}^{-1} = \frac{1}{K} \left( \sum_{c=1}^K \mathbb{E} [\boldsymbol{\mu}_c \boldsymbol{\mu}_c^T] - 2 \sum_{c=1}^K \mathbb{E} [\boldsymbol{\mu}_c] \mathbf{m}^T + K \mathbf{m} \mathbf{m}^T \right)$$

$$\mathbf{\Lambda}_{t+1}^{-1} = \frac{1}{n} \left( \mathbf{S} - 2 \sum_{c=1}^K \mathbf{F}_c \mathbb{E} [\boldsymbol{\mu}_c]^T - \sum_{c=1}^K n_c \mathbb{E} [\boldsymbol{\mu}_c \boldsymbol{\mu}_c^T] \right)$$

The expectations required to compute the updated model parameters can be obtained from the Gaussian density expectations. Since  $\mathbb{E} [\boldsymbol{\mu} \boldsymbol{\mu}^T] = \mathbb{E} [\boldsymbol{\mu}] \mathbb{E} [\boldsymbol{\mu}]^T + \text{Cov}(\boldsymbol{\mu})$ , for class mean  $\boldsymbol{M}_c$ , we have

$$\begin{aligned}\mathbb{E}_{\boldsymbol{M}_c | \boldsymbol{D}_c = \boldsymbol{\mathcal{D}}_c, \boldsymbol{m}_t, \boldsymbol{P}_t, \boldsymbol{\Lambda}_t} [\boldsymbol{\mu}_c] &= \boldsymbol{m}_{t,c} \\ \mathbb{E}_{\boldsymbol{M}_c | \boldsymbol{D}_c = \boldsymbol{\mathcal{D}}_c, \boldsymbol{m}_t, \boldsymbol{P}_t, \boldsymbol{\Lambda}_t} [\boldsymbol{\mu}_c \boldsymbol{\mu}_c^T] &= \boldsymbol{P}_{t,c}^{-1} + \boldsymbol{m}_{t,c} \boldsymbol{m}_{t,c}^T\end{aligned}$$

and we can compute the parameters update as in the previous slides

We should note that, to get a robust estimate of the matrix  $\boldsymbol{P}$ , we require a large number of classes

# Bayesian models

We can generalize the model assuming that a sample  $\mathbf{x}_{c,i}$  of class  $c$  corresponds to the sum of three terms

$$\mathbf{X}_{c,i} = \mathbf{m} + \mathbf{U}\mathbf{Y}_c + \mathbf{E}_{c,i}$$

where  $\mathbf{U}$  is a rectangular matrix

The prior for  $\mathbf{Y}_c$  is assumed to be

$$f_{\mathbf{Y}_c}(\mathbf{y}_c) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

and the noise terms  $\mathbf{E}_{c,i}$  are i.i.d distributed as

$$f_{\mathbf{E}_{c,i}|\mathbf{\Lambda}}(\varepsilon_{c,i}) = \mathcal{N}(\varepsilon_{c,i}|\mathbf{0}, \mathbf{\Lambda}^{-1})$$

Matrix  $\mathbf{U}$  can be used to constrain the model means to lie in a subspace of the original data

The model is fully specified even for rank-deficient matrix  $\mathbf{U}$ , and more robust estimates can be obtained when the number of classes is not large enough for the estimation of a full rank matrix  $\mathbf{P}$



# Bayesian models

In this case, the conditional likelihood for a sample of class  $c$  corresponds to

$$f_{X_{c,i}|Y_c,m,U,\Lambda}(\mathbf{x}_{c,i}|\mathbf{y}_c, \mathbf{m}, \mathbf{U}, \Lambda) = \mathcal{N}(\mathbf{x}_{c,i}|\mathbf{m} + \mathbf{U}\mathbf{y}_{c,i}, \Lambda^{-1})$$

The derivations are similar to the previous case

The model is also known as (simplified) **Probabilistic Linear Discriminant Analysis**

The model we have previously seen, and the PLDA generalization, are effective when we have **many classes** (enough to estimate a between-class covariance matrix  $\mathbf{P}^{-1}$ , or the rectangular matrix  $\mathbf{U}$ ), but **few samples per class**, so that we cannot reliably estimate class means

For example, this is often the case in biometric recognition systems (e.g. face, speaker or fingerprint recognition)

# Bayesian models

The Bayesian model allows also addressing multi-class recognition problems

Indeed, we can model samples from the **none-of-the-others** class as belonging to a class  $\mathcal{H}_N$  for which we have **no observations**

If we are using a **proper prior** (e.g. the one estimated with the empirical Bayes approach), we can still compute the predictive distribution for this case

The posterior distribution we use to compute the predictive density will be the **prior** distribution. For the Gaussian example:

$$\begin{aligned} f_{X_t|C}(\mathbf{x}_t|\mathcal{H}_N) &= \int f_{X_t|M}(\mathbf{x}_t|\boldsymbol{\mu})f_M(\boldsymbol{\mu})d\boldsymbol{\mu} = \int \mathcal{N}(\mathbf{x}_t|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \mathcal{N}(\boldsymbol{\mu}|\mathbf{m}, \mathbf{P}^{-1}) d\boldsymbol{\mu} \\ &= \mathcal{N}(\mathbf{x}_t|\mathbf{m}, \boldsymbol{\Lambda}^{-1} + \mathbf{P}^{-1}) \end{aligned}$$

The prior encodes the prior probability for a class mean. Since we have no observations for the class the sample belongs to, our beliefs on the corresponding class mean are encoded by the prior

# Bayesian models

PLDA is an example of a more general class of models, which are known as Factor Analysis

In general, FA models express observed samples as a combination of different, Gaussian distributed factors

$$\mathbf{x} = \mathbf{m} + \mathbf{U}\mathbf{y} + \mathbf{V}\mathbf{w} + \dots + \boldsymbol{\varepsilon}$$

Each term  $\mathbf{y}, \mathbf{w}, \dots$  is a **latent factor** that models some characteristics of the data

The matrices that constrain factors to be small-dimensional are also called **factor loading** matrices

Usually, the **residual noise** term  $\boldsymbol{\varepsilon}$  is assumed to have diagonal covariance matrix (PLDA is an exception, but we can represent PLDA with an alternative form which satisfies this assumption)

The factors are thus used to capture correlations between observation components

# Gaussian Linear Models

Let  $[x_1 \dots x_n]$  be a set of (unlabeled) data. A particular instance of FA that has several practical applications consists in the model

$$x_i = m + Wy_i + \varepsilon_i$$

with

$$\varepsilon_i \sim E_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$$

$$y_i \sim Y_i \sim \mathcal{N}(\mathbf{0}, I)$$

i.e., a single factor model with isotropic noise

The conditional likelihood is given by

$$f_{X_i|Y_i}(x_i|y_i) = \mathcal{N}(x_i|m + Wy_i, \sigma^2 I)$$

As for PLDA, we can compute the *posterior* distribution:

$$f_{Y_i|X_i}(\mathbf{y}_i|\mathbf{x}_i) = \mathcal{N}(\mathbf{y}_i|\boldsymbol{\mu}_{x,i}, \boldsymbol{\Lambda}_x)$$

with

$$\begin{aligned}\boldsymbol{\Lambda}_x &= \frac{1}{\sigma^2} (\mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I}) \\ \boldsymbol{\mu}_{x,i} &= (\mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I})^{-1} \mathbf{W}^T (\mathbf{x}_i - \mathbf{m})\end{aligned}$$

The marginal (predictive) distribution is

$$f_{X_i}(\mathbf{x}_i) = \mathcal{N}(\mathbf{x}_i|\mathbf{m}, \mathbf{C})$$

with  $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$

# Gaussian Linear Models

As we did for PLDA, we can use the EM algorithm for estimating the parameters of the model

However, for this particular model closed form solutions can be given for the model parameters by direct maximization of the observed data log-likelihood

$$\sum_{i=1}^n -\frac{1}{2} \log |\mathbf{C}| - \frac{1}{2} (\mathbf{x}_i - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x}_i - \mathbf{m})$$

It can be shown that a ML solution for the model is given by

$$\mathbf{m} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \quad \mathbf{W} = \mathbf{U}(\mathbf{L} - \sigma^2 \mathbf{I})^{\frac{1}{2}} \mathbf{R}, \quad \sigma^2 = \frac{1}{D - M} \sum_{i=M+1}^D \lambda_i$$

where  $D$  and  $M$  are the dimensions of the feature and latent ( $\mathbf{y}$ ) spaces, respectively

# Gaussian Linear Models

$U$  is a  $D \times M$  matrix whose columns are the first  $M$  eigenvectors of the data covariance matrix

$L$  is the diagonal matrix containing the corresponding eigenvalues  $\lambda_i$  for  $i \leq M$

$R$  is an arbitrary orthogonal matrix

The ML solution is not unique. In particular, we can arrange the column vectors of  $U$  so that the corresponding eigenvalues are in decreasing order

Furthermore, different choices of  $R$  correspond to the same predictive distribution<sup>2</sup>

Choosing  $R = I$ , the ML solution for  $W$  is then the matrix containing the first  $M$  eigenvectors of the covariance matrix, scaled by  $(\lambda_i - \sigma^2)^{\frac{1}{2}}$

---

<sup>2</sup>The model is not identifiable

Matrix  $W$  thus spans the same subspace as Principal Component Analysis

Indeed, the model is known as **Probabilistic Principal Component Analysis**

If we assume that  $\sigma$  is known, and we take the limit  $\sigma \rightarrow 0$ , then the posterior distribution mean becomes

$$\mu_x = (W^T W)^{-1} W^T (x - m)$$

corresponding to the PCA projection

An interesting property of the model is that we can still apply the EM algorithm even in the limit case  $\sigma \rightarrow 0$



The corresponding algorithm is known as EM-PCA

EM-PCA allows computing the PCA projection for problems characterized by both a large number of samples and large sample dimensionality, for which the computation and decomposition of the data covariance matrix is unfeasible

We have seen a brief introduction to Bayesian models

The main limitations of Bayesian approaches derive from the need to specify a prior (but we have seen some possible ways to select prior distributions) and to the difficulty of computing integrals for generic prior / likelihood combinations

The latter issue can be mitigated by using conjugate priors, however even simple model can quickly become intractable

There exist approximate methods to handle complex posteriors, that are based on devising an approximation of the posterior distribution which makes computations tractable

If we are able to compute the posterior mode and the corresponding Hessian matrix, we can use the Laplace approximation to approximate the posterior by a Gaussian distribution

We may also specify the approximate posterior distribution as belonging to a parametric family, and look for the parameters that best approximate the intractable posterior (e.g. Variational Bayes methods)

When dealing with multiple latent variables, we may assume an approximate posterior distribution that factorizes over different latent variables even in those cases where the latent variables are, a posteriori, dependent (mean-field Variational Bayes)

For further information you can refer to the course textbooks