# CS229 Lecture notes

#### Andrew Ng

# Part X Factor analysis

When we have data  $x^{(i)} \in \mathbb{R}^n$  that comes from a mixture of several Gaussians, the EM algorithm can be applied to fit a mixture model. In this setting, we usually imagine problems where we have sufficient data to be able to discern the multiple-Gaussian structure in the data. For instance, this would be the case if our training set size m was significantly larger than the dimension nof the data.

Now, consider a setting in which  $n \gg m$ . In such a problem, it might be difficult to model the data even with a single Gaussian, much less a mixture of Gaussian. Specifically, since the m data points span only a low-dimensional subspace of  $\mathbb{R}^n$ , if we model the data as Gaussian, and estimate the mean and covariance using the usual maximum likelihood estimators,

$$\mu = \frac{1}{m} \sum_{i=1}^{m} x^{(i)}$$
  

$$\Sigma = \frac{1}{m} \sum_{i=1}^{m} (x^{(i)} - \mu) (x^{(i)} - \mu)^{T},$$

we would find that the matrix  $\Sigma$  is singular. This means that  $\Sigma^{-1}$  does not exist, and  $1/|\Sigma|^{1/2} = 1/0$ . But both of these terms are needed in computing the usual density of a multivariate Gaussian distribution. Another way of stating this difficulty is that maximum likelihood estimates of the parameters result in a Gaussian that places all of its probability in the affine space spanned by the data,<sup>1</sup> and this corresponds to a singular covariance matrix.

<sup>&</sup>lt;sup>1</sup>This is the set of points x satisfying  $x = \sum_{i=1}^{m} \alpha_i x^{(i)}$ , for some  $\alpha_i$ 's so that  $\sum_{i=1}^{m} \alpha_1 = 1$ .

More generally, unless m exceeds n by some reasonable amount, the maximum likelihood estimates of the mean and covariance may be quite poor. Nonetheless, we would still like to be able to fit a reasonable Gaussian model to the data, and perhaps capture some interesting covariance structure in the data. How can we do this?

In the next section, we begin by reviewing two possible restrictions on  $\Sigma$ , ones that allow us to fit  $\Sigma$  with small amounts of data but neither of which will give a satisfactory solution to our problem. We next discuss some properties of Gaussians that will be needed later; specifically, how to find marginal and conditonal distributions of Gaussians. Finally, we present the factor analysis model, and EM for it.

#### **1** Restrictions of $\Sigma$

If we do not have sufficient data to fit a full covariance matrix, we may place some restrictions on the space of matrices  $\Sigma$  that we will consider. For instance, we may choose to fit a covariance matrix  $\Sigma$  that is diagonal. In this setting, the reader may easily verify that the maximum likelihood estimate of the covariance matrix is given by the diagonal matrix  $\Sigma$  satisfying

$$\Sigma_{jj} = \frac{1}{m} \sum_{i=1}^{m} (x_j^{(i)} - \mu_j)^2$$

Thus,  $\Sigma_{jj}$  is just the empirical estimate of the variance of the *j*-th coordinate of the data.

Recall that the contours of a Gaussian density are ellipses. A diagonal  $\Sigma$  corresponds to a Gaussian where the major axes of these ellipses are axisaligned.

Sometimes, we may place a further restriction on the covariance matrix that not only must it be diagonal, but its diagonal entries must all be equal. In this setting, we have  $\Sigma = \sigma^2 I$ , where  $\sigma^2$  is the parameter under our control. The maximum likelihood estimate of  $\sigma^2$  can be found to be:

$$\sigma^2 = \frac{1}{mn} \sum_{j=1}^n \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2.$$

This model corresponds to using Gaussians whose densities have contours that are circles (in 2 dimensions; or spheres/hyperspheres in higher dimensions). If we were fitting a full, unconstrained, covariance matrix  $\Sigma$  to data, it was necessary that  $m \ge n+1$  in order for the maximum likelihood estimate of  $\Sigma$  not to be singular. Under either of the two restrictions above, we may obtain non-singular  $\Sigma$  when  $m \ge 2$ .

However, restricting  $\Sigma$  to be diagonal also means modeling the different coordinates  $x_i$ ,  $x_j$  of the data as being uncorrelated and independent. Often, it would be nice to be able to capture some interesting correlation structure in the data. If we were to use either of the restrictions on  $\Sigma$  described above, we would therefore fail to do so. In this set of notes, we will describe the factor analysis model, which uses more parameters than the diagonal  $\Sigma$  and captures some correlations in the data, but also without having to fit a full covariance matrix.

## 2 Marginals and conditionals of Gaussians

Before describing factor analysis, we digress to talk about how to find conditional and marginal distributions of random variables with a joint multivariate Gaussian distribution.

Suppose we have a vector-valued random variable

$$x = \left[ \begin{array}{c} x_1 \\ x_2 \end{array} \right],$$

where  $x_1 \in \mathbb{R}^r$ ,  $x_2 \in \mathbb{R}^s$ , and  $x \in \mathbb{R}^{r+s}$ . Suppose  $x \sim \mathcal{N}(\mu, \Sigma)$ , where

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}.$$

Here,  $\mu_1 \in \mathbb{R}^r$ ,  $\mu_2 \in \mathbb{R}^s$ ,  $\Sigma_{11} \in \mathbb{R}^{r \times r}$ ,  $\Sigma_{12} \in \mathbb{R}^{r \times s}$ , and so on. Note that since covariance matrices are symmetric,  $\Sigma_{12} = \Sigma_{21}^T$ .

Under our assumptions,  $x_1$  and  $x_2$  are jointly multivariate Gaussian. What is the marginal distribution of  $x_1$ ? It is not hard to see that  $E[x_1] = \mu_1$ , and that  $Cov(x_1) = E[(x_1 - \mu_1)(x_1 - \mu_1)] = \Sigma_{11}$ . To see that the latter is true, note that by definition of the joint covariance of  $x_1$  and  $x_2$ , we have that

$$Cov(x) = \Sigma$$
  
=  $\begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$   
=  $E[(x - \mu)(x - \mu)^T]$   
=  $E\left[\begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}^T\right]$   
=  $E\left[\begin{pmatrix} (x_1 - \mu_1)(x_1 - \mu_1)^T & (x_1 - \mu_1)(x_2 - \mu_2)^T \\ (x_2 - \mu_2)(x_1 - \mu_1)^T & (x_2 - \mu_2)(x_2 - \mu_2)^T \end{bmatrix}$ 

Matching the upper-left subblocks in the matrices in the second and the last lines above gives the result.

Since marginal distributions of Gaussians are themselves Gaussian, we therefore have that the marginal distribution of  $x_1$  is given by  $x_1 \sim \mathcal{N}(\mu_1, \Sigma_{11})$ .

Also, we can ask, what is the conditional distribution of  $x_1$  given  $x_2$ ? By referring to the definition of the multivariate Gaussian distribution, it can be shown that  $x_1|x_2 \sim \mathcal{N}(\mu_{1|2}, \Sigma_{1|2})$ , where

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2), \qquad (1)$$

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}.$$
 (2)

When working with the factor analysis model in the next section, these formulas for finding conditional and marginal distributions of Gaussians will be very useful.

## 3 The Factor analysis model

In the factor analysis model, we posit a joint distribution on (x, z) as follows, where  $z \in \mathbb{R}^k$  is a latent random variable:

$$z \sim \mathcal{N}(0, I)$$
  
 $x|z \sim \mathcal{N}(\mu + \Lambda z, \Psi)$ 

Here, the parameters of our model are the vector  $\mu \in \mathbb{R}^n$ , the matrix  $\Lambda \in \mathbb{R}^{n \times k}$ , and the diagonal matrix  $\Psi \in \mathbb{R}^{n \times n}$ . The value of k is usually chosen to be smaller than n.

Thus, we imagine that each datapoint  $x^{(i)}$  is generated by sampling a k dimension multivariate Gaussian  $z^{(i)}$ . Then, it is mapped to a k-dimensional affine space of  $\mathbb{R}^n$  by computing  $\mu + \Lambda z^{(i)}$ . Lastly,  $x^{(i)}$  is generated by adding covariance  $\Psi$  noise to  $\mu + \Lambda z^{(i)}$ .

Equivalently (convince yourself that this is the case), we can therefore also define the factor analysis model according to

$$\begin{aligned} z &\sim \mathcal{N}(0, I) \\ \epsilon &\sim \mathcal{N}(0, \Psi) \\ x &= \mu + \Lambda z + \epsilon \end{aligned}$$

where  $\epsilon$  and z are independent.

Let's work out exactly what distribution our model defines. Our random variables z and x have a joint Gaussian distribution

$$\left[\begin{array}{c}z\\x\end{array}\right] \sim \mathcal{N}(\mu_{zx}, \Sigma).$$

We will now find  $\mu_{zx}$  and  $\Sigma$ .

We know that  $\mathbf{E}[z] = 0$ , from the fact that  $z \sim \mathcal{N}(0, I)$ . Also, we have that

$$E[x] = E[\mu + \Lambda z + \epsilon]$$
  
=  $\mu + \Lambda E[z] + E[\epsilon]$   
=  $\mu$ .

Putting these together, we obtain

$$\mu_{zx} = \left[ \begin{array}{c} \vec{0} \\ \mu \end{array} \right]$$

Next, to find,  $\Sigma$ , we need to calculate  $\Sigma_{zz} = \mathrm{E}[(z - \mathrm{E}[z])(z - \mathrm{E}[z])^T]$  (the upper-left block of  $\Sigma$ ),  $\Sigma_{zx} = \mathrm{E}[(z - \mathrm{E}[z])(x - \mathrm{E}[x])^T]$  (upper-right block), and  $\Sigma_{xx} = \mathrm{E}[(x - \mathrm{E}[x])(x - \mathrm{E}[x])^T]$  (lower-right block).

Now, since  $z \sim \mathcal{N}(0, I)$ , we easily find that  $\Sigma_{zz} = \text{Cov}(z) = I$ . Also,

$$E[(z - E[z])(x - E[x])^T] = E[z(\mu + \Lambda z + \epsilon - \mu)^T]$$
  
=  $E[zz^T]\Lambda^T + E[z\epsilon^T]$   
=  $\Lambda^T$ .

In the last step, we used the fact that  $E[zz^T] = Cov(z)$  (since z has zero mean), and  $E[z\epsilon^T] = E[z]E[\epsilon^T] = 0$  (since z and  $\epsilon$  are independent, and

hence the expectation of their product is the product of their expectations). Similarly, we can find  $\Sigma_{xx}$  as follows:

$$E[(x - E[x])(x - E[x])^{T}] = E[(\mu + \Lambda z + \epsilon - \mu)(\mu + \Lambda z + \epsilon - \mu)^{T}]$$
  
$$= E[\Lambda z z^{T} \Lambda^{T} + \epsilon z^{T} \Lambda^{T} + \Lambda z \epsilon^{T} + \epsilon \epsilon^{T}]$$
  
$$= \Lambda E[z z^{T}] \Lambda^{T} + E[\epsilon \epsilon^{T}]$$
  
$$= \Lambda \Lambda^{T} + \Psi.$$

Putting everything together, we therefore have that

$$\begin{bmatrix} z \\ x \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \vec{0} \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda\Lambda^T + \Psi \end{bmatrix}\right).$$
(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  $\{x^{(i)}; i = 1, \ldots, m\}$ , we can write down the log likelihood of the parameters:

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

To perform maximum likelihood estimation, we would like to maximize this quantity with respect to the parameters. But maximizing this formula explicitly is hard (try it yourself), and we are aware of no algorithm that does so in closed-form. So, we will instead use to the EM algorithm. In the next section, we derive EM for factor analysis.

#### 4 EM for factor analysis

The derivation for the E-step is easy. We need to compute  $Q_i(z^{(i)}) = p(z^{(i)}|x^{(i)}; \mu, \Lambda, \Psi)$ . By substituting the distribution given in Equation (3) into the formulas (1-2) used for finding the conditional distribution of a Gaussian, we find that  $z^{(i)}|x^{(i)}; \mu, \Lambda, \Psi \sim \mathcal{N}(\mu_{z^{(i)}|x^{(i)}}, \Sigma_{z^{(i)}|x^{(i)}})$ , where

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

So, using these definitions for  $\mu_{z^{(i)}|x^{(i)}}$  and  $\Sigma_{z^{(i)}|x^{(i)}}$ , we have

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

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

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

with respect to the parameters  $\mu$ ,  $\Lambda$ ,  $\Psi$ . We will work out only the optimization with respect to  $\Lambda$ , and leave the derivations of the updates for  $\mu$  and  $\Psi$ as an exercise to the reader.

We can simplify Equation (4) as follows:

$$\sum_{i=1}^{m} \int_{z^{(i)}} Q_i(z^{(i)}) \left[ \log p(x^{(i)} | z^{(i)}; \mu, \Lambda, \Psi) + \log p(z^{(i)}) - \log Q_i(z^{(i)}) \right] dz^{(i)}$$
(5)  
= 
$$\sum_{i=1}^{m} E_{z^{(i)} \sim Q_i} \left[ \log p(x^{(i)} | z^{(i)}; \mu, \Lambda, \Psi) + \log p(z^{(i)}) - \log Q_i(z^{(i)}) \right]$$
(6)

Here, the " $z^{(i)} \sim Q_i$ " subscript indicates that the expectation is with respect to  $z^{(i)}$  drawn from  $Q_i$ . In the subsequent development, we will omit this subscript when there is no risk of ambiguity. Dropping terms that do not depend on the parameters, we find that we need to maximize:

$$\begin{split} &\sum_{i=1}^{m} \mathbf{E} \left[ \log p(x^{(i)} | z^{(i)}; \mu, \Lambda, \Psi) \right] \\ &= \sum_{i=1}^{m} \mathbf{E} \left[ \log \frac{1}{(2\pi)^{n/2} |\Psi|^{1/2}} \exp \left( -\frac{1}{2} (x^{(i)} - \mu - \Lambda z^{(i)})^T \Psi^{-1} (x^{(i)} - \mu - \Lambda z^{(i)}) \right) \right] \\ &= \sum_{i=1}^{m} \mathbf{E} \left[ -\frac{1}{2} \log |\Psi| - \frac{n}{2} \log(2\pi) - \frac{1}{2} (x^{(i)} - \mu - \Lambda z^{(i)})^T \Psi^{-1} (x^{(i)} - \mu - \Lambda z^{(i)}) \right] \end{split}$$

Let's maximize this with respect to  $\Lambda$ . Only the last term above depends on  $\Lambda$ . Taking derivatives, and using the facts that tr a = a (for  $a \in \mathbb{R}$ ), trAB = trBA, and  $\nabla_A \text{tr}ABA^TC = CAB + C^TAB$ , we get:

$$\nabla_{\Lambda} \sum_{i=1}^{m} -E \left[ \frac{1}{2} (x^{(i)} - \mu - \Lambda z^{(i)})^{T} \Psi^{-1} (x^{(i)} - \mu - \Lambda z^{(i)}) \right]$$
  
= 
$$\sum_{i=1}^{m} \nabla_{\Lambda} E \left[ -\operatorname{tr} \frac{1}{2} z^{(i)^{T}} \Lambda^{T} \Psi^{-1} \Lambda z^{(i)} + \operatorname{tr} z^{(i)^{T}} \Lambda^{T} \Psi^{-1} (x^{(i)} - \mu) \right]$$
  
= 
$$\sum_{i=1}^{m} \nabla_{\Lambda} E \left[ -\operatorname{tr} \frac{1}{2} \Lambda^{T} \Psi^{-1} \Lambda z^{(i)} z^{(i)^{T}} + \operatorname{tr} \Lambda^{T} \Psi^{-1} (x^{(i)} - \mu) z^{(i)^{T}} \right]$$
  
= 
$$\sum_{i=1}^{m} E \left[ -\Psi^{-1} \Lambda z^{(i)} z^{(i)^{T}} + \Psi^{-1} (x^{(i)} - \mu) z^{(i)^{T}} \right]$$

Setting this to zero and simplifying, we get:

$$\sum_{i=1}^{m} \Lambda \mathbf{E}_{z^{(i)} \sim Q_{i}} \left[ z^{(i)} z^{(i)^{T}} \right] = \sum_{i=1}^{m} (x^{(i)} - \mu) \mathbf{E}_{z^{(i)} \sim Q_{i}} \left[ z^{(i)^{T}} \right].$$

Hence, solving for  $\Lambda$ , we obtain

$$\Lambda = \left(\sum_{i=1}^{m} (x^{(i)} - \mu) \mathbf{E}_{z^{(i)} \sim Q_i} \left[ z^{(i)^T} \right] \right) \left(\sum_{i=1}^{m} \mathbf{E}_{z^{(i)} \sim Q_i} \left[ z^{(i)} z^{(i)^T} \right] \right)^{-1}.$$
 (7)

It is interesting to note the close relationship between this equation and the normal equation that we'd derived for least squares regression,

$$"\theta^T = (y^T X)(X^T X)^{-1}."$$

The analogy is that here, the x's are a linear function of the z's (plus noise). Given the "guesses" for z that the E-step has found, we will now try to estimate the unknown linearity  $\Lambda$  relating the x's and z's. It is therefore no surprise that we obtain something similar to the normal equation. There is, however, one important difference between this and an algorithm that performs least squares using just the "best guesses" of the z's; we will see this difference shortly.

To complete our M-step update, let's work out the values of the expectations in Equation (7). From our definition of  $Q_i$  being Gaussian with mean  $\mu_{z^{(i)}|x^{(i)}}$  and covariance  $\Sigma_{z^{(i)}|x^{(i)}}$ , we easily find

The latter comes from the fact that, for a random variable Y,  $\operatorname{Cov}(Y) = \operatorname{E}[YY^T] - \operatorname{E}[Y]\operatorname{E}[Y]^T$ , and hence  $\operatorname{E}[YY^T] = \operatorname{E}[Y]\operatorname{E}[Y]^T + \operatorname{Cov}(Y)$ . Substituting this back into Equation (7), we get the M-step update for  $\Lambda$ :

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

It is important to note the presence of the  $\sum_{z^{(i)}|x^{(i)}}$  on the right hand side of this equation. This is the covariance in the posterior distribution  $p(z^{(i)}|x^{(i)})$ of  $z^{(i)}$  give  $x^{(i)}$ , and the M-step must take into account this uncertainty about  $z^{(i)}$  in the posterior. A common mistake in deriving EM is to assume that in the E-step, we need to calculate only expectation E[z] of the latent random variable z, and then plug that into the optimization in the M-step everywhere z occurs. While this worked for simple problems such as the mixture of Gaussians, in our derivation for factor analysis, we needed  $E[zz^T]$ as well E[z]; and as we saw,  $E[zz^T]$  and  $E[z]E[z]^T$  differ by the quantity  $\Sigma_{z|x}$ . Thus, the M-step update must take into account the covariance of z in the posterior distribution  $p(z^{(i)}|x^{(i)})$ .

Lastly, we can also find the M-step optimizations for the parameters  $\mu$  and  $\Psi$ . It is not hard to show that the first is given by

$$\mu = \frac{1}{m} \sum_{i=1}^{m} x^{(i)}.$$

Since this doesn't change as the parameters are varied (i.e., unlike the update for  $\Lambda$ , the right hand side does not depend on  $Q_i(z^{(i)}) = p(z^{(i)}|x^{(i)}; \mu, \Lambda, \Psi)$ , which in turn depends on the parameters), this can be calculated just once and needs not be further updated as the algorithm is run. Similarly, the diagonal  $\Psi$  can be found by calculating

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

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