

Visualisation and Dimension Reduction

Christophe Ambroise 

christophe.ambroise@univ-evry.fr

UEVE, UMR CNRS 8071

November 7, 2023

Factor Analyser

Factor Analysis

- Using discrete latent variables provides limited summary (clustering)
- An alternative is to use a vector of real-valued latent variables, $\mathbf{z} \in \mathbb{R}^L$.
- “Factor analysis (FA) is a statistical method used to describe variability among observed, correlated variables in terms of a potentially lower number of unobserved variables called factors.” Wikipedia quote.
- PCA and FA are related, but not identical.

The model of factor analysis

We consider the observation $\mathbf{x} \in \mathbb{R}^D$

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}$$

where

- the noise $\boldsymbol{\epsilon} \sim \mathcal{N}_D(\mathbf{0}, \boldsymbol{\Psi})$
- the hidden (latent) vector $\mathbf{z} \sim \mathcal{N}_L(\mathbf{0}, \mathbf{I}_L)$

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

the mean is a linear function of the (hidden) inputs

- \mathbf{W} is a $D \times L$ matrix, known as the factor loading matrix,
- Ψ is a $D \times D$ covariance matrix that we take to be diagonal

The special case in which $\Psi = \sigma^2 \mathbf{I}$ is called probabilistic principal components analysis or PPCA.

Reminder: Joint and conditional Gaussian distribution (see Murphy chapter 4)

Let us recall that if

- $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_{zz})$ and $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx})$ and $\text{cov}(\mathbf{z}, \mathbf{x}) = \boldsymbol{\Sigma}_{zx}$

then

$$p(\mathbf{z}, \mathbf{x}) = N\left(\begin{bmatrix} \mathbf{z} \\ \mathbf{x} \end{bmatrix} \mid \begin{bmatrix} \boldsymbol{\mu}_z \\ \boldsymbol{\mu}_x \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{zz} & \boldsymbol{\Sigma}_{zx} \\ \boldsymbol{\Sigma}_{xz} & \boldsymbol{\Sigma}_{xx} \end{bmatrix}\right)$$

and

$$p(\mathbf{z}, \mathbf{x}) = p(\mathbf{z}|\mathbf{x})p(\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_{z|x}, \boldsymbol{\Sigma}_{z|x})\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx})$$

where

- $\boldsymbol{\mu}_{z|x} = \boldsymbol{\mu}_z + \boldsymbol{\Sigma}_{zx}\boldsymbol{\Sigma}_{xx}^{-1}(\mathbf{x} - \boldsymbol{\mu}_x)$
- $\boldsymbol{\Sigma}_{z|x} = \boldsymbol{\Sigma}_{zz} - \boldsymbol{\Sigma}_{zx}\boldsymbol{\Sigma}_{xx}^{-1}\boldsymbol{\Sigma}_{xz}$

Marginal and posterior distribution

Marginal distribution

$$\mathbf{x} \sim \mathcal{N}_D(\boldsymbol{\mu}, \boldsymbol{\Sigma}_{xx} = \mathbf{W}\mathbf{W}^T + \boldsymbol{\Psi})$$

Posterior distribution

$$\mathbf{z}|\mathbf{x} \sim \mathcal{N}_L(\boldsymbol{\mu}_{z|x}, \boldsymbol{\Sigma}_{z|x})$$

where

- $\boldsymbol{\Sigma}_{z|x} = (\mathbf{I}_L + \mathbf{W}^T\boldsymbol{\Psi}^{-1}\mathbf{W})^{-1} = \mathbf{S}$
- $\boldsymbol{\mu}_{z|x} = \boldsymbol{\Sigma}_{z|x}\boldsymbol{\Sigma}_{xx}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = \mathbf{S}\mathbf{W}^T\boldsymbol{\Psi}^{-1}(\mathbf{x} - \boldsymbol{\mu})$

Exercise

Demonstrate the above formulas

Estimation

The mean μ

can be estimated by maximum likelihood

$$\mu_{mle} = \bar{x}$$

W and Ψ

are estimated using an EM algorithm

EM algorithm

Data

- Observed data : $\mathbf{x}_{1:n}$
- Missing (or hidden) data : $\mathbf{z}_{1:n}$

Principle

- Starting from θ^0
- At step q
 - E(xpectation) step:
 $Q(\theta, \theta^q) = E_{\mathbf{z}_{1:n}|\mathbf{x}_{1:n}}[\log P(\mathbf{x}_{1:n}, \mathbf{z}_{1:n}, \theta)]$
 - M(aximisation) step: $\theta^{q+1} = \operatorname{argmax}_{\theta} Q(\theta, \theta^q)$

EM for factor analysis

Let us assume that $\boldsymbol{\mu} = \mathbf{0}$ (centering of the \mathbf{x}_i), the complete log-likelihood is

$$\begin{aligned}\log p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \boldsymbol{\Psi}) &= \sum_i \log \mathcal{N}_L(z_i; \mathbf{0}, \mathbf{I}) + \log \mathcal{N}_D(\mathbf{x}_i; \\ &= -\frac{n}{2} \log |\mathbf{I}_L| - \frac{n}{2} \operatorname{Tr}(\hat{\boldsymbol{\Sigma}}_{zz}) \\ &\quad - \frac{n}{2} \log |\boldsymbol{\Psi}| - \frac{n}{2} \operatorname{Tr}(\hat{\boldsymbol{\Sigma}}_{xx} \boldsymbol{\Psi}^{-1}) +\end{aligned}$$

where

$$\hat{\Sigma}_{xx} = \frac{1}{n} \sum_i (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)(\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T$$

Exercise

Demonstrate the above formula

E step

The expectation of the complete log-likelihood requires

1. $\mathbb{E}_{z|x}[\mathbf{z}_i] = \mathbf{S} \mathbf{W}^T \boldsymbol{\Psi}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})$ where
 $\mathbf{S} = (\mathbf{I}_L + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1}$
2. $\mathbb{E}_{z|x}[\mathbf{z}_i \mathbf{z}_i^T] = \mathbb{E}_{z|x}[\mathbf{z}_i] \mathbb{E}_{z|x}[\mathbf{z}_i^T] + \mathbf{S}$

M step

Reminders

$$\frac{\partial(b^T a)}{\partial a} = b$$

$$\frac{\partial(a^T A a)}{\partial a} = (A + A^T)a$$

$$\frac{\partial}{\partial A} \text{tr}(BA) = B^T$$

$$\frac{\partial}{\partial A} \log |A| = (A^{-1})^T$$

$$\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$$

Thus if x is a vector

$$x^T A x = \text{tr}(x^T A x) = \text{tr}(A x x^T)$$

M step for Ψ

$$\mathbb{E}_{z|x} \left[\frac{\partial L(\mathbf{W}, \Psi)}{\partial \Psi^{-1}} \right] = \mathbb{E}_{z|x} \left[\frac{n}{2} (\Psi - \hat{\Sigma}_{xx}) \right] = \frac{n}{2} (\Psi - \mathbb{E}_{z|x}$$

where

$$\begin{aligned} \mathbb{E}_{z|x} \left[\hat{\Sigma}_{xx} \right] &= \frac{1}{n} \left(\sum_i \mathbf{x}_i \mathbf{x}_i^T + W \left(\sum_i \mathbb{E}_{z|x} [\mathbf{z}_i \mathbf{z}_i^T] \right) W^T - 2 \right. \\ &= \frac{1}{n} \left(\sum_i \mathbf{x}_i \mathbf{x}_i^T + W \left(\sum_i \mathbb{E}_{z|x} [\mathbf{z}_i \mathbf{x}_i^T] \right) - 2W \right. \\ &= \frac{1}{n} \left(\sum_i \mathbf{x}_i \mathbf{x}_i^T - W \sum_i \mathbf{E}_{z|x} [\mathbf{z}_i] \mathbf{x}_i^T \right) \end{aligned}$$

M step for \mathbf{W}

$$\mathbb{E}_{z|x} \left[\frac{\partial L(\mathbf{W}, \Psi)}{\partial \mathbf{W}} \right] = \mathbb{E}_{z|x} \left[-\Psi^{-1} \sum_i \mathbf{x}_i \mathbf{z}_i^T + \Psi^{-1} \mathbf{W} \sum_i \mathbf{z}_i \right]$$

M Step summary

Loading matrix

$$\mathbf{W}^{q+1} = \left(\sum_i (\mathbf{x}_i - \bar{\mathbf{x}}) \mathbb{E}_{z|x} [\mathbf{z}_i]^T \right) \left(\sum_i \mathbb{E}_{z|x} [\mathbf{z}_i \mathbf{z}_i^T] \right)^{-1}$$

Noise covariance matrix

$$\mathbf{\Psi}^{q+1} = \frac{1}{N} \text{diag} \left\{ \sum_i \mathbf{x}_i \mathbf{x}_i^T - \mathbf{W}^{q+1} \mathbb{E}_{z|x} [\mathbf{z}_i] \mathbf{x}_i^T \right\}$$

Log-likelihood

The log-likelihood can be computed using the EM decomposition

$$\log P(X; \Theta) = E_{Z_{1:n} | \mathbf{x}_{1:n}} [\log P(\mathbf{x}_{1:n}, \mathbf{z}_{1:n}; \theta)] - E_{Z_{1:n} | \mathbf{x}_{1:n}} [\log$$

Implementation of the algorithm

Initialisation via a PCA

```
1 initialisation.FA<-function(X,L=1){
2   # Return W and Psi
3   d<-ncol(X)
4   Sigmaxx<-var(X)
5   W<-eigen(Sigmaxx)$vectors[,1:L]
6   if (L==1) W<-cbind(W)
7   Psi<-rep(1,d)
8   return(list(W=W,Psi=Psi))
9 }
```

E step

```
1 FA.E.step<-function(X,W,Psi){
2   # X is assumed to be centered
3   # M contain the conditional expectation of the latent factor
4   # S contains the covariance of the latent factor
5   L<-ncol(W)
6   S <- solve(diag(L) + t(W)%*%diag(1/Psi)%*%W)
7   M<- X%*%diag(1/Psi)%*%W%*%S
```

```
8   return(list(S=S,M=M))
9 }
```

M Step

```
1 FA.M.step<-function(X,S,M,W,Psi){
2   n<-nrow(X)
3   Psi<-1/n*diag(t(X)%*%X -W%*%t(M)%*%X)
4   W<- (t(X)%*%M)%*%solve(n*S+t(M)%*%M)
5   return(list(Psi=Psi,W=W))
6 }
```

Computation of the criterion

```
1 log.likelihood.FA<-function(X,S,M,Psi,W){
2   n<-nrow(X)
3   Sigmax<-(t(X)%*%X-W%*%t(M)%*%X)
4   return(-(sum(diag(S+t(M)%*%M/n))
5           +log(det(diag(Psi))))+
6           log(det(S))+
7           1/n*sum(diag(Sigmax%*%diag(1/Psi))))))
8 }
```

Putting it all together

```
1 FA.EM<-function(X,L=1,max.iter=50){
2   X<-scale(X,scale=FALSE);mu<-attr(X,"scaled:center")
3   log.likelihood<-NULL; init<-initialisation.FA(X,L)
```

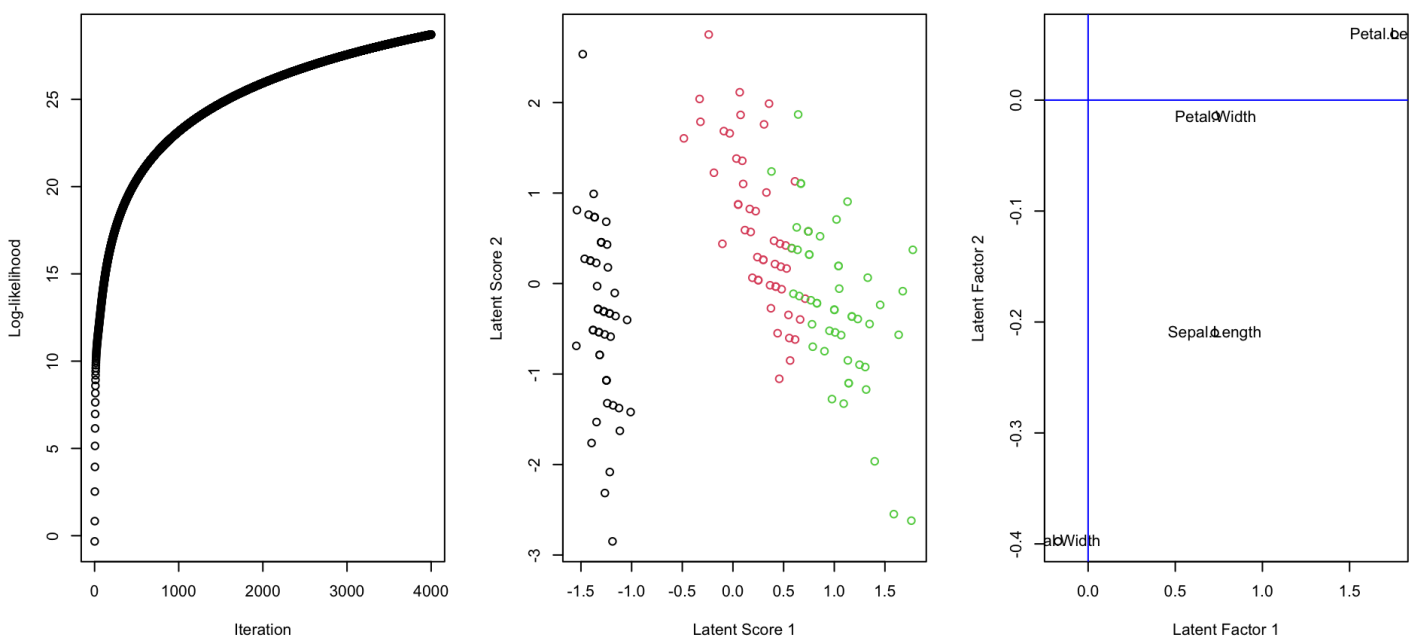
```

4 W<-init$W; Psi<-init$Psi; criterion<- Inf; iteration<-1;
5 log.likelihood[iteration]<--Inf
6 while ((criterion>1e-6)&&(iteration<=max.iter)){
7   E.step<-FA.E.step(X,W,Psi); E.step$S->S; E.step$M->M
8   M.step<-FA.M.step(X,S,M,W,Psi); M.step$Psi->Psi; M.step$W->W
9   iteration<-iteration+1
10  log.likelihood[iteration]<-log.likelihood.FA(X,S,M,Psi,W)
11  criterion<-abs((log.likelihood[iteration] - log.likelihood[iteration-1])/max(log
12 }
13 return(list( W=data.frame(W), Psi=Psi,
14             M=data.frame(M), S=S,mu=mu,
15             log.likelihood= log.likelihood[-1]))
16 }

```

1

Example with the Iris



1

Unidentifiability

- If we consider \mathbf{R} an orthogonal rotation matrix such that

$$\mathbf{R}\mathbf{R}^T = \mathbf{I}$$

It appears that $\tilde{\mathbf{W}} = \mathbf{W}\mathbf{R}$ produces the same log-likelihood.

- \mathbf{W} cannot be uniquely identified.

1

Possible rotations

1. Forcing \mathbf{W} to be orthogonal with columns ordered by decreasing variance
2. Forcing \mathbf{W} to be lower triangular (problem of founder variables)
3. Choosing an informative rotation matrix. For example the varimax rotation.
4. ...

Varimax

Varimax rotation maximizes the sum of the variance of the squared correlations between variables and factors

$$R_{\text{VARIMAX}} = \arg \max_R \left(\frac{1}{p} \sum_{j=1}^k \sum_{i=1}^p (WR)_{ij}^4 - \sum_{j=1}^k \left(\frac{1}{p} \sum_{i=1}^p (WR)_{ij}^2 \right) \right)$$

This results in high factor loadings for a small number of variables and low factor loadings for the rest.

1

Varimax

```
1 W.FA<-FA.result$W
2 W.Varimax<-varimax(as.matrix(FA.result$W))$loadings
3 print(W.Varimax)
```

Loadings:

	Latent Factor 1	Latent Factor 2
Sepal.Length	0.756	
Sepal.Width		-0.429
Petal.Length	1.683	0.509
Petal.Width	0.711	0.174

	Latent Factor 1	Latent Factor 2
SS loadings	3.916	0.473
Proportion Var	0.979	0.118
Cumulative Var	0.979	1.097

2

Mixture of factor analysers

- Factor analysis is a way to estimate a variance matrix with few parameters
- This property can be used in the context of Gaussian mixture model assuming the following parameterization for component densities:

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

where k is the component number and \mathbf{W}_k is a loading matrix defining the relation between the observation \mathbf{x}_i and the latent vector \mathbf{z}_i

- This approach is similar to the Banfield-Raftery idea of decomposing the component variance matrix $\boldsymbol{\Sigma}_k$ in volume, form and direction.

Relation to principal component analysis

Assumption

If

- $\Psi = \sigma^2 \mathbf{I}$
- \mathbf{W} is orthonal

and

- $\sigma^2 \rightarrow 0$

Then

Tipping, M. and C. Bishop (1999, Probabilistic principal component analysis. J. of Royal Stat. Soc. Series B 21(3), 611–622) showed that FA is equivalent to PCA

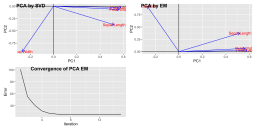
Criterion

$$J(\mathbf{W}, \mathbf{Z}) = \|\mathbf{X} - \mathbf{Z}\mathbf{W}^T\|_F^2$$

where $\mathbf{W}^T \mathbf{W} = \mathbf{I}$

A Constrained EM Algorithm for PCA (from Ahn, J.-H. and J.-H. Oh, 2003)

```
1 upper<-function(A){A[lower.tri(A,diag=FALSE)]<-0;return(A)}
2 lower<-function(A){A[upper.tri(A,diag=FALSE)]<-0;return(A)}
3 PCA.EM<-function(X,q=2){
4   p<-ncol(X); n<-nrow(X)
5   W<-diag(p)[,1:q]; M<-X%%W # Initialisation
6   Jold<-0; J<-1; iteration<-0; Error<-NULL
7   while ((abs(J - Jold)>1e-3)){
8     Jold<-sum((X-M%%t(W))^2)
9     S <- solve(upper(t(W)%%W)); M<- X%%W%%S # E-step
10    W<- (t(X)%%M)%%solve(lower(n*S+t(M)%%M))# M-step
11    W<-apply(W,2,function(x){x/sqrt(sum(x^2))})#orthogonalisation
12    J<-sum((X-M%%t(W))^2); Error[iteration<-iteration+1]<-J
13  }
14  return(list(W=data.frame(W),M=data.frame(M),Error=Error))}
```



Neural networks and unsupervised learning