

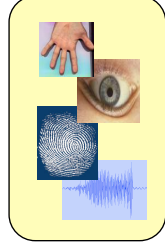
Factor Analysis

Leibny Paola García Perera.
 Carnegie Mellon University,
 Tecnológico de Monterrey, Campus Monterrey, Mexico
 Universidad de Zaragoza, Spain.
 Bhiksha Raj, Juan Arturo Nolasco Flores, Eduardo Lleida

Introduction

Problem: Lots of data with n -dimensions vectors.

Example:

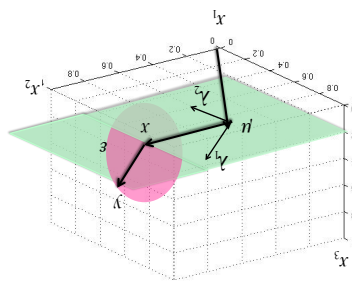


$$Y = \begin{bmatrix} d_{11} & d_{12} & d_{13} & \dots & d_{1p} \\ d_{21} & d_{22} & d_{23} & \dots & d_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ d_{N1} & d_{N2} & d_{N3} & \dots & d_{Np} \end{bmatrix}$$

Features

Can we reduce the number of dimensions? To reduce computing time, simplify process?
 © ESI

Factor Analysis (FA): Geometrical Representation

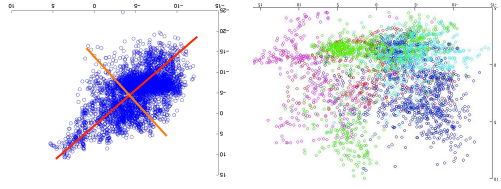


What is Factor Analysis?
 □ Analysis of the covariance in observed variables (Y).
 □ In terms of few (latent) common factors.
 □ Plus a specific error.

Factor Analysis (FA)

Introduction: Covariance matrix

- What can give us information of the data? (just for this special case)
- The covariance matrix
- Get rid of not important information.
- Think of continuous factors that control the data.



Agenda

- Introduction
- Motivation:
- Dimension reduction
- Modeling: covariance matrix
- Factor Analysis (FA)
- Geometrical explanation
- Formulation (The Equations)
- EM algorithm
- Comparison with PCA and PPCA.
- Example with numbers
- Applications
- Speaker Verification: Joint Factor Analysis (JFA)
- Some results
- References

Factor Analysis (FA): Formulation (the equations)

Now that we have checked the matrices dimensions.

The model:

$$d(x) = N(x|0, I)$$

$$d(y|x, \theta) = N(y|n + \Lambda x, \Psi)$$

Quick notes:

- Are Gaussians!!
- $d(x, y)$
- $d(x)$
- $d(y|x)$

Factor Analysis (FA): Formulation (the equations)

Form

$$y - \mu = \Lambda x + \epsilon$$

$$E(\Lambda^T) = I$$

$$\Psi_{11} = \begin{bmatrix} 0 & 0 \\ 0 & \psi_{pp} \end{bmatrix}$$

$$E(\epsilon \epsilon^T) = \Psi = \begin{bmatrix} 0 & 0 \\ 0 & \psi_{pp} \end{bmatrix}$$

Assumptions

- $y \leftarrow P \times 1$ data vector
- $\mu \leftarrow P \times 1$ mean vector
- $\Lambda \leftarrow P \times R$ loading Matrix
- $x \leftarrow R \times 1$ factor vector
- $\epsilon \leftarrow P \times 1$ error vector

$$E(y, x) = \Lambda$$

$$Z = E(y y^T) = \Lambda \Lambda^T + \Psi \text{ Full rank!!}$$

Factor Analysis (FA): Formulation (the equations)

So, factor analysis is a constrained covariance Gaussian Model!!!

$$d(y|x, \theta) = N(y|n + \Lambda x, \Psi)$$

So, what is the covariance?

$$\text{COV}(y) = \begin{bmatrix} \psi_{11} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \psi_{pp} \end{bmatrix} + \Lambda \Lambda^T$$

Factor Analysis (FA): Formulation (the equations)

Now, we can compute:

$$d(y|x, \theta) = \int d(x) d(y|x, \theta) dx = N(y|n + \Lambda x, \Psi)$$

This marginal is... a Gaussian!!

Compute the expected value and covariance:

$$E(y) = E(n + \Lambda x + \epsilon) = E(n) + \Lambda E(x) + E(\epsilon) = n$$

$$\text{Cov}(y) = E[(y - n)(y - n)^T] = E[(\Lambda x + \epsilon)(\Lambda x + \epsilon)^T] = \Lambda E[x x^T] \Lambda^T + E[\epsilon \epsilon^T] = \Lambda \Lambda^T + \Psi$$

Factor Analysis (FA): Formulation (the equations)

So we need sufficient statistics...

mean:

$$\sum_{t=1}^n y^t$$

covariance:

$$\sum_{t=1}^n (y^t - n)(y^t - n)^T$$

Factor Analysis (FA): Formulation (the equations)

How can we compute the likelihood function?

$$\ell(\theta, D) = -\frac{1}{N} \log |V \Lambda^T + \Psi| - \frac{1}{N} \sum_{t=1}^n (y^t - n - \Lambda x^t)^T (V \Lambda^T + \Psi)^{-1} (y^t - n - \Lambda x^t)$$

$$\ell(\theta, D) = -\frac{1}{N} \log |\Sigma| - \frac{1}{2N} \sum_{t=1}^n (y^t - n)^T (\Sigma^{-1} - \frac{1}{2} \Lambda^T \Lambda) (y^t - n)$$

Conclusion:

S is the sample data covariance Matrix.

Constrained model close to the Sample covariance!

Factor Analysis (FA): Expectation Maximization

- How to estimate μ ?
- Just compute the mean of the data.
- For the rest of the parameters Λ, Ψ ?
- Expectation Maximization

Factor Analysis (FA): Expectation Maximization

- Advantages
- Focuses on maximizing the likelihood
- Disadvantages
- Need to know the distribution
- No analytical solution

Factor Analysis (FA): Expectation Maximization

Remember EM algorithm?

- E-step: $q_{t+1}^n = d(x_n, \theta^t)$
- M-step: $\theta^{t+1} = \arg \max_{\theta} \sum_x \int d^u \log p(x_n, \theta)$

Factor Analysis (FA): Expectation Maximization

What do we need?

- E-step: Conditional probability!!!
- M-step: $q_{t+1}^n = d(x_n, \theta^t) = N(x_n | m_n, \Sigma_n)$

Log of the complete data for:

- $V_{t+1} = \arg \max_V \sum d^u \log p(x_n, y_n, V)$
- $\Psi_{t+1} = \arg \max_{\Psi} \sum d^u \log p(x_n, y_n, \Psi)$

Factor Analysis (FA): Expectation Maximization

What else is needed? $d(x|x)$

Let's start with: $d = \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} V \\ \Psi \end{pmatrix} + \begin{pmatrix} 0 \\ I \end{pmatrix} \begin{pmatrix} m \\ x \end{pmatrix}$

Remember that: $\text{cov}(x, y) = E(x - \mu)(x - \mu)^T = E(VV^T + \Psi\Psi^T + V\Psi^T + \Psi V^T)$

Factor Analysis (FA): Expectation Maximization

Now,

$d(x|x) = N(x|m, \Lambda)$

$V = V_{t+1} - \Psi_{t+1} \Psi_{t+1}^{-1} (V_{t+1} - \Psi_{t+1}^{-1} V_{t+1} \Psi_{t+1}^{-1} + \Psi_{t+1}^{-1} \Psi_{t+1}^{-1} V_{t+1} \Psi_{t+1}^{-1})$

Remember inversion lemma:

$\Lambda^{-1} = \Lambda^{-1} + \Psi_{t+1}^{-1} \Psi_{t+1}^{-1} \Psi_{t+1}^{-1}$

Remembering Gaussian conditioning formulas

Inverting this matrix is much more efficient $O(M^2)$ instead of $O(P^2)$, thanks to the lemma.

$$E[\mu_j] = \frac{1}{N} \sum_{i=1}^N \mu_j^i$$

$$E[\sigma_j^2] = \frac{1}{N} \sum_{i=1}^N \sigma_j^{i2}$$

And the expectations with respect to μ_j

$$\frac{\partial \log \ell(\mu, \Sigma)}{\partial \mu_j} = -\frac{1}{\sigma_j^2} \sum_{i=1}^N (\mu_j^i - \mu_j)$$

$$\frac{\partial \log \ell(\mu, \Sigma)}{\partial \sigma_j^2} = -\frac{1}{2\sigma_j^4} \sum_{i=1}^N (\mu_j^i - \mu_j)^2$$

Now, let's compute the M step! (Almost there!)
 We need to calculate the derivatives of the log likelihood

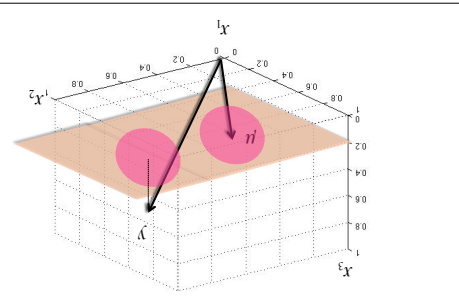
Factor Analysis (FA): Expectation Maximization

$$\mu_j = \frac{1}{N} \sum_{i=1}^N \mu_j^i$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N \sigma_j^{i2}$$

Finally, set the derivatives to zero and solve!

Factor Analysis (FA): Expectation Maximization



How does it look?

Factor Analysis (FA): Expectation Maximization

$$S = \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T$$

$$\frac{\partial \log \ell(\mu, \Sigma)}{\partial \mu} = -\frac{1}{\Sigma} \sum_{i=1}^N (x_i - \mu)$$

$$\frac{\partial \log \ell(\mu, \Sigma)}{\partial \Sigma} = -\frac{1}{2\Sigma^2} \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T$$

Let's subtract the mean for our computation.

Factor Analysis (FA): Expectation Maximization

$$p(x|y) = N(x|m, \Lambda)$$

$$\Lambda = (I - \Phi \Phi^T)^{-1}$$

$$m = \Phi \mu + \Lambda^{-1} y$$

We finally obtain:

Factor Analysis (FA): Expectation Maximization

$$p(x|y) = N(x|m, \Lambda)$$

$$\Lambda = (I - \Phi \Phi^T)^{-1}$$

$$m = \Phi \mu + \Lambda^{-1} y$$

Means that the posterior mean is just a linear operation!!!
 And the covariance does not depend on the observed data!!!

Some nice observations:

Factor Analysis (FA): Expectation Maximization

Factor Analysis (FA): Expectation Maximization

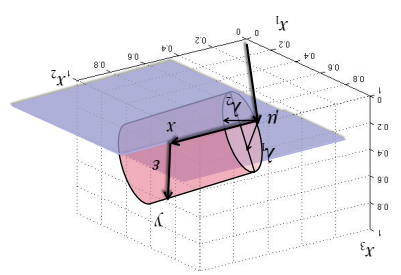
What are the final equations?

$$\begin{aligned}
 \textcircled{1} \mu &\rightarrow \text{Sample mean (Subtract the mean from data).} \\
 \textcircled{2} q_{t+1}^n &= p(x^n | y^n, \theta^t) = N(x^n | m_n, V_n) \\
 V_n &= (I - V V^T \Psi^{-1} V)^{-1} \\
 m_n &= V V^T \Psi^{-1} (y - \mu) \\
 V_{t+1} &= \left(\sum_{n=1}^m x^n x^{nT} \right)^{-1} \left(\sum_{n=1}^m x^n y^n \right) \\
 \Psi_{t+1} &= \frac{N}{I} \text{diag} \left(\sum_{n=1}^m y^n y^{nT} + V V^T + \sum_{n=1}^m m_n m_n^T \right)
 \end{aligned}$$

ⓐ M-step

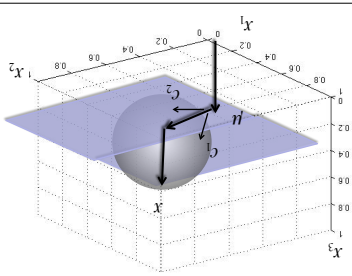
Factor Analysis (FA): Geometrical Representation

How does FA really look like?



Factor Analysis (FA): Comparison

Nice, isn't it? ☺

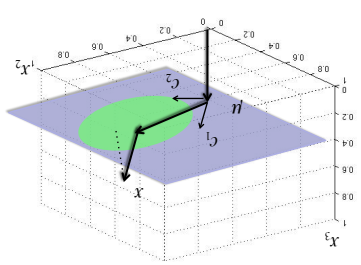


$$d(x) = N(x | 0, I) \quad d(y | x, \theta) = N(y | n + \Lambda x, \sigma^2 I)$$

What is PPCA? Just a quick intuition.

Factor Analysis (FA): Comparison

Nice! ☺



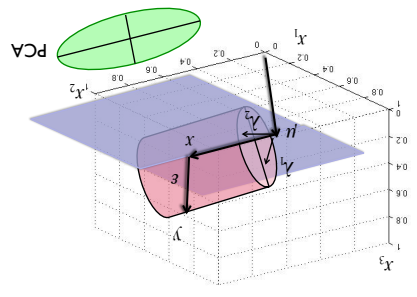
$$d(x) = N(x | 0, I) \quad d(y | x, \theta) = N(y | n + \Lambda x, 0) \quad \sigma^2 \rightarrow 0$$

What about PCA? Just a quick intuition.

Factor Analysis (FA): Comparison

- Final notes:
- FA is invariant if we change the scale.
- FA looks for correlation of the data.
- PCA is invariant if we rotate the data.
- PCA looks for direction of large variance.

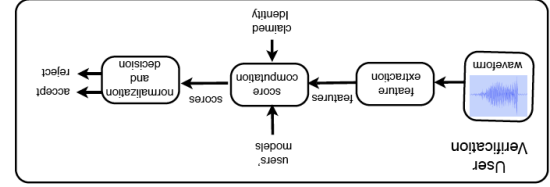
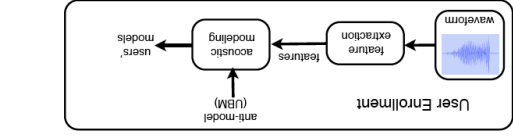
Factor Analysis (FA): Comparison



Speaker Verification

- Each speaker has its own model, known as target model λ_i
- And its antimodel λ_j
- The target model is the prototype of each speaker in the training.
- The antimodel is the impostor's prototype.
- When all the impostors share the same model, the final model is called: UBM Universal Background Model.

Speaker Verification



Motivation of using JFA

- ② Speaker model generation:
 - The amount of speech is quite small for an optimal estimation.
 - It is not possible to use rely on EM
- Solution: MAP (maximum a posteriori)**
- $$\theta^{MAP} = \arg \max_{\lambda, \theta} p(X|\lambda; \theta) d(\lambda; \theta),$$

Speaker Verification

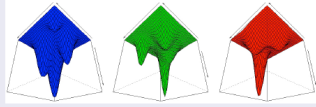
This is how we represent each speaker:

Gaussian Mixture Model (GMM)

- Let X be the acoustic feature vector.

$$P(X; \lambda) = \sum_{k=1}^K w_k \mathcal{N}(X; \mu_k, \Sigma_k), \quad \lambda = (w_k, \mu_k, \Sigma_k)$$

Three sample GMMs for a 3D feature vector:



- The GMM characterizes the set of mechanical configurations of a person's vocal tract.

Motivation of using JFA

- ① UBM Generation
 - We take all the data available and model a GMM (independent to the target speakers).
 - The technique used is: Expectation Maximization (EM).
- Traditional systems are based on the estimation of the probability density functions (GMM in this case).

Motivation of using JFA

- What is the real problem?
 - Speaker data trained over different channels.
 - MAP doesn't work. It does assume conventional conjugate priors.
 - What is the solution for non-ideal cases?
 - Provides priors for the parameters.
 - Separates the speaker and the channel factors.
 - The channel factors don't give information of the speaker so they can be marginalized out when computing score.
- JFA!!!

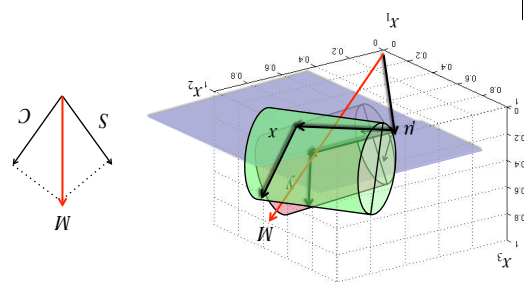
Motivation

- ② Speaker model generation JFA Joint Factor Analysis
- Is it possible to include a new latent variable? YES!!!
- What is the new model?

$$M = m + Vy + Ux + Dz,$$

- $m \rightarrow CF \times 1$ supervector
- $V \rightarrow$ low rank matrix eigenvalues
- $y \rightarrow$ speaker factors
- $z \rightarrow$ normally distributed random vector
- $U \rightarrow$ low rank matrix eigenchannels
- $x \rightarrow$ channel factors
- $D \rightarrow$ diagonal matrix

Factor Analysis (FA): Geometrical representation



Algorithm

We may use a variable change in order to get an estimation of the Vy, Ux and Dz contributions with the Factor Analysis estimating methods.

$$Data1' = m + Dz + Ux$$

$$Data = Data1' + Vy$$

$$Data2' = m + Vy + Dz$$

$$Data = Data2' + Ux$$

$$Data3' = m + Vy + Ux$$

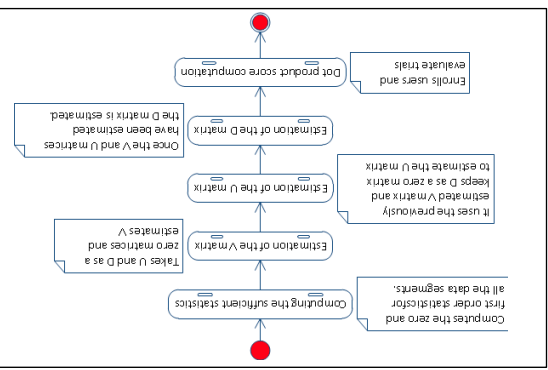
$$Data = Data3' + Dz$$

Algorithm

- ① Compute Sufficient Statistics
- ② Compute V and Y
- ③ Compute U and X
- ④ Compute D and Z



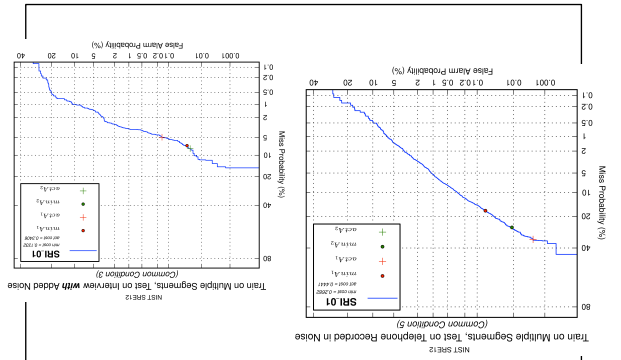
JFA Algorithm



History

What happened next?
 Researchers discovered that the channel factors contained information of the speaker. ☺
 Go back to factor analysis!!! Now is called: I-vectors!!!
 Important notes:
 □ JFA is actually used to build a model of the data
 □ I-vectors are used as feature extractor:
 Obtains the important information of the speakers and transforms it into vectors.

Some results.. Last week. Best system!



- **Saul and Rahim.** Maximum Likelihood and Minimum Classification Error Factor Analysis for Automatic Speech Recognition
- **D'Souza.** Derivation of Maximum Likelihood Factor Analysis using EM
- **Johnson and Wichein.** Applied Multivariate Statistical Analysis
- **Dehak, N., Kenny, P., Dehak, R., Dumouchel, P and Ouellet,** P. Front-End Factor Analysis for Speaker Verification
- **Kenny, P** Joint factor analysis of speaker and session variability : Theory and algorithms - Technical report CRIM-06/08-13 Montreal, CRIM, 2005

References: