

CS 5135/6035 Learning Probabilistic Models

Lecture 12: Factor Analysis

Gowtham Atluri

October 12, 2018

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

1 / 36

Factor Analysis

- Problem definition
- Benefits
- Issues

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

3 / 36

Factor Analysis

- Hidden variable(s) is (are) continuous
- First studied by Charles Spearman in 1904
 - to explain hidden structure of human intelligence
- Spearman observed that schoolchildren's grades in different subjects were correlated with each other
- He explained:
 - the reason grades in math, English, history, etc. are correlated with each other, is they are all correlated with something else
 - a **common factor**, which he named "general intelligence"

ID	math	English	history	Intelligence
1	77.4	82.6	81.6	?
2	58.5	61.3	52.2	?
.	.	.	.	?
100	76.7	72.4	72.2	?

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

5 / 36

Reading Material

- Andrew Ng, Factor Analysis
<http://cs229.stanford.edu/notes/cs229-notes9.pdf>
- Cosma Shalizi, Factor Analysis
<https://www.stat.cmu.edu/~cshalizi/350/lectures/12/lecture-12.pdf>

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

2 / 36

Latent Variables

Mixture Models

$$z_i \sim \text{Multinomial}(\lambda_1, \lambda_2, \dots, \lambda_k)$$

$$x_i | z_i \sim \mathcal{N}(\mu_{z_i}, \sigma_{z_i}^2)$$

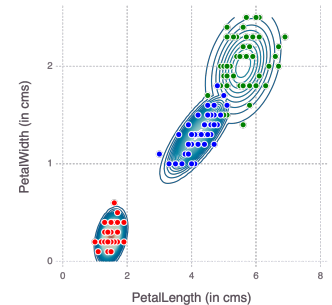
Note: Latent var. z_i is discrete; Observed var. x_i is continuous

Parameter estimation using MLE

EM Algorithm

E-step: Estimate $p(z_i | x_i)$, assuming μ_i and σ_i^2 ($i \in \{1, \dots, k\}$) are available

M-step: Estimate μ_i and σ_i^2 ($i \in \{1, \dots, k\}$), assuming $p(z_i | x_i)$ is available



Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

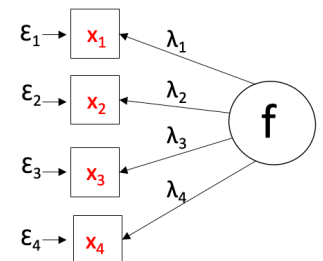
October 12, 2018

4 / 36

Spearman's factor analysis model

$$x_i = \lambda_i f + \epsilon_i$$

- x_i are the **observed variables**
 - e.g., x_1, x_2 , and x_3 are exam scores obtained by a student in math, English and history.
- f is the underlying **common factor**
 - e.g., student's intelligence
- λ_i are the **factor loadings**
 - e.g., how much is the contribution of intelligence to exam score
- ϵ_i are **unique factors** or residuals or random noise terms
 - e.g., how much result differs from student's general ability



Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

6 / 36

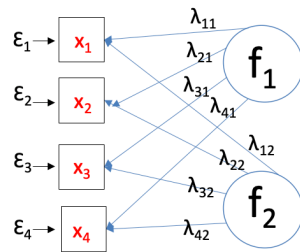
When there are multiple factors...

$$x_i = \lambda_{i1}f_1 + \lambda_{i2}f_2 + \dots + \lambda_{ik}f_k + \epsilon_i$$

Example: Source-apportionment of fine particulate matter (PM $\leq 2.5\mu m$)

Factors include:

- Vehicle emissions
- Coal combustion
- Biomass combustion
- residential heaters
- industrial processes



For multiple samples

$$x_{i,1} = \lambda_{11}f_{i,1} + \lambda_{12}f_{i,2} + \epsilon_{i,1}$$

$$x_{i,2} = \lambda_{21}f_{i,1} + \lambda_{22}f_{i,2} + \epsilon_{i,2}$$

$$x_{i,3} = \lambda_{31}f_{i,1} + \lambda_{32}f_{i,2} + \epsilon_{i,3}$$

$$x_{i,4} = \lambda_{41}f_{i,1} + \lambda_{42}f_{i,2} + \epsilon_{i,4}$$

- $x_{i,j}$ is the j^{th} observed variable in the i^{th} sample
 - Total number of variables are denoted by d
- $f_{i,j}$ is the j^{th} factor for sample i
 - Total number of factors are denoted by k
- $\lambda_{i,j}$ is the loading of the j^{th} factor on the i^{th} observed variable
- $\epsilon_{i,j}$ is the unique factor contributing to the j^{th} observed variable in the i^{th} sample

For multiple samples

$$x_{i,1} = \lambda_{11}f_{i,1} + \lambda_{12}f_{i,2} + \epsilon_{i,1}$$

$$x_{i,2} = \lambda_{21}f_{i,1} + \lambda_{22}f_{i,2} + \epsilon_{i,2}$$

$$x_{i,3} = \lambda_{31}f_{i,1} + \lambda_{32}f_{i,2} + \epsilon_{i,3}$$

$$x_{i,4} = \lambda_{41}f_{i,1} + \lambda_{42}f_{i,2} + \epsilon_{i,4}$$

Matrix notation

$$\begin{bmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \\ x_{i,4} \end{bmatrix} = \begin{bmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{21} & \lambda_{22} \\ \lambda_{31} & \lambda_{32} \\ \lambda_{41} & \lambda_{42} \end{bmatrix} \begin{bmatrix} f_{i,1} \\ f_{i,2} \end{bmatrix} + \begin{bmatrix} \epsilon_{i,1} \\ \epsilon_{i,2} \\ \epsilon_{i,3} \\ \epsilon_{i,4} \end{bmatrix}$$

$$\mathbf{x} = \mathbf{\Lambda}\mathbf{f} + \boldsymbol{\epsilon}$$

Advantages: Applications

- Matrix Factorization
 - [Genes \times Conditions] = [Genes \times Pathways] * [Pathways \times Conditions]
 - [Users \times Movies] = [Users \times Genres] * [Genres \times Movies]
- Dimensionality reduction
 - Work with \mathbf{f} (k-dim) instead of \mathbf{x} (d-dim) ($k < d$)

For all the samples $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} \end{bmatrix} = \begin{bmatrix} f_{11} & f_{12} \\ f_{21} & f_{22} \\ f_{31} & f_{32} \\ \vdots & \vdots \\ f_{n1} & f_{n2} \end{bmatrix} \begin{bmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{21} & \lambda_{22} \\ \lambda_{31} & \lambda_{32} \\ \lambda_{41} & \lambda_{42} \end{bmatrix}^T + \begin{bmatrix} \epsilon_{11} & \epsilon_{12} & \epsilon_{13} & \epsilon_{14} \\ \epsilon_{21} & \epsilon_{22} & \epsilon_{23} & \epsilon_{24} \\ \epsilon_{31} & \epsilon_{32} & \epsilon_{33} & \epsilon_{34} \\ \vdots & \vdots & \vdots & \vdots \\ \epsilon_{n1} & \epsilon_{n2} & \epsilon_{n3} & \epsilon_{n4} \end{bmatrix}$$

Problem of indeterminacy, Imposing Constraints

$$\mathbf{x} = \mathbf{\Lambda}\mathbf{f} + \boldsymbol{\epsilon}$$

- Problem of indeterminacy
 - No values on the R.H.S. are known
 - Many solutions are possible
- We impose constraints to reduce indeterminacy
 - Unique factors all have mean zero ($E(\boldsymbol{\epsilon}) = 0$)
 - Common factors all have mean zero ($E(\mathbf{f}) = 0$)
 - Let $E(\mathbf{x}) = \boldsymbol{\mu}$

$$\mathbf{x} = \boldsymbol{\mu} + \mathbf{\Lambda}\mathbf{f} + \boldsymbol{\epsilon}$$

More Constraints

$$\mathbf{x} = \boldsymbol{\mu} + \mathbf{\Lambda}\mathbf{f} + \boldsymbol{\epsilon}$$

- More constraints
 - Common factors are standardized uncorrelated random variables $E(\mathbf{f}\mathbf{f}^T) = \mathbf{I}$
 - their variance 1, covariance 0.
 - Unique factors are uncorrelated and heteroscedastic $E(\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T) = \Psi$

$$\Psi = \begin{bmatrix} \psi_{11} & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & \psi_{rr} & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \psi_{dd} \end{bmatrix}$$

- Unique factors are uncorrelated with the common factors

$$E(\mathbf{f}\boldsymbol{\epsilon}^T) = 0 \text{ and } E(\boldsymbol{\epsilon}\mathbf{f}^T) = 0$$

Problem of Identifiability

Identifiability

A statistical model is said to be identifiable if a unique solution exists to the estimation problem.

- The factor model $\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{f} + \boldsymbol{\epsilon}$ is not identifiable
 - there is an infinite number of different matrices $\boldsymbol{\Lambda}$ that can generate the same \mathbf{x} values
- The model can be equivalently rewritten as

$$\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{f} + \boldsymbol{\epsilon} = \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{f} + \boldsymbol{\epsilon} = \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{G}\mathbf{G}^T + \boldsymbol{\epsilon}$$
- If we set, $\boldsymbol{\Lambda}^* = \boldsymbol{\Lambda}\mathbf{G}$ and $\mathbf{f}^* = \mathbf{G}^T\mathbf{f}$, the factor model becomes

$$\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Lambda}^*\mathbf{f}^* + \boldsymbol{\epsilon}$$
- The two factor models are completely equivalent and indistinguishable
 - They have the same properties
 - $E[\mathbf{f}^*] = 0$
 - $E[\mathbf{f}^*\mathbf{f}^{*T}] = \mathbf{I}$
 - $E[\mathbf{f}^*\boldsymbol{\epsilon}^T] = 0$

Probabilistic Modelling

- Formulation
- Geometric interpretation
- Joint distribution
- Max. Likelihood Estimation
 - EM approach

Probabilistic Modelling

Generative model for factor analysis

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

$$\mathbf{x}|\mathbf{f} \sim \mathcal{N}(\boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{f}, \boldsymbol{\Psi})$$

Alternatively,

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

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi})$$

$$\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{f} + \boldsymbol{\epsilon}$$

Parameters of this model are:

- Vector $\boldsymbol{\mu} \in \mathbb{R}^d$
- Matrix $\boldsymbol{\Lambda} \in \mathbb{R}^{d \times k}$
 - usually $k < d$
- Diagonal matrix $\boldsymbol{\Psi} \in \mathbb{R}^{d \times d}$

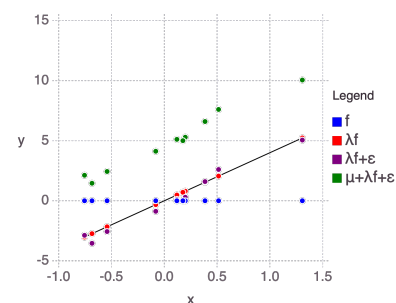
This assumption that data follows a Gaussian distr. need not suit all cases.
- Can be a disadvantage.

Geometric Interpretation in 2D

Data generating process:

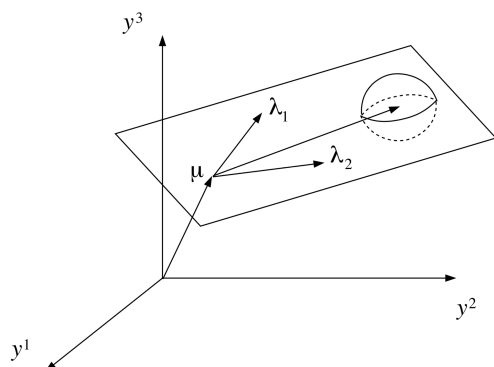
$$\mathbf{f} \sim \mathcal{N}(0, 1)$$

$$\mathbf{x}|\mathbf{f} \sim \mathcal{N}(\boldsymbol{\mu} + \lambda\mathbf{f} + \boldsymbol{\epsilon}, \sigma)$$



- Data points lie in a close to linear subspace.
 - Observed variables lie in 2D
 - Factors in 1D

Geometric Interpretation in 3D



- Data points lie in a close to linear subspace.

Joint Distribution

Our Model

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

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi})$$

$$\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{f} + \boldsymbol{\epsilon}$$

RVs \mathbf{x} and \mathbf{f} have a joint Gaussian distribution

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{x} \end{bmatrix} \sim \mathcal{N}(\boldsymbol{\mu}_{f,x}, \boldsymbol{\Sigma})$$

We will now determine $\boldsymbol{\mu}_{f,x}, \boldsymbol{\Sigma}$

$$E[\mathbf{x}] = E[\boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{f} + \boldsymbol{\epsilon}]$$

$$= \boldsymbol{\mu} + \boldsymbol{\Lambda}E[\mathbf{f}] + E[\boldsymbol{\epsilon}]$$

$$= \boldsymbol{\mu} \text{ (as } E[\mathbf{f}] = 0 \text{ and } E[\boldsymbol{\epsilon}] = 0 \text{ from constraints)}$$

From this, we have

$$\boldsymbol{\mu}_{f,x} = \begin{bmatrix} \mathbf{0} \\ \boldsymbol{\mu} \end{bmatrix}$$

Joint Distribution

$$\Sigma = \begin{bmatrix} \Sigma_{f,f} & \Sigma_{f,x} \\ \Sigma_{x,f} & \Sigma_{x,x} \end{bmatrix} = \begin{bmatrix} E[(f - E[f])(f - E[f])^T] & E[(f - E[f])(x - E[x])^T] \\ E[(x - E[x])(f - E[f])^T] & E[(x - E[x])(x - E[x])^T] \end{bmatrix}$$

Since $f \sim \mathcal{N}(\mu, \Lambda)$, $\Sigma_{f,f} = \text{Cov}(f) = \Lambda$.

$$\begin{aligned} E[(f - E[f])(x - E[x])^T] &= E[(\mu + \Lambda f + \epsilon - \mu)^T] \\ &= E[f^T \Lambda + \epsilon^T] \\ &= \Lambda^T \end{aligned}$$

$$\begin{aligned} E[(x - E[x])(x - E[x])^T] &= E[(\mu + \Lambda f + \epsilon - \mu)^T(\mu + \Lambda f + \epsilon - \mu)^T] \\ &= E[\Lambda f f^T \Lambda^T + \epsilon f^T \Lambda^T + \Lambda f \epsilon^T + \epsilon \epsilon^T] \\ &= \Lambda E[f f^T] \Lambda^T + E[\epsilon \epsilon^T] \\ &= \Lambda \Lambda^T + \Psi \end{aligned}$$

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

19 / 36

Joint Distribution

- Putting it all together...

$$\begin{bmatrix} f \\ x \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda \Lambda^T + \Psi \end{bmatrix}\right)$$

- From this, marginal distribution of x is $x \sim \mathcal{N}(\mu, \Lambda \Lambda^T + \Psi)$
- Given a dataset of samples $D = \{x_1, x_2, \dots, x_n\}$, we can write the log likelihood of the parameters as:

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

- To perform MLE, we need to maximize this quantity w.r.t. the parameters.
 - No closed form solution exists
 - We will use Expectation-Maximization (EM) algorithm

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

20 / 36

EM algorithm

- For any choice of $q(\mathbf{Z})$, the following decomposition holds

$$\log p(\mathbf{X}|\theta) = \mathcal{L}(q, \theta) + KL(q||p)$$

where we define

$$\mathcal{L}(q, \theta) = \sum_{\mathbf{z}} q(\mathbf{z}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\theta)}{q(\mathbf{Z})} \right\}$$

$$KL(q||p) = -\sum_{\mathbf{z}} q(\mathbf{z}) \log \left\{ \frac{p(\mathbf{Z}|\theta)}{q(\mathbf{Z})} \right\}$$

- As $KL(p||q) \geq 0$, $\mathcal{L}(q, \theta) \leq \log p(\mathbf{X}|\theta)$.
 - $\mathcal{L}(q, \theta)$ is the lower bound on $\log p(\mathbf{X}|\theta)$.
- In E-Step: Lower bound $\mathcal{L}(q, \theta)$ is maximized w.r.t. $q(\mathbf{Z})$, fixing θ^{old}
 - Essentially substituting $q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}, \theta)$
- In M-Step: $\mathcal{L}(q, \theta)$ is maximized w.r.t. θ to give some new value θ^{new}

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

21 / 36

EM Algorithm: E-Step

In the E-step, we need to compute $Q_i(f_i) = p(f_i|x_i; \mu, \Lambda, \Psi)$

$$\begin{bmatrix} f \\ x \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda \Lambda^T + \Psi \end{bmatrix}\right)$$

We know for MV Gaussians...

$$\begin{bmatrix} a \\ b \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}, \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}\right)$$

the conditional $p(a|b) = \mathcal{N}(a|\mu_a + \Sigma_{ab}\Sigma_{bb}^{-1}(b - \mu_b), \Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})$

Using this property, we can write $p(f_i|x_i; \mu, \Lambda, \Psi) = \mathcal{N}(\mu_{f_i|x_i}, \Sigma_{f_i|x_i})$

where, $\mu_{f_i|x_i} = \Lambda^T(\Lambda \Lambda^T + \Psi)^{-1}(x_i - \mu)$

$$\Sigma_{f_i|x_i} = I - \Lambda^T(\Lambda \Lambda^T + \Psi)^{-1}\Lambda$$

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

22 / 36

EM Algorithm: M-Step

We need to maximize the following function w.r.t. parameters μ, Λ, Ψ

$$\begin{aligned} &\sum_{i=1}^n \int_{f_i} q(f_i) \log \frac{p(x_i, f_i; \mu, \Lambda, \Psi)}{q(f_i)} df_i \\ &= \sum_{i=1}^n \int_{f_i} q(f_i) [\log p(x_i|f_i; \mu, \Lambda, \Psi) + \log p(f_i) - \log q(f_i)] df_i \\ &= \sum_{i=1}^n E_{f_i \sim q} [\log p(x_i|f_i; \mu, \Lambda, \Psi) + \log p(f_i) - \log q(f_i)] \end{aligned}$$

First, computing gradient w.r.t. Λ (we only need the first log term)

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

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

23 / 36

EM Algorithm: M-Step

- Some results from multivariate calculus

$$\begin{aligned} &\bullet \text{ } tra = a \text{ for } a \in \mathbb{R}; \text{ } trAB = trBA; \nabla_a trABA^T C = CAB + C^T AB \\ &\sum_{i=1}^n E[\log p(x_i|f_i; \mu, \Lambda, \Psi)] \end{aligned}$$

$$= \sum_{i=1}^n E\left[-\frac{1}{2} \log |\Psi| - \frac{n}{2} \log(2\pi) - \frac{1}{2}(x_i - \mu - \Lambda f_i)^T \Psi^{-1} (x_i - \mu - \Lambda f_i)\right]$$

$$\nabla_{\Lambda} \sum_{i=1}^n E\left[\frac{1}{2}(x_i - \mu - \Lambda f_i)^T \Psi^{-1} (x_i - \mu - \Lambda f_i)\right]$$

$$= \sum_{i=1}^n \nabla_{\Lambda} E\left[tr \frac{1}{2} f_i^T \Lambda^T \Psi^{-1} \Lambda f_i + tr f_i^T \Lambda^T \Psi^{-1} (x_i - \mu)\right]$$

$$= \sum_{i=1}^n \nabla_{\Lambda} E\left[tr \frac{1}{2} \Lambda^T \Psi^{-1} \Lambda f_i f_i^T + tr \Lambda^T \Psi^{-1} (x_i - \mu) f_i^T\right]$$

$$= \sum_{i=1}^n E\left[-\Psi^{-1} \Lambda f_i f_i^T + \Psi^{-1} (x_i - \mu) f_i^T\right]$$

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

24 / 36

EM Algorithm: M-Step

Gradient w.r.t. Λ is

$$\sum_{i=1}^n E \left[-\Psi^{-1} \Lambda f_i f_i^T + \Psi^{-1} (x_i - \mu) f_i^T \right]$$

Setting this to 0 and simplifying,

$$\sum_{i=1}^n \Lambda E_{f_i \sim q} [f_i f_i^T] = \sum_{i=1}^n (x_i - \mu) E_{f_i \sim q} [f_i^T]$$

Solving for Λ , we get

$$\Lambda = \left(\sum_{i=1}^n (x_i - \mu) E_{f_i \sim q} [f_i^T] \right) \left(\sum_{i=1}^n E_{f_i \sim q} [f_i f_i^T] \right)^{-1}$$

EM Algorithm: M-Step

$$\Lambda = \left(\sum_{i=1}^n (x_i - \mu) E_{f_i \sim q} [f_i^T] \right) \left(\sum_{i=1}^n E_{f_i \sim q} [f_i f_i^T] \right)^{-1}$$

We know from our definition of q being Gaussian and $\text{Cov}(Y) = E[YY^T] - E[Y]E[Y^T]$,

$$E_{f_i \sim q} [f_i^T] = \mu_{f_i|x_i}^T$$

$$E_{f_i \sim q} [f_i f_i^T] = \mu_{f_i|x_i} \mu_{f_i|x_i}^T + \Sigma_{f_i|x_i}$$

$$\Lambda = \left(\sum_{i=1}^n (x_i - \mu) \mu_{f_i|x_i}^T \right) \left(\sum_{i=1}^n \mu_{f_i|x_i} \mu_{f_i|x_i}^T + \Sigma_{f_i|x_i} \right)^{-1}$$

EM Algorithm: M-Step

Similarly, we can find M-Step optimizations for μ and Ψ .

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\Phi = \frac{1}{n} \sum_{i=1}^n x_i x_i^T - x_i \mu_{f_i|x_i}^T \Lambda^T - \Lambda \mu_{f_i|x_i} x_i^T + \Lambda (\mu_{f_i|x_i} \mu_{f_i|x_i}^T + \Sigma_{f_i|x_i}) \Lambda^T$$

and setting $\Psi_{ij} = \Phi_{ij}$, i.e., letting Ψ be the diagonal matrix containing only diagonal elements of Φ .

EM algorithm

E Step

$$\mu_{f_i|x_i} = \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} (x_i - \mu)$$

$$\Sigma_{f_i|x_i} = I - \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} \Lambda$$

M Step

$$\Lambda = \left(\sum_{i=1}^n (x_i - \mu) \mu_{f_i|x_i}^T \right) \left(\sum_{i=1}^n \mu_{f_i|x_i} \mu_{f_i|x_i}^T + \Sigma_{f_i|x_i} \right)^{-1}$$

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\Phi = \frac{1}{n} \sum_{i=1}^n x_i x_i^T - x_i \mu_{f_i|x_i}^T \Lambda^T - \Lambda \mu_{f_i|x_i} x_i^T + \Lambda (\mu_{f_i|x_i} \mu_{f_i|x_i}^T + \Sigma_{f_i|x_i}) \Lambda^T$$

and setting $\Psi_{ij} = \Phi_{ij}$, i.e., letting Ψ be the diagonal matrix containing only diagonal elements of Φ .

Julia Implementation

- E-Step and M-Step
- Computing Log-likelihood
- EM algorithm
- Application on a synthetic dataset

E-Step

$$\mu_{f_i|x_i} = \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} (x_i - \mu)$$

$$\Sigma_{f_i|x_i} = I - \Lambda^T (\Lambda \Lambda^T + \Psi)^{-1} \Lambda$$

```
function E_Step(X,mu,Lambda,Psi,k)
    mu_f_by_x = (X - repmat(mu',size(X,1),1))*
                (Lambda'*inv(Lambda*Lambda' + Psi))';
    Sig_f_by_x = eye(k) - Lambda'*
                inv(Lambda*Lambda' + Psi)*Lambda;
    return mu_f_by_x,Sig_f_by_x;
end
```

E_Step (generic function with 1 method)

M-Step

```
function M_Step(X,mu_f_by_x,Sig_f_by_x,k)
    nrows, ncols = size(X);
    #Computing mu
    mu = mean(X,1)';
    #Computing Lambda
    Lambda_term1 = zeros(ncols,k); Lambda_term2 = zeros(k,k);
    for i=1:nrows
        Lambda_term1 += ((X[i,:] - mu)*mu_f_by_x[i,:])';
        Lambda_term2 += (mu_f_by_x[i,:]*mu_f_by_x[i,:]')
            + Sig_f_by_x;
    end
    Lambda = Lambda_term1*inv(Lambda_term2);
    #Computing Psi
    Phi = zeros(ncols,ncols);
    for i=1:nrows
        Phi += (X[i,:]*X[i,:]' - X[i,:]*mu_f_by_x[i,:]'*Lambda'
            - Lambda*mu_f_by_x[i,:]*X[i,:]'
            + Lambda*(mu_f_by_x[i,:]*mu_f_by_x[i,:]' + Sig_f_by_x)*Lambda')
    end
    Psi = diagm(diag(Phi./nrows));
    return mu, Lambda, Psi
end
```

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

31 / 36

Computing Log-Likelihood

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

```
function compute_llh(X,mu,Lambda,Psi)
    llh = 0;
    for i=1:size(X,1)
        llh = llh + log(pdf(MvNormal(vec(mu),
            (Lambda*Lambda')+Psi),X[i,:])));
    end
    return llh;
end
```

compute_llh (generic function with 1 method)

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

32 / 36

Overall EM approach for Factor Analysis

```
function fa_em(X,k)
    max_iter = 100;
    eps = 0.0001;
    llh = -Inf*ones(max_iter+1);
    mu = mean(X,1)';
    Lambda = rand(size(X,2),k);
    Psi = diagm(rand(size(X,2)));
    llh[1] = compute_llh(X,mu,Lambda,Psi);
    for i=1:max_iter
        mu_f_by_x,Sig_f_by_x = E_Step(X,mu,Lambda,Psi,k);
        mu_new, Lambda_new, Psi_new = M_Step(X,mu_f_by_x,Sig_f_by_x,k);
        llh[i+1] = compute_llh(X,mu_new,Lambda_new,Psi_new);
        if (sum(abs.(mu_new-mu))<eps &&
            sum(abs.(Lambda_new-Lambda))<eps && sum(abs.(Psi_new-Psi))<eps)
            break;
        end
        mu = mu_new; Lambda = Lambda_new; Psi = Psi_new;
    end
    return mu, Lambda, Psi, llh;
end
```

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

33 / 36

Generating sample data

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

$$\mathbf{x} | \mathbf{f} \sim \mathcal{N}(\boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{f}, \boldsymbol{\Psi})$$

```
mu = [0 0 0 0]';
Lambda =
[1.0 0
 1.0 0
 0 0.99
 0 1.0];
Psi = diagm([0.1, 0.1, 0.1, 0.1]);
d1 = MvNormal([0,0],ones(2));
X = zeros(10000,4);
for i=1:10000
    f = rand(d1,1);
    d2 = MvNormal(vec(mu+ Lambda*f),Psi);
    x = rand(d2,1);
    X[i,:] = x';
end
```

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

34 / 36

Using EM to learn Lambda

```
mu_est, Lambda_est, Psi_est, llh = fa_em(X,2);
```

mu_est

```
## 4x1 Array{Float64,2}:
## -0.0092927
## -0.0126923
## -0.00481124
## 0.00146455
```

Psi_est

```
## 4x4 Array{Float64,2}:
## 0.094361 0.0 0.0 0.0
## 0.0 0.10286 0.0 0.0
## 0.0 0.0 0.100012 0.0
## 0.0 0.0 0.0 0.099
```

Lambda_est

```
## 4x2 Array{Float64,2}:
## 2.73949 1.78666
## 2.71648 1.76701
## -2.61815 3.98344
## -2.62815 4.01916
```

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

35 / 36

Summary

- Factor analysis model: $\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{f} + \boldsymbol{\epsilon}$
 - Latent variable \mathbf{f} is continuous
- Factor Analysis is useful to model data when it lies in a 'nearly' linear subspace
- Probabilistic formulation
 - $\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - $\mathbf{x} | \mathbf{f} \sim \mathcal{N}(\boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{f}, \boldsymbol{\Psi})$
- One approach for parameter estimation is MLE
 - Closed form solution is not possible
 - EM approach for modeling continuous latent variables
- Advantages:** Identifying hidden factors, Dimensionality reduction
- Disadvantages:** Indentifiability/Rotation problem, Data must follow Gaussian dist.

Gowtham Atluri

CS 5135/6035 Learning Probabilistic Models

October 12, 2018

36 / 36