

Week 7: The EM-Algorithm

MATH-517 Statistical Computation and Visualization

Tomas Masak

November 4th 2022

Section 1

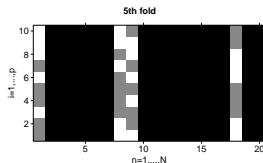
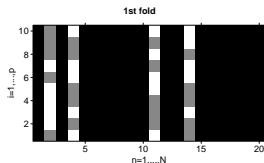
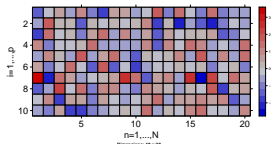
Motivation From the Last Week

CV for PCA Repaired

Assume that data $\mathbf{x}_n \in \mathbb{R}^p$ are i.i.d. realizations of $X \sim \mathcal{N}(\mu, \Sigma)$.

- split data into K folds J_1, \dots, J_K
- **for** $k = 1, \dots, K$
 - estimate μ and Σ empirically using all but the k -th fold J_k , but truncate Σ to be rank- r
 - **for** $n \in J_k$
 - split \mathbf{x}_n a “missing” part \mathbf{x}^{miss} that will be used for validation and an “observed” part \mathbf{x}^{obs}
 - predict \mathbf{x}_n^{miss} from \mathbf{x}_n^{obs} as discussed on the previous slide
 - **end for**
 - calculate $Err_k(r) = \sum_{n \in J_k} \|\mathbf{x}_n^{obs} - \hat{\mathbf{x}}_n^{obs}\|_2^2$
- **end for**
- choose $\hat{r} = \arg \min_r \sum_{k=1}^K |J_k|^{-1} Err_k(r)$

CV for PCA Repaired

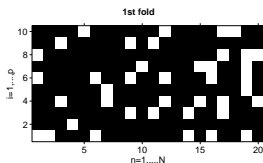


For every fold:

- use **black** entries to obtain $\hat{\mu}$ and $\hat{\Sigma}$
- predict **white** entries using **grey** entries and $\hat{\mu}$ and $\hat{\Sigma}$
- check the quality of your prediction

Improvements?

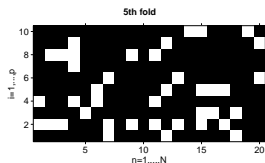
- Grey entries provide information on μ and Σ , shouldn't we use it?
- Isn't it awkward to first split rows and then columns? Why not to just split the bivariate index set?



.

.

.



To cope with this, we need to know how to do **MLE with missing data**.

Section 2

Expectation-Maximization (EM) Algorithm

EM Algorithm

Iterative algorithm for calculating Maximum-Likelihood-Estimators (MLEs) in situations, where

- there is **missing data** complicating the calculations (Example 1 and 3 below) or
- it is beneficial to think of our data as if there were some components missing (Example 2 below)
 - when knowing that missing components would render the problem simple

We will assume that solving MLE with the **complete data** is simple.

EM will allow us to act like if we knew everything – even when we don't or when we cannot use all the information.

Notation and the Algorithm

- \mathbf{X}_{obs} are the **observed** random variables
- \mathbf{X}_{miss} are the **missing** random variables
- $\ell_{comp}(\theta)$ is the **complete** log-likelihood of $\mathbf{X} = (\mathbf{X}_{obs}, \mathbf{X}_{miss})$
 - maximizing this to obtain MLE is supposed to be *simple*
 - θ denotes all the parameters, e.g. contains μ and Σ

Our task is to maximize $\ell_{obs}(\theta)$, the **observed** log-likelihood of \mathbf{X}_{obs} .

EM Algorithm: Start from an initial estimate $\hat{\theta}^{(0)}$ and for $l = 1, 2, \dots$ iterate the following two steps until convergence:

- **E-step:** calculate $\mathbb{E}_{\hat{\theta}^{(l-1)}} [\ell_{comp}(\theta) | \mathbf{X}_{obs} = \mathbf{x}_{obs}] =: Q(\theta, \hat{\theta}^{(l-1)})$
- **M-step:** optimize $\arg \max_{\theta} Q(\theta, \hat{\theta}^{(l-1)})$

Ex.1: Censored Observations

Suppose you want to estimate the mean waiting time at an EPFL food truck:

- observed waiting times $\mathbf{x}_{obs} = (x_{obs}^1, \dots, x_{obs}^{N_{obs}})^\top$ for \mathbf{X}_{obs} .
- food truck closes when N_{miss} individuals are still queuing, such that $\mathbf{X}_{miss} = (X_{miss}^1, \dots, X_{miss}^{N_{miss}})^\top$ are not observed but only a vector of right-censored waiting times $\tilde{\mathbf{x}}_{miss}$ with $\forall n : X_{miss}^{(n)} > \tilde{x}_{miss}^{(n)}$.
- overall $N = N_{obs} + N_{miss}$ individuals considered.

\Rightarrow Apply EM-algorithm assuming waiting times iid. distributed following an exponential distribution with density $f(x) = \lambda \exp(-\lambda x)$.

Ex.1: Censored Observations – E-step

- **E-step:** calculate

$$\mathbb{E}_{\hat{\theta}^{(l-1)}} [\ell_{comp}(\theta) | \mathbf{X}_{obs} = \mathbf{x}_{obs}, \forall n : X_{miss}^{(n)} > \tilde{x}_{miss}^{(n)}] =: Q(\theta, \hat{\theta}^{(l-1)})$$

For iterations $l = 1, \dots$:

$$\mathbb{E}_{\hat{\lambda}^{(l-1)}} [\ell_{comp}(\theta) | \mathbf{x}_{obs}, \tilde{\mathbf{x}}_{miss}] =$$

$$= \mathbb{E}_{\hat{\lambda}^{(l-1)}} \left[\underbrace{N \log(\lambda) - \lambda \sum_{n=1}^{N_{obs}} X_{obs}^{(n)} - \lambda \sum_{n=1}^{N_{miss}} X_{miss}^{(n)}}_{\log(\prod_{n=1}^{N_{obs}} f(X_{obs}^{(n)}) \cdot \prod_{n=1}^{N_{miss}} f(X_{miss}^{(n)}))} \mid \mathbf{x}_{obs}, \tilde{\mathbf{x}}_{miss} \right]$$

$$= N \log(\lambda) - \lambda \sum_{n=1}^{N_{obs}} x_{obs}^{(n)} - \lambda \sum_{n=1}^{N_{miss}} \underbrace{\mathbb{E}_{\hat{\lambda}^{(l-1)}} [X_{miss}^{(n)} \mid \tilde{\mathbf{x}}_{miss}]}_{\substack{X \sim \text{Exponential}(\hat{\lambda}^{(l-1)}) \\ = \\ \text{"memoryless"} \\ 1/\hat{\lambda}^{(l-1)} + \tilde{x}_{miss}^{(n)}}}$$

$$= N \log(\lambda) - \lambda \left(N_{obs} \bar{x}_{obs} + N_{miss} \frac{1}{\hat{\lambda}^{(l-1)}} + N_{miss} \bar{\tilde{x}}_{miss} \right) = Q(\lambda, \hat{\lambda}^{(l-1)})$$

Ex.1: Censored observations – M-step

- **M-step:** optimize $\arg \max_{\theta} Q(\theta, \hat{\theta}^{(l-1)})$

$$Q(\lambda, \hat{\lambda}^{(l-1)}) = N \log(\lambda) - \lambda \left(N_{obs} \bar{x}_{obs} + N_{miss} \frac{1}{\hat{\lambda}^{(l-1)}} + N_{miss} \bar{\bar{x}}_{miss} \right)$$

$$\Rightarrow \frac{\partial Q}{\partial \lambda}(\lambda, \hat{\lambda}^{(l-1)}) = \frac{N}{\lambda} - \left(N_{obs} \bar{x}_{obs} + N_{miss} \frac{1}{\hat{\lambda}^{(l-1)}} + N_{miss} \bar{\bar{x}}_{miss} \right) \stackrel{!}{=} 0$$

$$\Rightarrow \frac{1}{\hat{\lambda}^{(l)}} = \frac{N_{obs} \bar{x}_{obs} + N_{miss} \frac{1}{\hat{\lambda}^{(l-1)}} + N_{miss} \bar{\bar{x}}_{miss}}{N}$$

Ex.2: Mixture distributions

One of the most popular applications of the EM-algorithm:
Estimating mixture distributions for modelling multimodality

Mixture of two Gaussian distributions:

Let $X^{(1)}, \dots, X^{(N)}$ be iid. distributed as X with probability density

$$f(x) = (1 - \tau) \varphi\left(\frac{x - \mu_1}{\sigma_1}\right) + \tau \varphi\left(\frac{x - \mu_2}{\sigma_2}\right)$$

where

- φ is the density of a standard normal, and
- $\mu_1 < \mu_2$ and σ_1^2, σ_2^2 are the means and variances of the mixture components, and
- $\tau \in (0, 1)$ is the share of the second component stacked in a vector $\theta = (\tau, \mu_1, \mu_2, \sigma_1^2, \sigma_2^2)^\top$

Ex.2: Mixture distributions – factorization via latent variables

However, log-likelihood has no nice form:

$$\ell_{obs}(\theta) = \sum_{n=1}^N \log \left((1 - \tau) \varphi\left(\frac{x - \mu_1}{\sigma_1}\right) + \tau \varphi\left(\frac{x - \mu_2}{\sigma_2}\right) \right)$$

Trick: add latent iid. component indicators $Z_n \sim \text{Bernoulli}(\tau)$ such that $X^{(n)} \mid Z^{(n)} = 0 \sim N(\mu_1, \sigma_1^2)$ and $X^{(n)} \mid Z^{(n)} = 1 \sim N(\mu_2, \sigma_2^2)$.

Given $Z^{(n)} = z^{(n)}$, $n = 1, \dots, N$, the joint likelihood can be written as

$$L_{comp}(\theta) = (1 - \tau)^{N_1} \tau^{N_2} \prod_{n=1}^N \varphi\left(\frac{X^{(n)} - \mu_1}{\sigma_1}\right)^{(1-Z^{(n)})} \varphi\left(\frac{X^{(n)} - \mu_2}{\sigma_2}\right)^{Z^{(n)}}$$

with $N_2 = \sum_{n=1}^N Z^{(n)}$ and $N_1 = N - N_2$.

Ex.2: Mixture distributions – E-step – Part I

- **E-step:** calculate $\mathbb{E}_{\hat{\theta}^{(l-1)}}[\tilde{\ell}(\theta)|\mathbf{X} = \mathbf{x}] =: Q(\theta, \hat{\theta}^{(l-1)})$

$$\begin{aligned}\ell_{comp}(\theta) = \log L_{comp}(\theta) &= N_1 \log(1 - \tau) + N_2 \log(\tau) + \\ &+ \sum_{n=1}^N (1 - Z^{(n)}) \log \varphi\left(\frac{X^{(n)} - \mu_1}{\sigma_1}\right) + \sum_{n=1}^N Z^{(n)} \log \varphi\left(\frac{X^{(n)} - \mu_2}{\sigma_2}\right)\end{aligned}$$

such that, we obtain

$$\begin{aligned}\mathbb{E}_{\hat{\theta}^{(l-1)}}[\ell_{comp}(\theta)|\mathbf{X} = \mathbf{x}] &= \log(1 - \tau)(N - \sum_{n=1}^N p_n^{(l-1)}) + \log(\tau) \sum_{n=1}^N p_n^{(l-1)} + \\ &+ \sum_{n=1}^N (1 - p_n^{(l-1)}) \log \varphi\left(\frac{x^{(n)} - \mu_1}{\sigma_1}\right) + \sum_{n=1}^N p_n^{(l-1)} \log \varphi\left(\frac{x^{(n)} - \mu_2}{\sigma_2}\right)\end{aligned}$$

$$\text{with } p_n^{(l-1)} = \mathbb{E}_{\hat{\theta}^{(l-1)}}[Z^{(n)}|X^{(n)} = x^{(n)}] \stackrel{\text{Bayes}}{=} \frac{\varphi\left(\frac{x^{(n)} - \hat{\mu}_2^{(l-1)}}{\hat{\sigma}_2^{(l-1)}}\right)^{\hat{\tau}^{(l-1)}}}{f_{\hat{\theta}^{(l-1)}}(x^{(n)})}.$$

Ex.2: Mixture distributions – M-step

- **M-step:** optimize $\arg \max_{\theta} Q(\theta, \hat{\theta}^{(l-1)})$

Hence, $Q(\theta, \hat{\theta}^{(l-1)})$ nicely splits into three parts

$$Q(\theta, \hat{\theta}^{(l-1)}) =$$

$$\mathbf{A} : \quad \log(1 - \tau)(N - \sum_{n=1}^N p_n^{(l-1)}) + \log(\tau) \sum_{n=1}^N p_n^{(l-1)} +$$

$$\mathbf{B} : \quad + \sum_{n=1}^N (1 - p_n^{(l-1)}) \log \varphi \left(\frac{x^{(n)} - \mu_1}{\sigma_1} \right) +$$

$$\mathbf{C} : \quad + \sum_{n=1}^N p_n^{(l-1)} \log \varphi \left(\frac{x^{(n)} - \mu_2}{\sigma_2} \right)$$

which can be optimized separately, where **A** has the form of a binomial and **B** and **C** of (weighted) Gaussian log-likelihood \Rightarrow optimize accordingly.

Ex.3: Multivariate Gaussian with Missing Entries

Let $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(N)}$ iid. p -variate normally distributed with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$.

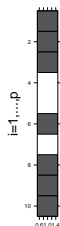
For each n , only a realization $\mathbf{x}_{obs}^{(n)}$ of a subvector $\mathbf{X}_{obs}^{(n)}$ of $\mathbf{X}^{(n)}$ is observed.

The goal is to estimate $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ from the incomplete measurements.

Ex.3: Multivariate Gaussian with Missing Entries

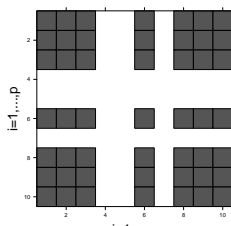
Let further denote $\mu_{obs}^{(n)}$ and $\Sigma_{obs}^{(n)}$ the mean and the covariance of $\mathbf{x}_{obs}^{(n)}$, i.e. $\mu_{obs}^{(n)}$ is just a sub-vector of μ and $\Sigma_{obs}^{(n)}$ is a sub-matrix of Σ .

$\mu_{obs}^{(n)}$



Dimensions: 10 x 1

$\Sigma_{obs}^{(n)}$



Dimensions: 10 x 10

Ex.3: Multivariate Gaussian with Missing Entries

Recall the density $f(\mathbf{x})$ of a p -variate Gaussian:

$$f(\mathbf{x}^{(n)}) \propto \det(\Sigma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x}^{(n)} - \mu)^\top \Sigma^{-1}(\mathbf{x}^{(n)} - \mu)\right),$$

and, hence, log-likelihood are given by

$$\ell_{obs}(\mu, \Sigma) = \text{const} - \frac{1}{2} \sum_{n=1}^N \log \det(\Sigma_{obs}^{(n)}) -$$

$$- \sum_{n=1}^N \frac{1}{2} (\mathbf{x}_{obs}^{(n)} - \mu_{obs}^{(n)}) (\Sigma_{obs}^{(n)})^{-1} (\mathbf{x}_{obs}^{(n)} - \mu_{obs}^{(n)})$$

$$\ell_{comp}(\mu, \Sigma) = \text{const} - \frac{N}{2} \log \det(\Sigma) - \sum_{n=1}^N \frac{1}{2} \underbrace{(\mathbf{x}^{(n)} - \mu)^\top \Sigma^{-1} (\mathbf{x}^{(n)} - \mu)}_{\text{tr}\left((\mathbf{x}^{(n)} - \mu)(\mathbf{x}^{(n)} - \mu)^\top \Sigma^{-1}\right)}.$$

Optimizing ℓ_{comp} easier than optimizing ℓ_{obs} . \Rightarrow EM-Algorithm.

Ex.3: Multivariate Gaussian with Missing Entries – E-step

- **E-step:** calculate $\mathbb{E}_{\hat{\theta}^{(l-1)}} [\ell_{comp}(\theta) | \forall n : \mathbf{X}_{obs}^{(n)} = \mathbf{x}_{obs}^{(n)}] =: Q(\theta, \hat{\theta}^{(l-1)})$ with $\theta = (\mu, \Sigma)$.

$$Q(\theta | \hat{\theta}^{(l-1)}) = \text{const} - \frac{N}{2} \log \det(\Sigma) - \\ - \sum_{n=1}^N \frac{1}{2} \text{tr} \left(\underbrace{\mathbb{E}_{\theta^{(l-1)}} \left[(\mathbf{X}^{(n)} - \mu)(\mathbf{X}^{(n)} - \mu)^\top \right] | \forall n : \mathbf{X}_{obs}^{(n)} = \mathbf{x}_{obs}^{(n)}}_{\substack{\text{some calculation} \\ \equiv (\hat{\mathbf{x}}^{(n)(l-1)} - \mu)(\hat{\mathbf{x}}^{(n)(l-1)} - \mu)^\top + \mathbf{C}^{(n)}}} \Sigma^{-1} \right)$$

with $\hat{\mathbf{x}}^{(n)(l-1)} = \mathbb{E}_{\hat{\theta}^{(l-1)}} [\mathbf{X}^{(n)} | \forall n : \mathbf{X}_{obs}^{(n)} = \mathbf{x}_{obs}^{(n)}]$ and

$$\mathbf{C}^{(n)} = \left\{ \text{Cov}_{\hat{\theta}^{(l-1)}} \left(X_i^{(n)}, X_j^{(n)} \mid \mathbf{X}_{obs}^{(n)} = \mathbf{x}_{obs}^{(n)} \right) \right\}_{i,j}.$$

Ex.3: Multivariate Gaussian with Missing Entries – M-step

- **M-step:** optimize $\arg \max_{\theta} Q(\theta, \hat{\theta}^{(l-1)})$

$$Q(\theta, \hat{\theta}^{(l-1)}) = \text{const} - \frac{N}{2} \log \det(\Sigma) - \\ - \sum_{n=1}^N \frac{1}{2} \text{tr} \left((\hat{\mathbf{x}}^{(n)(l-1)} - \mu)(\hat{\mathbf{x}}^{(n)(l-1)} - \mu)^{\top} \Sigma^{-1} \right) - \frac{1}{2} \text{tr}(\mathbf{C} \Sigma^{-1}),$$

has a similar form as a multivariate normal and estimators can be derived accordingly, resulting in

$$\hat{\mu}^{(l)} = N^{-1} \sum_{n=1}^N \hat{\mathbf{x}}^{(n)(l-1)}$$

and

$$\hat{\Sigma}^{(l)} = \frac{1}{N} \sum_{n=1}^N [(\hat{\mathbf{x}}^{(n)(l-1)} - \mu)(\hat{\mathbf{x}}^{(n)(l-1)} - \mu)^{\top} + \mathbf{C}^{(n)}].$$

Recap

Example 1:

- part of data missing but their censored versions carry some information
- the likelihood is linear (w.r.t. observations) and thus the **E-step** coincides with imputation (missing data replaced by their expectations)
 - this is rare!

Example 2:

- there is no true missing data here, but it is beneficial to imagine it
- the likelihood is linear w.r.t. the imagined observations \Rightarrow simplification

Example 3:

- likelihood of observed data easy to formulate, yet harder to optimize directly
- no linearity in log-likelihood \Rightarrow no imputation, more effort to compute expected likelihood

- Dempster, A. P., N. M. Laird & D. B. Rubin. (1977) “Maximum likelihood from incomplete data via the EM algorithm.” *Journal of the Royal Statistical Society: Series B (Methodological)* 39.1: 1-22.
 - one of the most cited papers in statistics of all time
- Little, R. J., & Rubin, D. B. (2019). *Statistical analysis with missing data*. 3rd Edition.

Assignment 5 [5 %]

Simulate data (using Manual 9) from the mixture of two Gaussian distributions and implement the EM algorithm from Example 2 above. Use absolute change in ℓ of observed data as convergence criterion.

The following should naturally be done, though it is not mandatory:

- Visualize the resulting parametric density estimators.
- Try running your algorithm from different starting points.
 - How sensitive is your algorithm to your choice of starting point?
 - Can you find a bad starting point where your algorithm fails?

Project Ideas

- Cross-validation for PCA
 - A simulation study to compare the advantages of EM compared to what we did last week.
 - Two more approaches to CV for PCA to be considered and compared.
- Comparison of local regression implementations in different R packages.
 - [Wickham \(2011\)](#) examines several packages for KDE in R. Not only there are huge differences in terms of speed, but some of the packages are even inconsistent! Makes one wonder what is the situation with local regression.
- Various simulation studies on bandwidth selection in KDE or local linear regression.
- Comparison of variable selectors in regression.
 - [Hastie et al. \(2020\)](#) have some surprising results in their simulation study, but one important method (adaptive lasso) is omitted. Try to recreate the study with adaptive lasso included.