

HMM - Preliminaries

Covariance structures

Expectation-maximization (EM)

Parameters estimation in GMM... in 1893

III. Contributions to the Mathematical Theory of Evolution.

By KARL PEARSON, University College, London.

Communicated by Professor HENRICI, F.R.S.

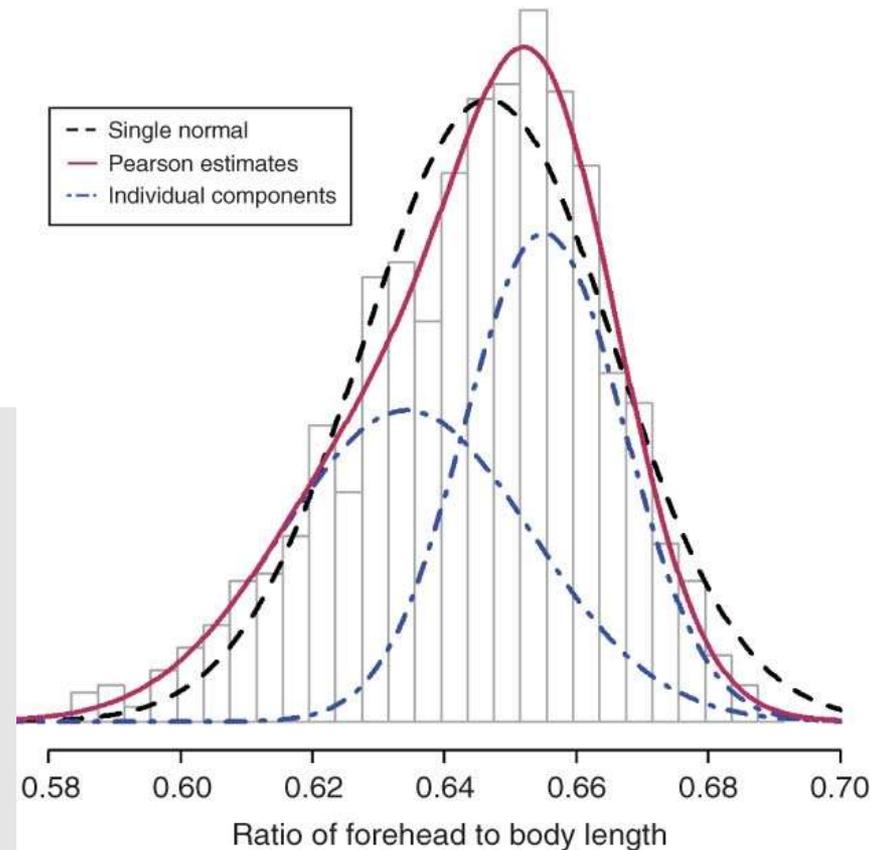
Received October 18,—Read November 16, 1893.

[PLATES 1—5.]

CONTENTS.

	Page.
—On the Dissection of Asymmetrical Frequency-Curves. General Theory, §§ 1–8.	71–85
Example: Professor WELDON'S measurements of the "Forehead" of Crabs.	
§§ 9–10	85–90
—On the Dissection of Symmetrical Frequency-Curves. General Theory, §§ 11–12	

Pearson's crab data and fitted mixtures



Parameters estimation in GMM... in 1893

54 pages!
Proposed solution:
Moment-based approach
requiring to solve a
polynomial of degree 9...

... which does not mean that moment-based approaches are old-fashioned!
They are actually today popular again with new developments related to spectral decomposition.

Gaussian Mixture Model (GMM)

K Gaussians
 N datapoints of
dimension D

$$\mathcal{P}(\boldsymbol{\xi}_t) = \sum_{i=1}^K \pi_i \mathcal{N}(\boldsymbol{\xi}_t | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

$$\mathcal{N}(\boldsymbol{\xi}_t | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) = \frac{1}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}_i|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}_i^{-1}(\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)\right)$$

$\boldsymbol{\xi} \in \mathbb{R}^{D \times N}$ Observations

$\pi_i \in \mathbb{R}$ Mixing coefficient

$\boldsymbol{\mu}_i \in \mathbb{R}^D$ Center (me.

$\boldsymbol{\Sigma}_i \in \mathbb{R}^{D \times D}$ Covariance

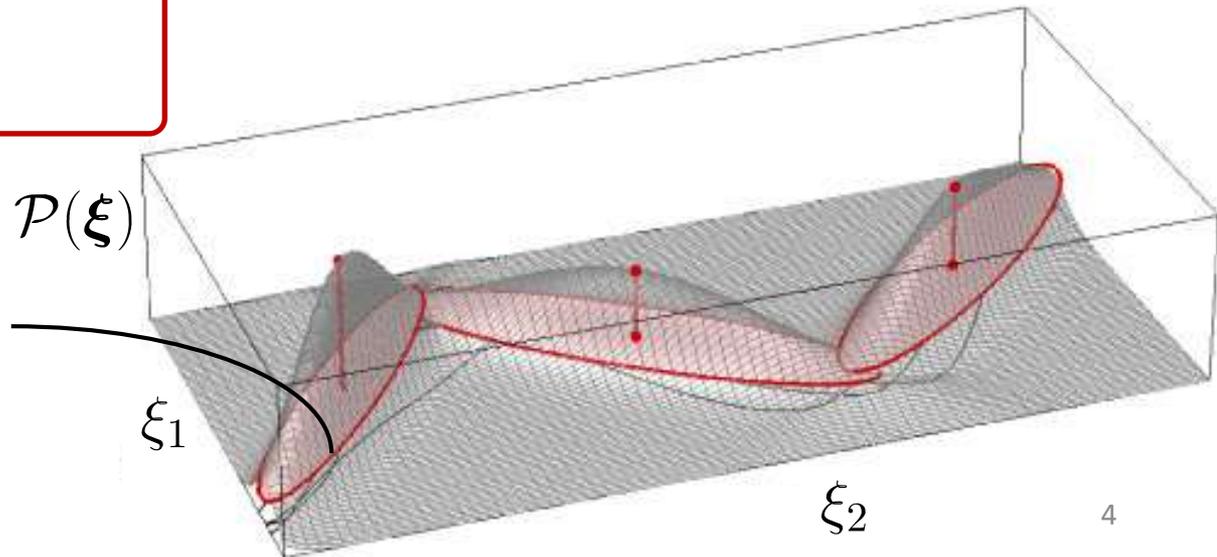
Parameters $\Theta^{\text{GMM}} = \{\pi_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^K$

Equidensity contour of
one standard deviation

$\mathcal{P}(\boldsymbol{\xi})$

ξ_1

ξ_2



Expectation-maximization (EM)

$z_{t,i} = 1$ if ξ_t is part of cluster i . It is 0 otherwise.

Each datapoint ξ_t is associated with a hidden/missing variable z_t .
The goal is to maximize the log-likelihood of the observed data

$$\mathcal{L}(\Theta) = \sum_{t=1}^N \log \mathcal{P}(\xi_t | \Theta) = \sum_{t=1}^N \log \left(\sum_{z_t} \mathcal{P}(\xi_t, z_t | \Theta) \right)$$

which is hard to optimize (“log cannot be pushed inside the sum”).

We can get around this problem by instead employing the expected complete data log-likelihood

$$Q(\Theta, \Theta^{\text{old}}) = \mathbb{E} \left[\sum_{t=1}^N \log \mathcal{P}(\xi_t, z_t | \Theta) \mid \xi, \Theta^{\text{old}} \right]$$

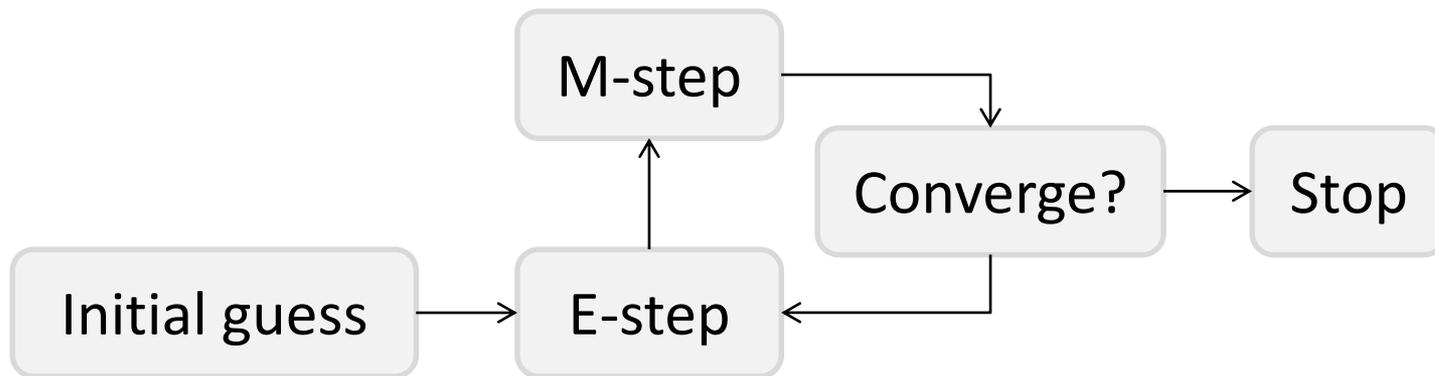
where $Q(\Theta, \Theta^{\text{old}})$ is called the auxiliary function.

Expectation-maximization (EM)

The expectation is taken with respect to the old model parameters Θ^{old} and the observed dataset ξ .

The *E-step* computes the terms in $Q(\Theta, \Theta^{\text{old}})$ of which the likelihood depends on, known as the expected sufficient statistics.

The *M-step* then optimizes Q with respect to Θ .



EM for GMM

Setting

$$\frac{\partial Q(\Theta, \Theta^{\text{old}})}{\partial \pi_i} = 0 \quad \frac{\partial Q(\Theta, \Theta^{\text{old}})}{\partial \mu_i} = 0 \quad \frac{\partial Q(\Theta, \Theta^{\text{old}})}{\partial \Sigma_i} = 0$$

and solving for π_i , μ_i and Σ_i results in an EM procedure to compute the maximum likelihood estimate of the parameters.

EM for GMM: Resulting procedure

K Gaussians
 N datapoints

E-step:

$$h_{t,i} = \frac{\pi_i \mathcal{N}(\boldsymbol{\xi}_t | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{\sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{\xi}_t | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$

M-step:

$$\pi_i \leftarrow \frac{\sum_{t=1}^N h_{t,i}}{N},$$

$$\boldsymbol{\mu}_i \leftarrow \frac{\sum_{t=1}^N h_{t,i} \boldsymbol{\xi}_t}{\sum_{t=1}^N h_{t,i}},$$

$$\boldsymbol{\Sigma}_i \leftarrow \frac{\sum_{t=1}^N h_{t,i} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)(\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top}{\sum_{t=1}^N h_{t,i}}$$

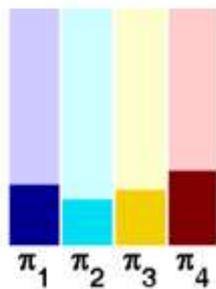
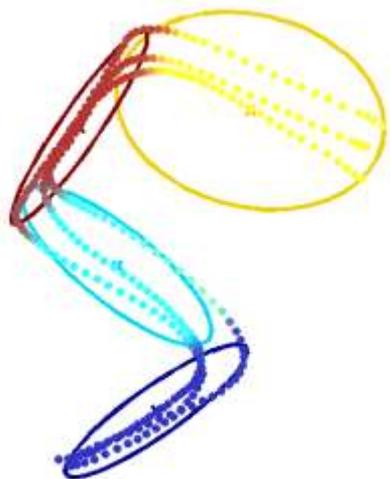
These results can be intuitively interpreted in terms of normalized counts.

EM provides a systematic approach to derive such procedure.

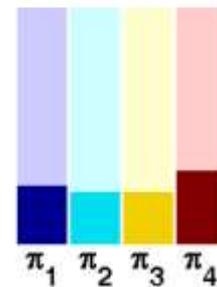
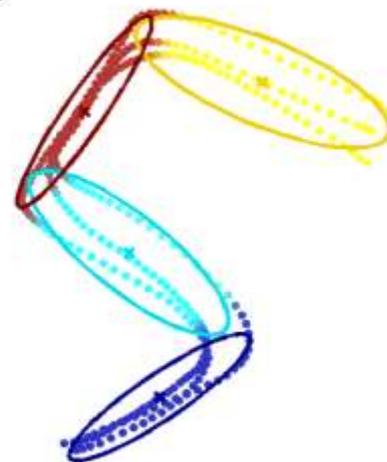
→ **Weighted averages taking into account the responsibility of each datapoint in each cluster.**

EM for GMM

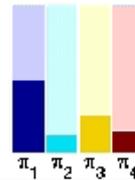
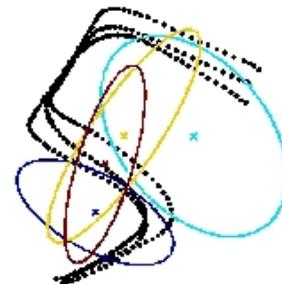
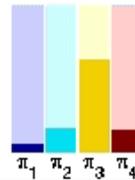
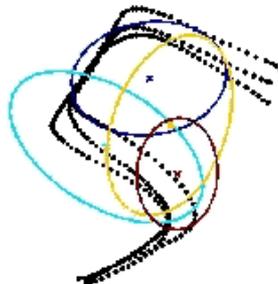
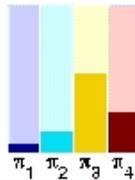
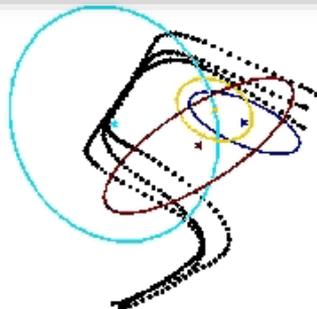
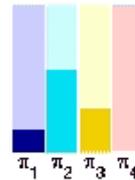
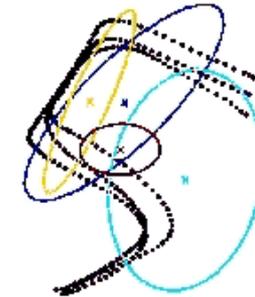
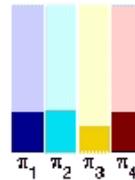
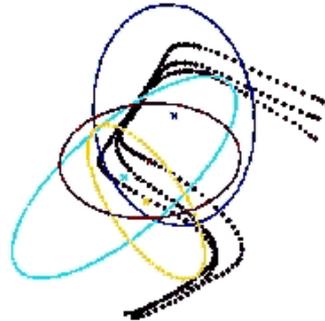
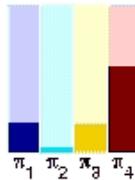
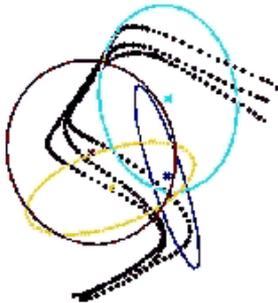
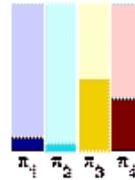
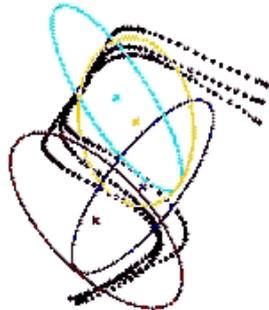
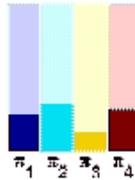
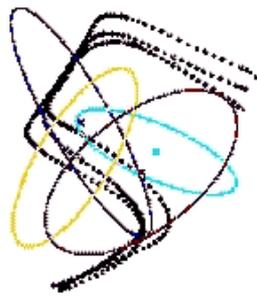
E-step



M-step

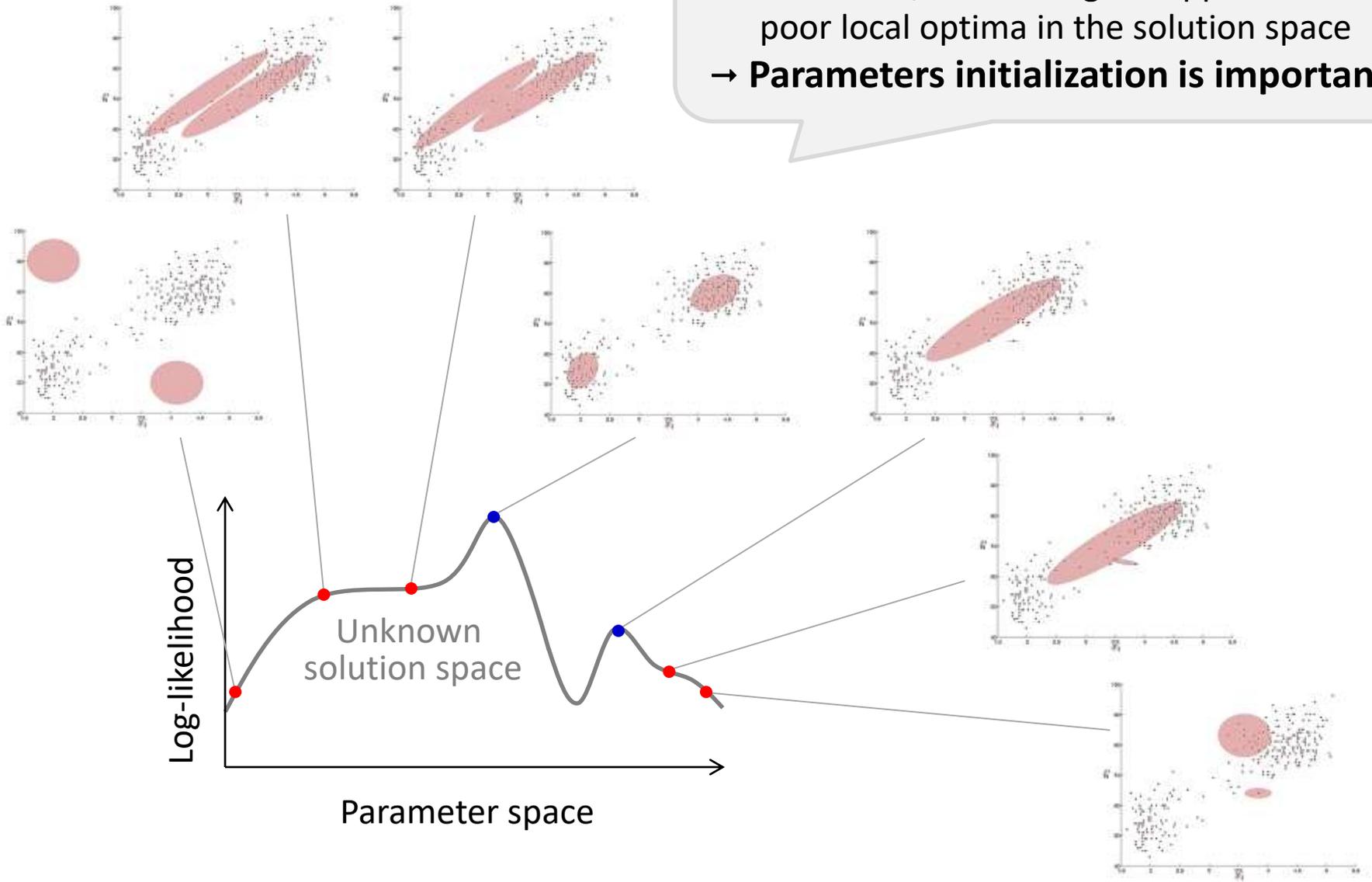


EM for GMM: Local optima issue



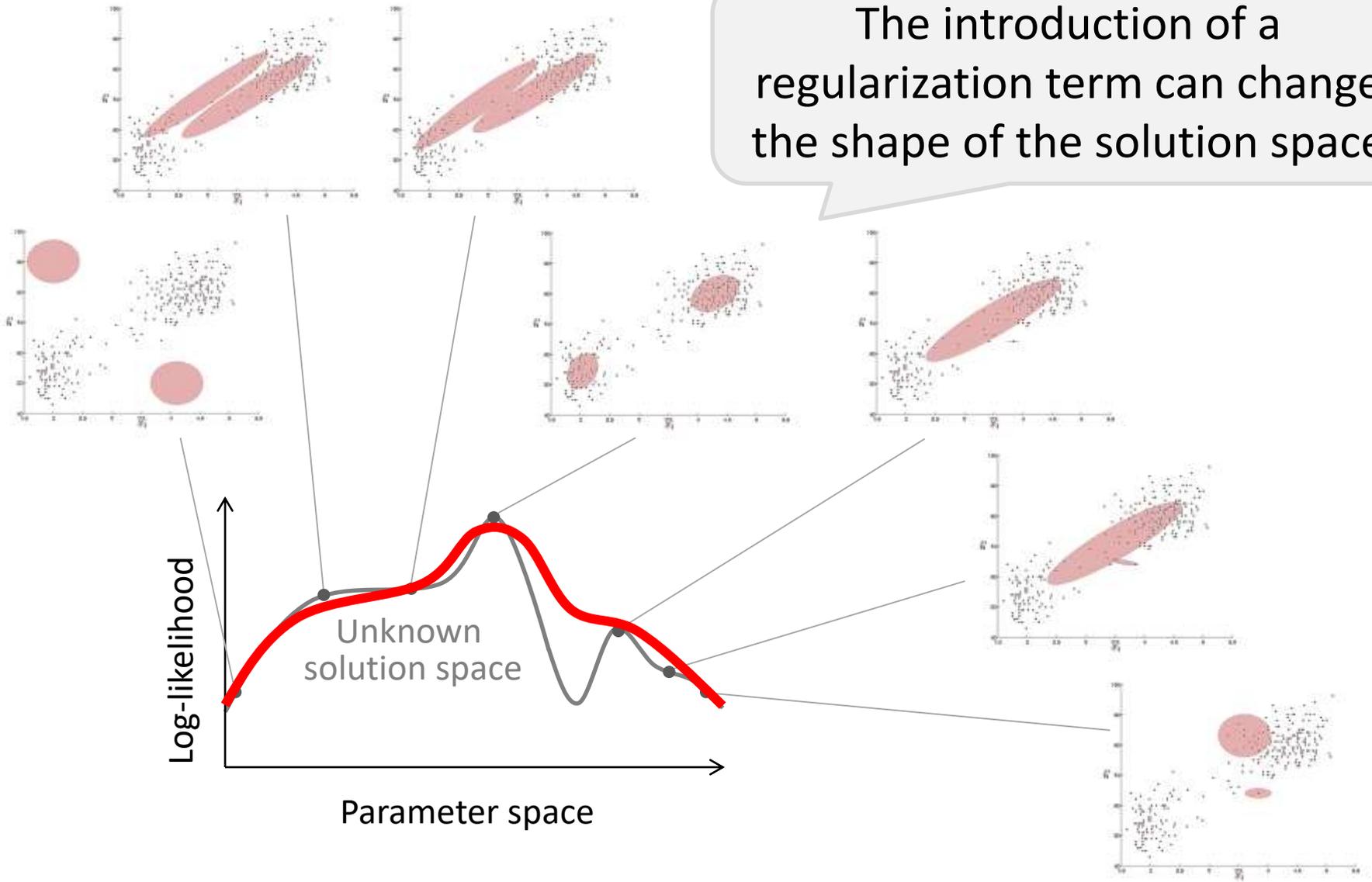
Local optima in EM

EM will improve the likelihood at each iteration, but it can get trapped into poor local optima in the solution space
→ **Parameters initialization is important!**



Regularization of the GMM parameters

The introduction of a regularization term can change the shape of the solution space



Regularization of the GMM parameters

Regularization with minimal admissible eigenvalue:

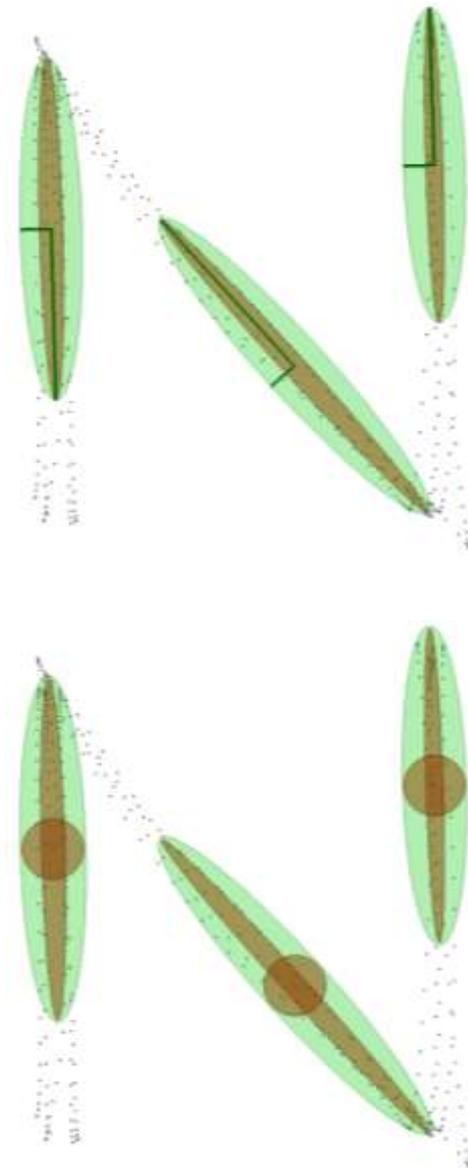
$$\Sigma_i \leftarrow V_i \tilde{D}_i V_i^\top$$

$$\text{with } \tilde{D}_i = \begin{bmatrix} \tilde{\lambda}_{i,1}^2 & 0 & \cdots & 0 \\ 0 & \tilde{\lambda}_{i,2}^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \tilde{\lambda}_{i,D}^2 \end{bmatrix}$$

$$\text{and } \tilde{\lambda}_{i,j}^2 = \max(\lambda_{i,j}^2, \lambda_{\min}^2) \quad \forall j \in \{1, \dots, D\}$$

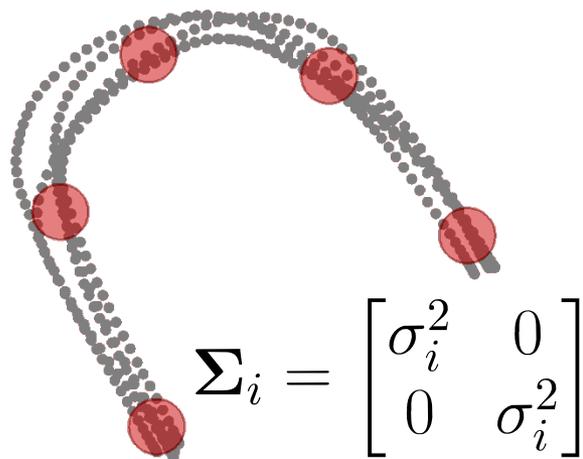
Tikhonov regularization with isotropic covariance:

$$\Sigma_i \leftarrow \Sigma_i + I \lambda_{\min}^2$$

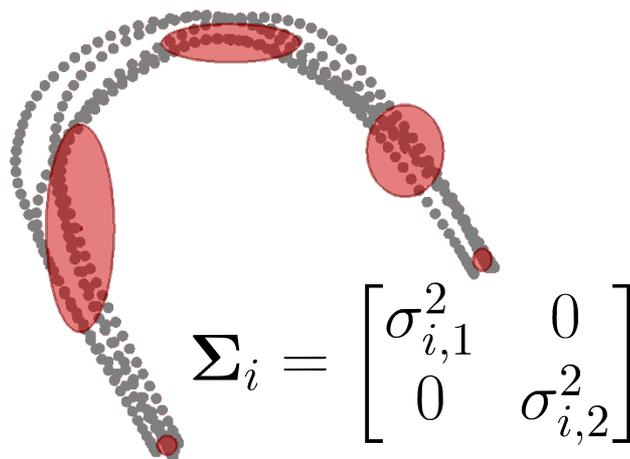


Covariance structures in GMM

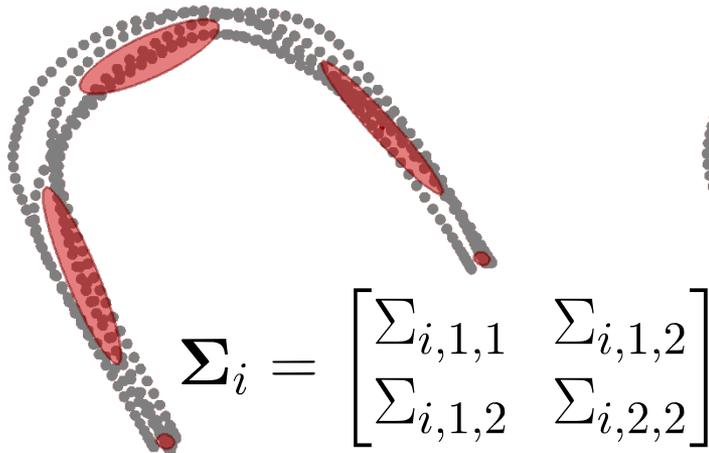
Isotropic



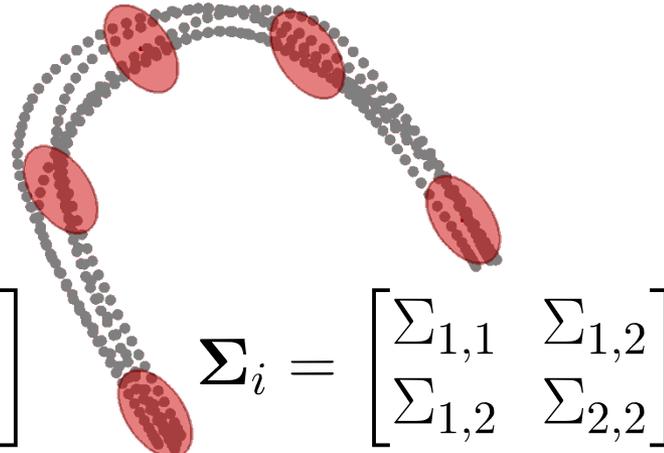
Diagonal



Full



Tied



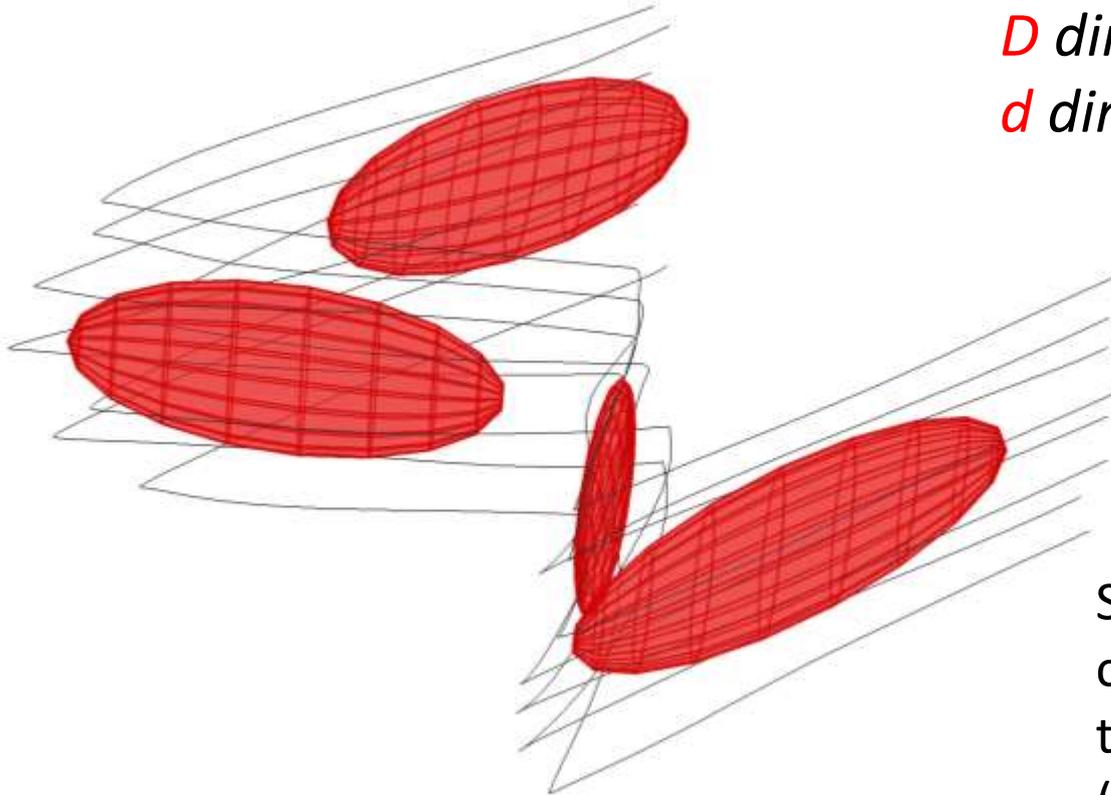
Subspace clustering

K clusters

N datapoints

D dimensions (original space)

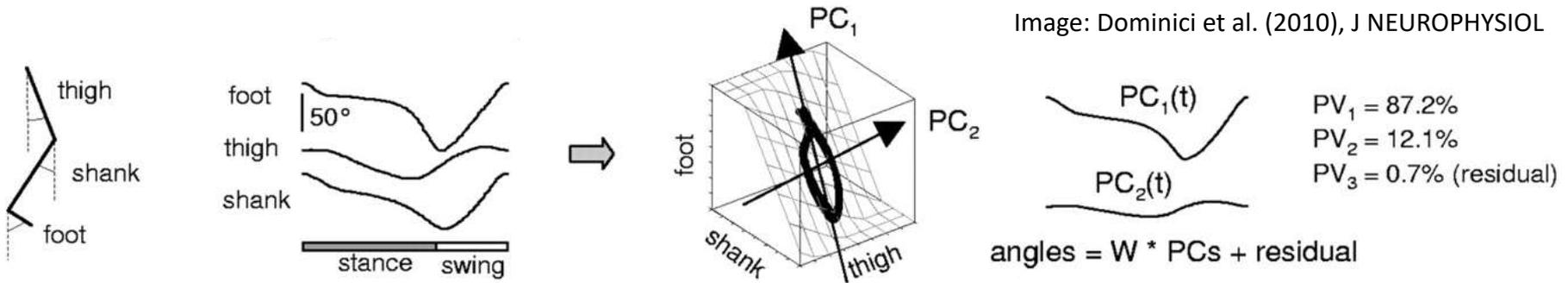
d dimensions (latent space)



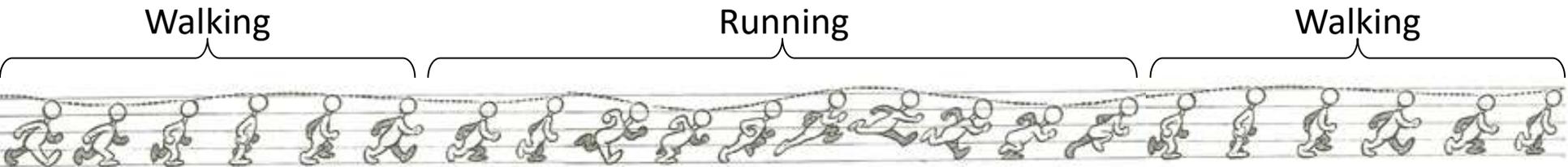
Subspace clustering aims at clustering data while reducing the dimension of each cluster (cluster-dependent subspace)

Considering the two problems separately (clustering, then subspace projection) can be inefficient and can produce poor local optima, especially when datapoints of high dimensions are considered.

Example of application: Whole body motion

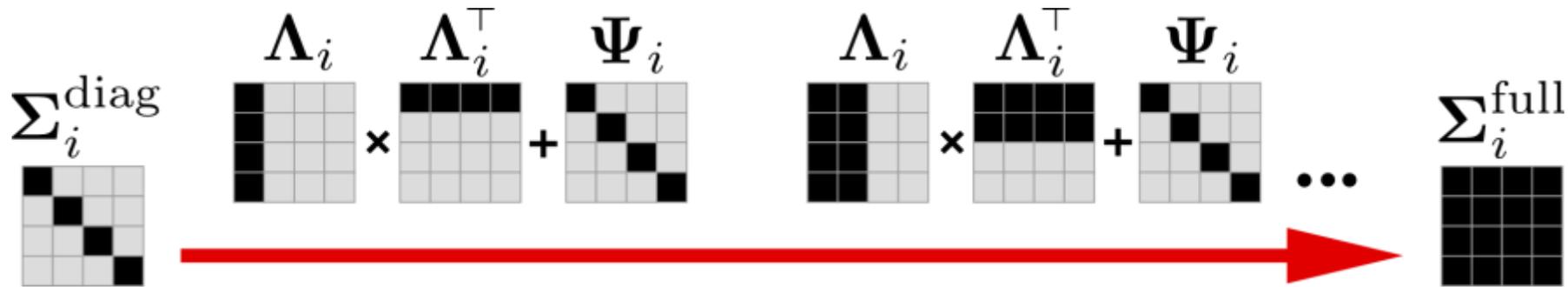


- About 90% of variance in walking motion can be explained by 2 principal components
- Each type of periodic motion can be characterized by a different subspace



- Requires clustering of the complete motion into different locomotion phases
- Requires extraction of coordination patterns for each cluster

Mixture of factor analyzers (MFA)



MFA assumes for each covariance i a structure of the form

$$\Sigma_i = \Lambda_i \Lambda_i^\top + \Psi_i$$

where $\Lambda_i \in \mathbb{R}^{D \times d}$, known as the *factor loading matrix*, typically has $d < D$ (providing a parsimonious representation of the data), and a diagonal noise matrix Ψ_i .

The *mixture of probabilistic principal component analyzers* (MPPCA) is a special case of MFA with the distribution of the errors assumed to be isotropic with $\Psi_i = \mathbf{I} \sigma_i^2$.

A taxonomy of parsimonious GMMs

D is used in this lecture

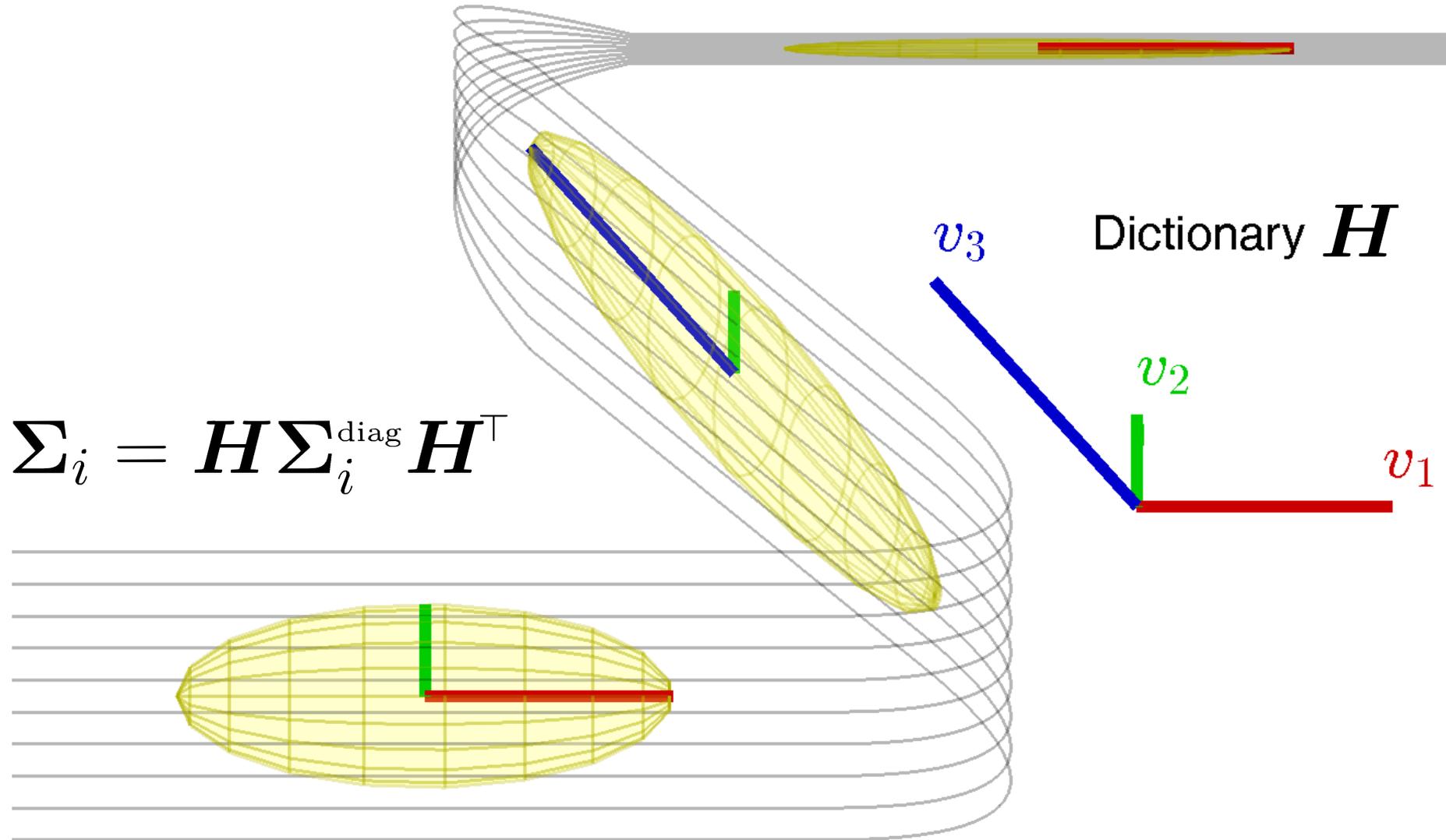
$K = 4, d = 3$
 $p = 100$

Model name	Cov. structure	Nb. of parameters	
UUUU - UUU	$S_k = \Lambda_k \Lambda_k^t + \Psi_k$	$(K - 1) + Kp + Kd[p - (d - 1)/2] + Kp$	1991
UUCU -	$S_k = \Lambda_k \Lambda_k^t + \omega_k \Delta_k$	$(K - 1) + Kp + Kd[p - (d - 1)/2] + [1 + K(p - 1)]$	1988
UCUU -	$S_k = \Lambda_k \Lambda_k^t + \omega_k \Delta$	$(K - 1) + Kp + Kd[p - (d - 1)/2] + [K + (p - 1)]$	1694
UCCU - UCU	$S_k = \Lambda_k \Lambda_k^t + \Psi$	$(K - 1) + Kp + Kd[p - (d - 1)/2] + p$	1691
UCUC - UUC	$S_k = \Lambda_k \Lambda_k^t + \psi_k \mathbf{I}_p$	$(K - 1) + Kp + Kd[p - (d - 1)/2] + K$	1595
UCCC - UCC	$S_k = \Lambda_k \Lambda_k^t + \psi \mathbf{I}_p$	$(K - 1) + Kp + Kd[p - (d - 1)/2] + 1$	1592
CUUU - CUU	$S_k = \Lambda \Lambda^t + \Psi_k$	$(K - 1) + Kp + d[p - (d - 1)/2] + Kp$	1100
CUCU -	$S_k = \Lambda \Lambda^t + \omega \Delta_k$	$(K - 1) + Kp + d[p - (d - 1)/2] + [1 + K(p - 1)]$	1097
CCUU -	$S_k = \Lambda \Lambda^t + \omega_k \Delta$	$(K - 1) + Kp + d[p - (d - 1)/2] + [K + (p - 1)]$	803
CCCU - CCU	$S_k = \Lambda \Lambda^t + \Psi$	$(K - 1) + Kp + d[p - (d - 1)/2] + p$	800
CCUC - CUC	$S_k = \Lambda \Lambda^t + \psi_k \mathbf{I}_p$	$(K - 1) + Kp + d[p - (d - 1)/2] + K$	704
CCCC - CCC	$S_k = \Lambda \Lambda^t + \psi \mathbf{I}_p$	$(K - 1) + Kp + d[p - (d - 1)/2] + 1$	701

where $\omega_k \in \mathbb{R}^+$ and $|\Delta_k| = 1$.

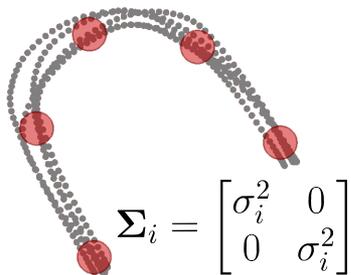
[C. Bouveyron and C. Brunet. Model-based clustering of high-dimensional data: A review. Computational Statistics and Data Analysis, 71:52–78, March 2014]

Sharing of parameters in mixture models

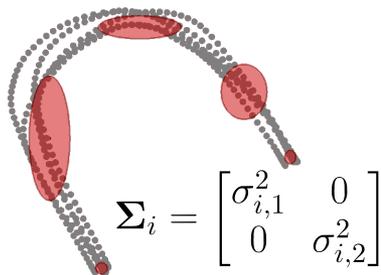


Summary of relevant covariance structures

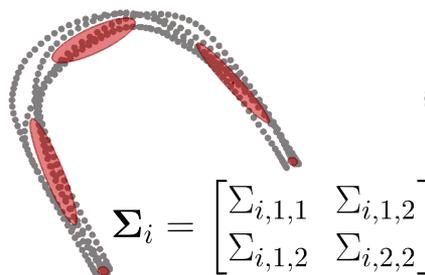
Isotropic



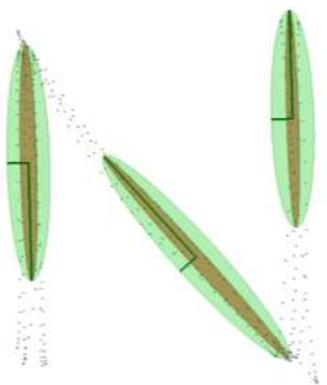
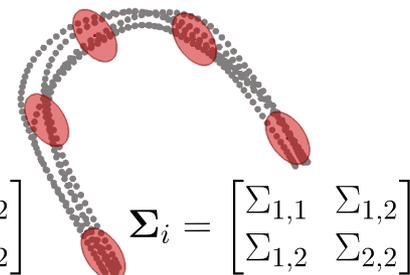
Diagonal



Full



Tied

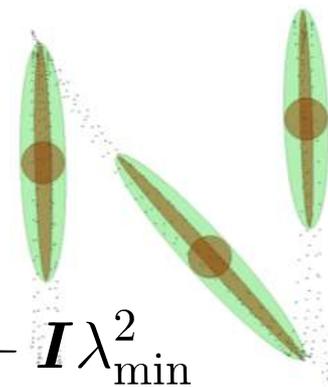


$$\Sigma_i \leftarrow V_i \tilde{D}_i V_i^\top \text{ with}$$

$$\tilde{D}_i = \text{diag}(\tilde{\lambda}_{i,1}^2, \dots, \tilde{\lambda}_{i,D}^2) \text{ and}$$

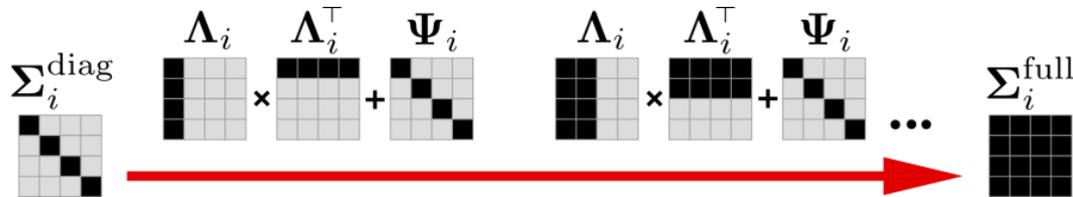
$$\tilde{\lambda}_{i,j}^2 = \max(\lambda_{i,j}^2, \lambda_{\min}^2)$$

$$\Sigma_i \leftarrow \Sigma_i + I \lambda_{\min}^2$$

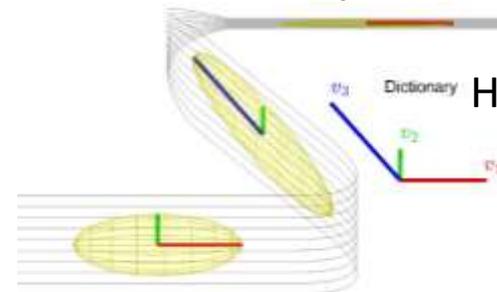


$$\text{MFA: } \Sigma_i = \Lambda_i \Lambda_i^\top + \Psi_i$$

$$\text{MPPCA: } \Sigma_i = \Lambda_i^\top \Lambda_i + I \sigma_i^2$$



$$\Sigma_i = H \Sigma_i^{\text{diag}} H^\top$$



References

Parsimonious GMM

C. Bouveyron and C. Brunet. Model-based clustering of high-dimensional data: A review. *Computational Statistics and Data Analysis*, 71:52–78, March 2014

P. D. McNicholas and T. B. Murphy. Parsimonious Gaussian mixture models. *Statistics and Computing*, 18(3):285–296, September 2008

MFA

G. J. McLachlan, D. Peel, and R. W. Bean. Modelling high-dimensional data by mixtures of factor analyzers. *Computational Statistics and Data Analysis*, 41(3-4):379–388, 2003

G. E. Hinton, P. Dayan, and M. Revow. Modeling the manifolds of images of handwritten digits. *IEEE Trans. on Neural Networks*, 8(1):65–74, 1997

MPPCA

M. E. Tipping and C. M. Bishop. Mixtures of probabilistic principal component analyzers. *Neural Computation*, 11(2):443–482, 1999

GMM with semi-tied covariances

M. J. F. Gales. Semi-tied covariance matrices for hidden Markov models. *IEEE Trans. on Speech and Audio Processing*, 7(3):272–281, 1999

Labs

Teguh Lembono



Python notebooks and labs exercises:

<https://github.com/teguhSL/ee613-python>

Branch: master ▾ ee613-python / python_notebooks / linear_regression_1 / Create new file Find file History

 teguhSL minor edits Latest commit c1da0e0 9 hours ago

..		
 Ex1.ipynb	minor edits	9 hours ago
 Ex2.ipynb	minor edits	9 hours ago
 Ex3.ipynb	minor edits	9 hours ago
 demo_LS.ipynb	minor edits	9 hours ago
 demo_LS_polFit.ipynb	minor edits	9 hours ago
 demo_LS_recursive.ipynb	minor edits	9 hours ago
 demo_LS_weighted.ipynb	minor edits	9 hours ago

Appendix

EM for GMM

When applied to GMM, the auxiliary function $Q(\Theta, \Theta^{\text{old}})$ takes the form

$$\begin{aligned}
 Q(\Theta, \Theta^{\text{old}}) &= \mathbb{E} \left[\sum_{t=1}^N \log \mathcal{P}(\xi_t, z_t | \Theta) \mid \xi, \Theta^{\text{old}} \right] \\
 &= \sum_{t=1}^N \mathbb{E} \left[\log \left(\prod_{i=1}^K (\pi_i \mathcal{N}(\xi_t | \mu_i, \Sigma_i))^{z_{t,i}} \right) \mid \xi, \Theta^{\text{old}} \right] \\
 &= \sum_{t=1}^N \sum_{i=1}^K \mathbb{E} \left[\log \left((\pi_i \mathcal{N}(\xi_t | \mu_i, \Sigma_i))^{z_{t,i}} \right) \mid \xi, \Theta^{\text{old}} \right] \\
 &= \sum_{t=1}^N \sum_{i=1}^K \mathbb{E}[z_{t,i} \mid \xi, \Theta^{\text{old}}] \log \left(\pi_i \mathcal{N}(\xi_t | \mu_i, \Sigma_i) \right) \\
 &= \sum_{t=1}^N \sum_{i=1}^K h_{t,i} \left(\log(\pi_i) + \log \left(\mathcal{N}(\xi_t | \mu_i, \Sigma_i) \right) \right) \\
 &= \sum_{t=1}^N \sum_{i=1}^K h_{t,i} \left(\log(\pi_i) - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\xi_t - \mu_i)^\top \Sigma_i^{-1} (\xi_t - \mu_i) - \frac{D}{2} \log(2\pi) \right)
 \end{aligned}$$

$z_{t,i} = 1$ if ξ_i is part of cluster i .
 It is 0 otherwise.

e.g. $\prod_{i=1}^3 \pi_i^{z_i} = \pi_1^{z_1} \cdot \pi_2^{z_2} \cdot \pi_3^{z_3} = \pi_1^0 \cdot \pi_2^0 \cdot \pi_3^1 = 1 \cdot 1 \cdot \pi_3$

$\log(ab) = \log(a) + \log(b)$

$\log(a^b) = b \log(a)$

$\log(\exp(a)) = a$

$\mathcal{N}(\xi_t | \mu_i, \Sigma_i) = (2\pi)^{-\frac{D}{2}} |\Sigma_i|^{-\frac{1}{2}} \exp \left(-\frac{1}{2} (\xi_t - \mu_i)^\top \Sigma_i^{-1} (\xi_t - \mu_i) \right)$

where $h_{t,i}$ is the responsibility that cluster i takes for datapoint ξ_t .

EM for GMM

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{t=1}^N \sum_{i=1}^K h_{t,i} \left(\log(\pi_i) - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i) - \frac{D}{2} \log(2\pi) \right)$$

By using the linear algebra relations

$$\frac{\partial}{\partial \mathbf{A}} \log |\mathbf{A}| = (\mathbf{A}^\top)^{-1} \quad \frac{\partial}{\partial \mathbf{A}} \mathbf{x}^\top \mathbf{A} \mathbf{x} = \mathbf{x} \mathbf{x}^\top \quad \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{A} \mathbf{x} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

(= 2 $\mathbf{A}\mathbf{x}$ if \mathbf{A} symmetric)

and the derivation chain rule, we obtain

$$\frac{\partial Q(\Theta, \Theta^{\text{old}})}{\partial \boldsymbol{\mu}_i} = \frac{1}{2} \sum_{t=1}^N h_{t,i} 2 \Sigma_i^{-1} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i) = \Sigma_i^{-1} \sum_{t=1}^N h_{t,i} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i) = 0$$

$$\iff \boldsymbol{\mu}_i = \frac{\sum_{t=1}^N h_{t,i} \boldsymbol{\xi}_t}{\sum_{t=1}^N h_{t,i}}$$

$$\frac{\partial Q(\Theta, \Theta^{\text{old}})}{\partial \Sigma_i} = \frac{1}{2} \Sigma_i \sum_{t=1}^N h_{t,i} - \frac{1}{2} \sum_{t=1}^N h_{t,i} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i) (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top = 0$$

$$\iff \Sigma_i = \frac{\sum_{t=1}^N h_{t,i} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i) (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top}{\sum_{t=1}^N h_{t,i}}$$

EM for GMM $Q(\Theta, \Theta^{\text{old}}) = \sum_{t=1}^N \sum_{i=1}^K h_{t,i} \left(\log(\pi_i) - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i) - \frac{D}{2} \log(2\pi) \right)$

For π_i , we need to ensure the constraint $\sum_{i=1}^K \pi_i = 1$, which can be achieved through a Lagrange multiplier λ , yielding

$$\frac{\partial}{\partial \pi_i} \left[Q(\Theta, \Theta^{\text{old}}) - \lambda \left(\sum_{i=1}^K \pi_i - 1 \right) \right] = \frac{1}{\pi_i} \sum_{t=1}^N h_{t,i} - \lambda = 0$$

The sum over K of the above relation provides

$$\sum_{t=1}^N \sum_{i=1}^K h_{t,i} = \lambda \sum_{i=1}^K \pi_i \quad \iff \quad \lambda = N$$

$\sum_{i=1}^K h_{t,i} = 1, \sum_{i=1}^K \pi_i = 1$

which can be reintroduced in the equation to find

$$\frac{1}{\pi_i} \sum_{t=1}^N h_{t,i} - N = 0 \quad \iff \quad \pi_i = \frac{\sum_{t=1}^N h_{t,i}}{N}$$

Mixture of factor analyzers (MFA) $\Sigma_i = \Lambda_i \Lambda_i^\top + \Psi_i$

$$\Sigma = VD^{\frac{1}{2}}(VD^{\frac{1}{2}})^\top$$

$$\xi \sim \mu + VD^{\frac{1}{2}} \mathcal{N}(\mathbf{0}, \mathbf{I})$$



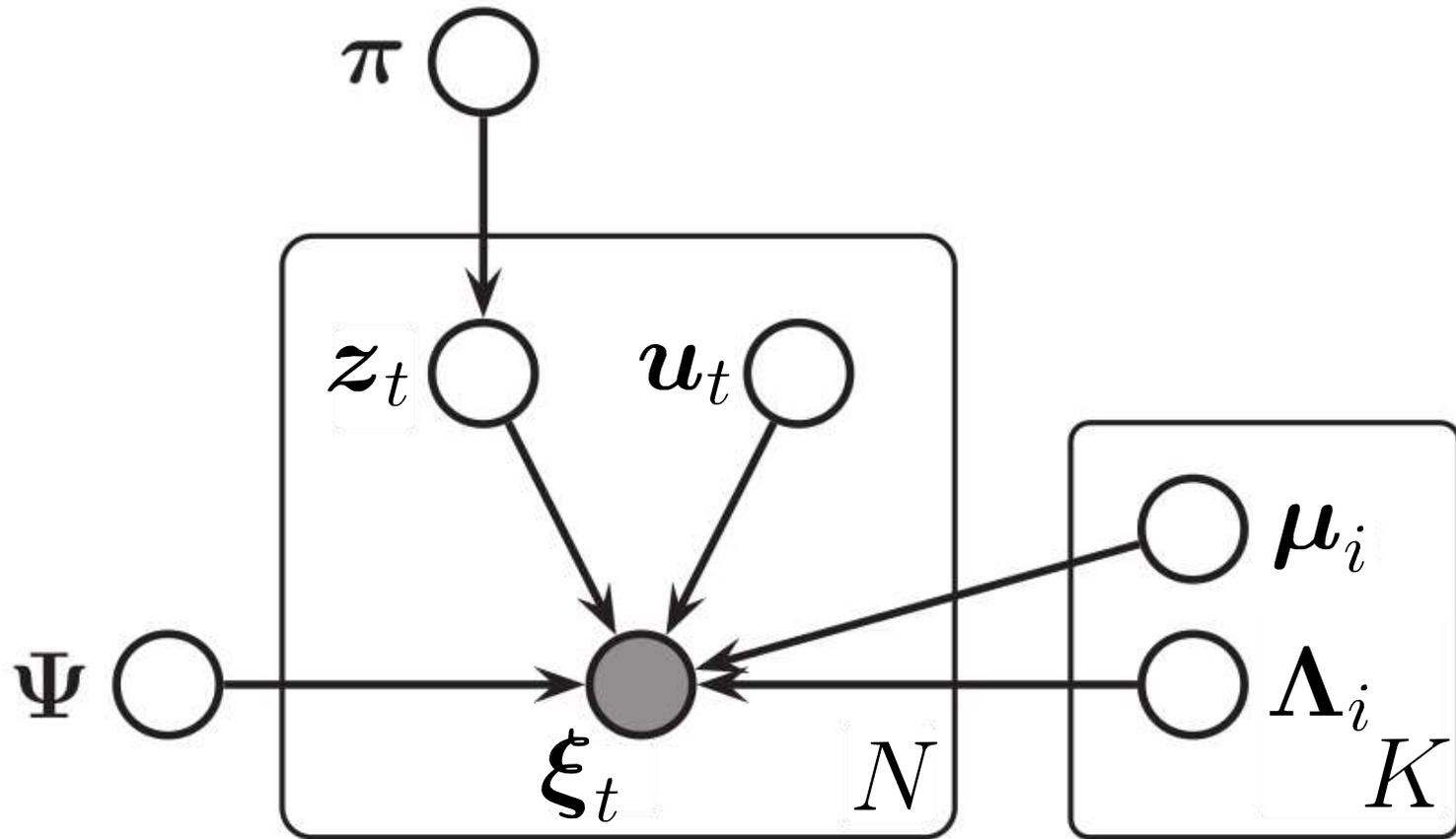
In MFA, the generative model for the i -th mixture component assumes that a D -dimensional random vector ξ is modeled using a d -dimensional vector of latent (unobserved) factors \mathbf{u}

$$\xi = \Lambda_i \mathbf{u} + \mu_i + \epsilon_i$$

where $\mu_i \in \mathbb{R}^D$ is the mean vector of the i -th factor analyzer, $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ (the factors are assumed to be distributed according to a zero-mean normal with unit variance), and $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \Psi_i)$ is a normal noise with diagonal covariance Ψ_i .

This diagonality is a key assumption in factor analysis. Namely, the observed variables are independent given the factors, and the goal of MFA is to best model the covariance structure of ξ .

Mixture of factor analyzers (MFA): graphical model



For MFA with covariance structure $\Sigma_i = \Lambda_i \Lambda_i^\top + \Psi$
(for $\Sigma_i = \Lambda_i \Lambda_i^\top + \Psi_i$, Ψ_i is moved to the right)

Mixture of factor analyzers (MFA) $\xi = \Lambda_i \mathbf{u} + \mu_i + \epsilon_i$

It follows from this model that the marginal distribution of ξ for the i -th component is

$$\xi \sim \mathcal{N}(\mu_i, \Lambda_i \Lambda_i^\top + \Psi_i)$$

and the joint distribution of ξ and \mathbf{u} is

$$\begin{bmatrix} \xi \\ \mathbf{u} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_i \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \Lambda_i \Lambda_i^\top + \Psi_i & \Lambda_i \\ \Lambda_i^\top & \mathbf{I} \end{bmatrix} \right)$$

To make some parallels with PCA, the above can be used to show that the d factors are informative projections of the data, which can be computed by Gaussian conditioning, corresponding to the affine projection

$$\mathbf{u} | \xi \sim \mathcal{N} \left(\mathbf{B}_i (\mu_i - \xi), \mathbf{I} - \mathbf{B}_i \Lambda_i \right) \quad \text{with} \quad \mathbf{B}_i = \Lambda_i^\top (\Lambda_i \Lambda_i^\top + \Psi_i)^{-1}$$

Mixture of factor analyzers (MFA) $\Sigma_i = \Lambda_i \Lambda_i^\top + \Psi_i$

This can be used to estimate the second moment of the factors

$$\begin{aligned}\mathbb{E}(\mathbf{u}\mathbf{u}^\top | \boldsymbol{\xi}) &= \text{cov}(\mathbf{u} | \boldsymbol{\xi}) + \mathbb{E}(\mathbf{u} | \boldsymbol{\xi}) \mathbb{E}(\mathbf{u} | \boldsymbol{\xi})^\top \\ &= \mathbf{I} - \mathbf{B}_i \Lambda_i + \mathbf{B}_i (\boldsymbol{\mu}_i - \boldsymbol{\xi}) (\boldsymbol{\mu}_i - \boldsymbol{\xi})^\top \mathbf{B}_i^\top\end{aligned}$$

which provides a measure of uncertainty in the factors that has no analogue in PCA.

This relation is exploited to derive an EM algorithm to train an MFA model of K components with parameters

$$\Theta^{\text{MFA}} = \{\pi_i, \boldsymbol{\mu}_i, \Lambda_i, \Psi_i\}_{i=1}^K$$

In the special case of a single cluster, it is worth noting that, in contrast to PPCA, FA also requires an EM algorithm to estimate

$$\Theta^{\text{FA}} = \{\boldsymbol{\mu}, \Lambda, \Psi\}$$

Estimation of parameters in MFA

In the case of MFA, it is considered that each datapoint ξ_t is associated with hidden variables z_t and u_t , and the goal is to maximize

$$\mathcal{L}(\Theta) = \sum_{t=1}^N \log \mathcal{P}(\xi_t | \Theta) = \sum_{t=1}^N \log \left(\sum_{z_t} \mathcal{P}(\xi_t, z_t, u_t | \Theta) \right)$$

which is, as seen before in the case of GMM, hard to optimize.

We can get around this problem by instead employing the expected complete data log-likelihood

$$Q(\Theta, \Theta^{\text{old}}) = \mathbb{E} \left[\sum_{t=1}^N \log \mathcal{P}(\xi_t, z_t, u_t | \Theta) \mid \xi, \Theta^{\text{old}} \right]$$

with $Q(\Theta, \Theta^{\text{old}})$ the auxiliary function.

Alternating Expectation Conditional Maximization (AECM)

In AECM, each iteration consists of the two cycles:

Cycle 1

Estimate μ_i and π_i with missing variables z_t based on auxiliary function $Q_1(\Theta, \Theta^{\text{old}})$.

Cycle 2

Estimate Λ_i and Ψ_i with missing variables z_t and u_t based on auxiliary function $Q_2(\Theta, \Theta^{\text{old}})$.

Each cycle has an E-step and a CM-step.

AECM guarantees convergence of the likelihood to the closest local optimum.

AEEM for MFA (UUU model in McNicholas and Murphy, 2008)

$$\Sigma_i = \Lambda_i \Lambda_i^\top + \Psi_i$$

The auxiliary function $Q_2(\Theta, \Theta^{\text{old}})$ to estimate Λ_i and Ψ_i becomes (see *McNicholas and Murphy (2008)* for details of computation)

$$Q_2(\Theta, \Theta^{\text{old}}) = \sum_{t=1}^N \sum_{i=1}^K h_{t,i} \left(\frac{1}{2} \log |\Psi_i^{-1}| - \text{tr}(\Psi_i^{-1} S_i) + \text{tr}(\Psi_i^{-1} \Lambda_i B_i S_i) - \frac{1}{2} \text{tr}(\Lambda_i^\top \Psi_i^{-1} \Lambda_i \theta_i) \right) + C$$

$$\mathbf{x}^\top \mathbf{S} \mathbf{x} = \text{tr}(\mathbf{S} \mathbf{x} \mathbf{x}^\top)$$

$$\text{with } S_i = \frac{\sum_{t=1}^N h_{t,i} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)(\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top}{\sum_{t=1}^N h_{t,i}}, \quad B_i = \Lambda_i^\top (\Lambda_i \Lambda_i^\top + \Psi_i)^{-1}$$

covariance as in GMM

$$\text{and } \boldsymbol{\theta}_i = \mathbf{I} - B_i \Lambda_i + B_i S_i B_i^\top$$

AEEM for MFA (UUU model in McNicholas and Murphy, 2008)

E-step:

$$\Sigma_i = \Lambda_i \Lambda_i^\top + \Psi_i$$

$$h_{t,i} = \frac{\pi_i \mathcal{N}(\xi_t \mid \mu_i, \Lambda_i \Lambda_i^\top + \Psi_i)}{\sum_{k=1}^K \pi_k \mathcal{N}(\xi_t \mid \mu_k, \Lambda_k \Lambda_k^\top + \Psi_k)}$$

CM-step:

Same as standard GMM

$$\pi_i \leftarrow \frac{\sum_{t=1}^N h_{t,i}}{N}$$

$$\Lambda_i \leftarrow \mathbf{S}_i \mathbf{B}_i^\top \overbrace{(\mathbf{I} - \mathbf{B}_i \Lambda_i + \mathbf{B}_i \mathbf{S}_i \mathbf{B}_i^\top)^{-1}}^{\theta_i^{-1}}$$

$$\mu_i \leftarrow \frac{\sum_{t=1}^N h_{t,i} \xi_t}{\sum_{t=1}^N h_{t,i}}$$

$$\Psi_i \leftarrow \text{diag} \{ \mathbf{S}_i - \Lambda_i \mathbf{B}_i \mathbf{S}_i \}$$

computed with the help of the intermediary variables

$$\mathbf{S}_i = \frac{\sum_{t=1}^N h_{t,i} (\xi_t - \mu_i) (\xi_t - \mu_i)^\top}{\sum_{t=1}^N h_{t,i}}$$

$$\mathbf{B}_i = \Lambda_i^\top (\Lambda_i \Lambda_i^\top + \Psi_i)^{-1}$$

covariance as in GMM

Mixture of probabilistic PCA (MPPCA)

$$\Sigma_i = \Lambda_i \Lambda_i^\top + \Psi_i$$

For comparison, the CM-step in MPPCA is given by

$$\tilde{\Lambda}_i \leftarrow \mathbf{S}_i \Lambda_i (\mathbf{I} \sigma_i^2 + \mathbf{M}_i^{-1} \Lambda_i^\top \mathbf{S}_i \Lambda_i)^{-1}$$

$$\Psi_i \leftarrow \mathbf{I} \sigma_i^2$$

computed with the help of the intermediary variables

covariance as in GMM

$$\mathbf{S}_i = \frac{\sum_{t=1}^N h_{t,i} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)(\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top}{\sum_{t=1}^N h_{t,i}}$$

$$\mathbf{M}_i = \Lambda_i^\top \Lambda_i + \mathbf{I} \sigma_i^2$$

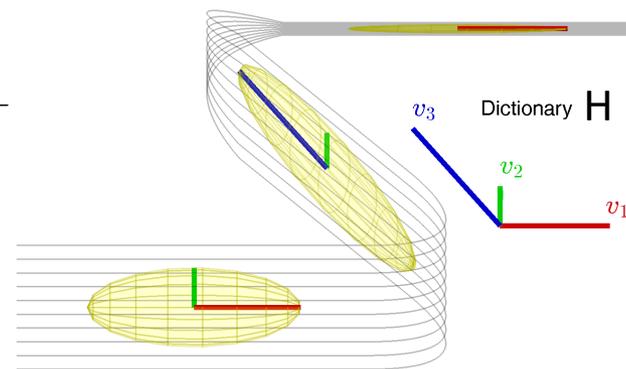
$$\sigma_i^2 = \frac{1}{D} \text{tr}(\mathbf{S}_i - \mathbf{S}_i \Lambda_i \mathbf{M}_i^{-1} \tilde{\Lambda}_i^\top)$$

where Λ_i is replaced by $\tilde{\Lambda}_i$ at each iteration.

GMM with semi-tied covariance matrices

The covariances share the same set of parameters for the latent feature space, where each covariance is composed of a common latent feature matrix $\mathbf{H} \in \mathbb{R}^{D \times D}$ and a component-specific diagonal covariance $\Sigma_i^{\text{diag}} \in \mathbb{R}^{D \times D}$ with

$$\Sigma_i = \mathbf{H} \Sigma_i^{\text{diag}} \mathbf{H}^\top$$



The latent feature matrix encodes the most relevant synergistic directions/basis vectors that are shared among all components, with the diagonal matrix representing the convex combination of basis vectors.

In other words, the aim is to find a global linear transformation of the data such that the transformed data can be modeled by a mixture of diagonal covariance matrices only.

GMM with semi-tied covariance matrices

$$\Sigma_i = \mathbf{H} \Sigma_i^{\text{diag}} \mathbf{H}^\top$$

The parameters of a GMM with semi-tied covariances are

$\Theta^{\text{tiedGMM}} = \{\mathbf{H}, \{\pi_i, \boldsymbol{\mu}_i, \Sigma_i^{\text{diag}}\}_{i=1}^K\}$. By setting $\mathbf{B} = \mathbf{H}^{-1}$, we have

$$\log |\mathbf{B}^{-1} \Sigma_i^{\text{diag}} \mathbf{B}^{-\top}| = \log \left(\frac{|\Sigma_i^{\text{diag}}|}{|\mathbf{B}|^2} \right) = \log |\Sigma_i^{\text{diag}}| - 2 \log |\mathbf{B}|$$

and the auxiliary function $Q(\Theta, \Theta^{\text{old}})$ of the standard GMM can be rewritten as

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{t=1}^N \sum_{i=1}^K h_{t,i} \left(\log(\pi_i) - \frac{1}{2} \log |\Sigma_i| - \frac{1}{2} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i) - \frac{D}{2} \log(2\pi) \right)$$

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{t=1}^N \sum_{i=1}^K h_{t,i} \left(\log(\pi_i) + \log |\mathbf{B}| - \frac{1}{2} \log |\Sigma_i^{\text{diag}}| - \frac{1}{2} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top \mathbf{B}^\top \Sigma_i^{(\text{diag})^{-1}} \mathbf{B} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i) - \frac{D}{2} \log(2\pi) \right).$$

GMM with semi-tied covariance matrices

$$\Sigma_i = \mathbf{B}^{-1} \Sigma_i^{\text{diag}} \mathbf{B}^{-\top}$$

Setting $\frac{\partial Q(\Theta, \Theta^{\text{old}})}{\partial \mathbf{B}}$ and $\frac{\partial Q(\Theta, \Theta^{\text{old}})}{\partial \Sigma_i^{\text{diag}}}$ equal to 0, and solving for \mathbf{B} and Σ_i^{diag} results in an expectation-maximization procedure to compute the maximum likelihood estimate of the parameters.

Following this, we get a row-by-row optimisation of \mathbf{B} , with \mathbf{b}_d (d -th row of \mathbf{B}) related to all other rows by the cofactor of \mathbf{B}

$$\begin{aligned} \mathbf{B}^{-1} &= \frac{\text{cof}(\mathbf{B})^\top}{|\mathbf{B}|} \\ \Leftrightarrow \begin{bmatrix} \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_D \end{bmatrix} &= |\mathbf{B}| (\mathbf{B}^\top)^{-1} \end{aligned} \quad \mathbf{b}_d = \mathbf{c}_d \mathbf{G}_d^{-1} \sqrt{\frac{\sum_{t=1}^T \sum_{i=1}^K h_{t,i}}{\mathbf{c}_d \mathbf{G}_d^{-1} \mathbf{c}_d^\top}}$$

where \mathbf{c}_d is the d -th row of cofactors of \mathbf{B} recomputed after each update of \mathbf{b}_d , and

$$\mathbf{G}_d = \sum_{i=1}^K \frac{1}{\Sigma_{i,d}^{\text{diag}}} \mathbf{S}_i \sum_{t=1}^T h_{t,i}$$

GMM with semi-tied covariance matrices

$$\Sigma_i = \mathbf{B}^{-1} \Sigma_i^{\text{diag}} \mathbf{B}^{-\top}$$

$\Sigma_{i,d}^{\text{diag}}$ is the d -th diagonal element of the i -th Gaussian, and \mathbf{S}_i is the full sample covariance matrix given by

covariance as in GMM

$$\mathbf{S}_i = \frac{\sum_{t=1}^T h_{t,i} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)(\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top}{\sum_{t=1}^T h_{t,i}}$$

The corresponding maximum likelihood estimate of Σ_i^{diag} is computed as

$$\Sigma_i^{\text{diag}} = \text{diag} \{ \mathbf{B} \mathbf{S}_i \mathbf{B}^\top \}$$

Note the variational nature of optimisation where the current estimate of Σ_i^{diag} is dependent on \mathbf{B} and vice versa.

Both \mathbf{B} and Σ_i^{diag} are iteratively improved in each EM step and the likelihood is guaranteed to increase at each step until convergence.