

EE613
Machine Learning for Engineers

SUBSPACE CLUSTERING

Sylvain Calinon
Robot Learning & Interaction Group
Idiap Research Institute
Oct. 25, 2017

SUBSPACE CLUSTERING (Wed, Oct. 25)

HIDDEN MARKOV MODELS (Wed, Nov. 1)

LINEAR REGRESSION (Thu, Nov. 9)

GAUSSIAN MIXTURE REGRESSION (Wed, Dec. 13)

GAUSSIAN PROCESS REGRESSION (Wed, Dec. 20)



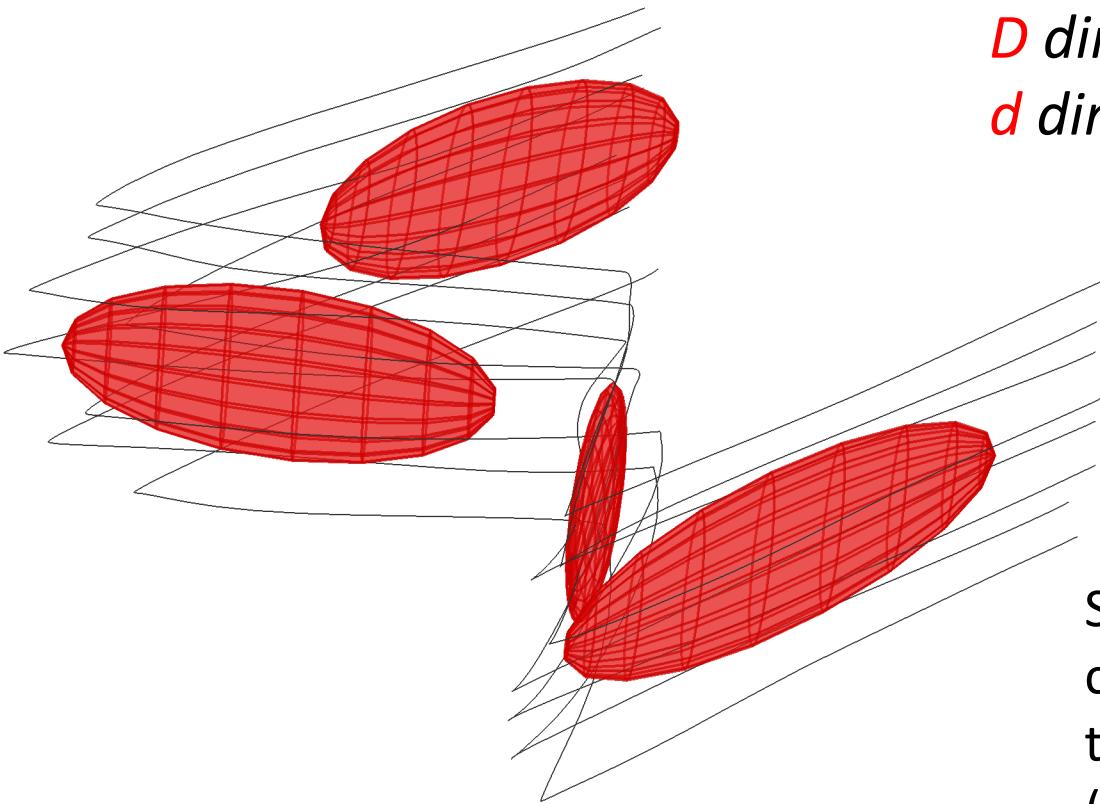
Time series analysis and synthesis,

Multivariate data processing

Outline

- High-dimensional data clustering (HDDC)
Matlab code: *demo_HDDC01.m*
- Mixture of factor analyzers (MFA)
Matlab code: *demo_MFA01.m*
- Mixture of probabilistic principal component analyzers (MPPCA)
Matlab code: *demo_MPPCA01.m*
- GMM with semi-tied covariance matrices
Matlab code: *demo_semitiedGMM01.m*

Introduction



K clusters

N datapoints

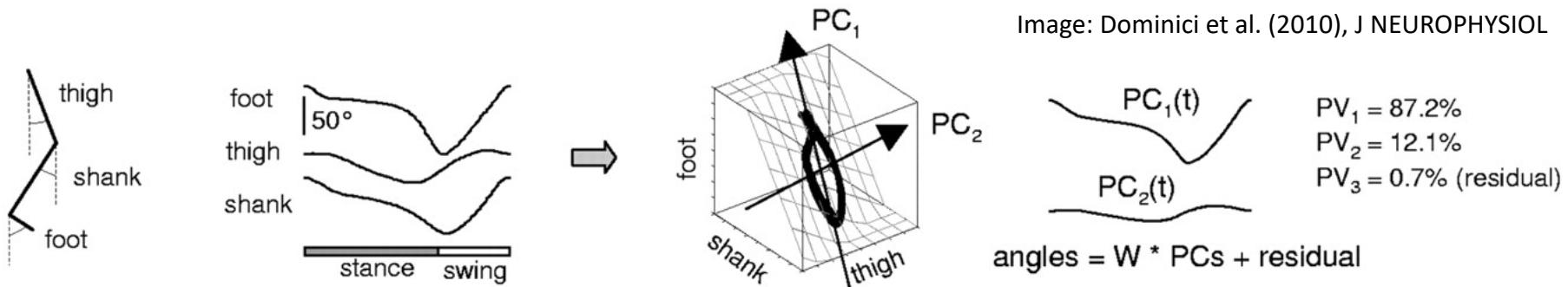
D dimensions (original space)

d dimensions (latent space)

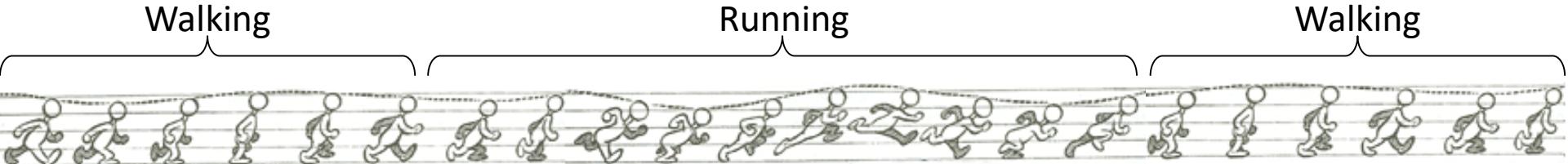
Subspace clustering aims at clustering data while reducing the dimension of each cluster (cluster-dependent subspace)

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

Curse of dimensionality in GMM encoding

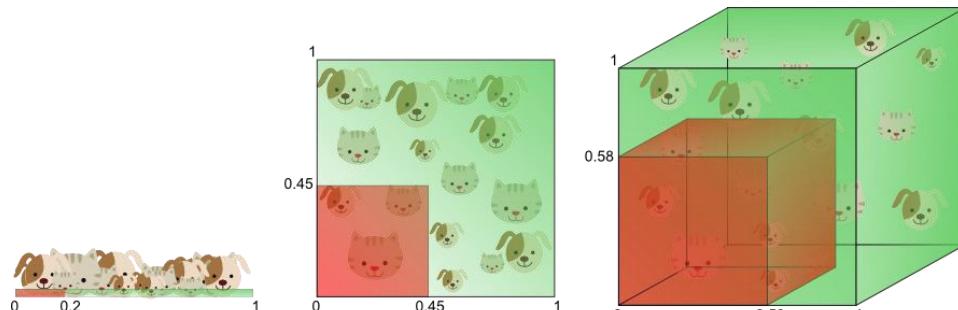


Image: datasciencecentral.com

K clusters
 N datapoints
 D dimensions (original space)
 d dimensions (latent space)

Classical Gaussian mixture models (GMM) tend to perform poorly in high-dimensional spaces if too few datapoints are available.

For a dataset $\{\boldsymbol{\xi}_t\}_{t=1}^N$ with $\boldsymbol{\xi}_t \in \mathbb{R}^D$, the *curse of dimensionality* appears if the dimension of the data D is too large compared to the size of the training set N .

In particular, the problem can affect the full covariance matrices $\boldsymbol{\Sigma}_i \in \mathbb{R}^{D \times D}$ because the number of parameters to be estimated grows quadratically with D .

Curse of dimensionality

Some characteristics of high-dimensional spaces can ease the classification of data. Indeed, having different groups living in different subspaces may be a useful property for discriminating the groups.

Subspace clustering exploits the phenomenon that high-dimensional spaces are mostly empty to ease the discrimination between groups of points.

→ **Curse of dimensionality or...
blessing of dimensionality?**

Curse of dimensionality

N datapoints

D dimensions (original space)

d dimensions (latent space)

Bouveyron and Brunet (2014, COMPUT STAT DATA AN) reviewed various ways of handling the problem of high-dimensional data in clustering problems:

1. Since D is too large w.r.t. N , a global dimensionality reduction should be applied as a pre-processing step to reduce D .
2. Since D is too large w.r.t. N , the solution space contains many poor local optima. The solution space should be smoothed by introducing ridge or lasso regularization in the estimation of the covariance (avoiding numerical problem and singular solutions when inverting the covariances). A simple form of regularization can be achieved after the maximization step of each EM loop.
3. Since D is too large w.r.t. N , the model is probably over-parametrized, and a more parsimonious model should be used (thus estimating a fewer number of parameters).

Gaussian Mixture Model (GMM)

K Gaussians
 N datapoints of dimension D

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

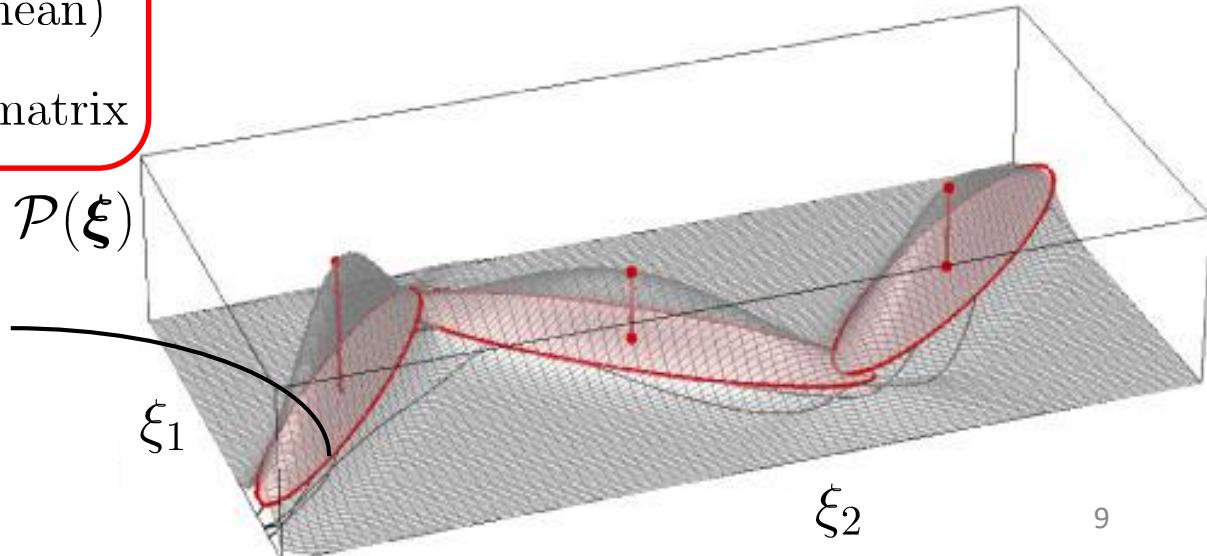
$$\mathcal{N}(\boldsymbol{\xi}_t | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) = \frac{1}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}_i|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}_i^{-1} (\boldsymbol{\xi}_t - \boldsymbol{\mu}_i) \right)$$

$\boldsymbol{\xi} \in \mathbb{R}^{D \times N}$ Observations ($N = \sum_{m=1}^M T_m$, the m -th trajectory has T_m datapoints)

$\pi_i \in \mathbb{R}$	Mixing coefficient
$\boldsymbol{\mu}_i \in \mathbb{R}^D$	Center (or mean)
$\boldsymbol{\Sigma}_i \in \mathbb{R}^{D \times D}$	Covariance matrix

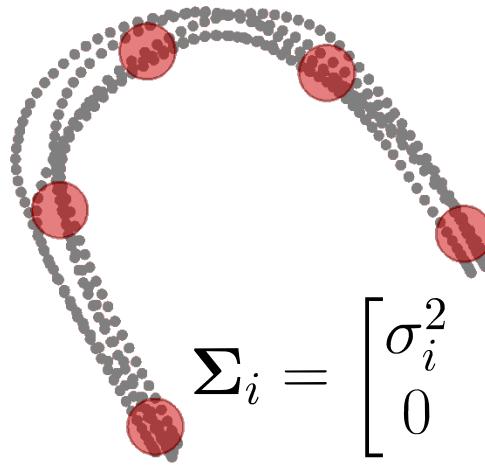
Parameters $\boldsymbol{\Theta}^{\text{GMM}} = \{\pi_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^K$

Equidensity contour of one standard deviation

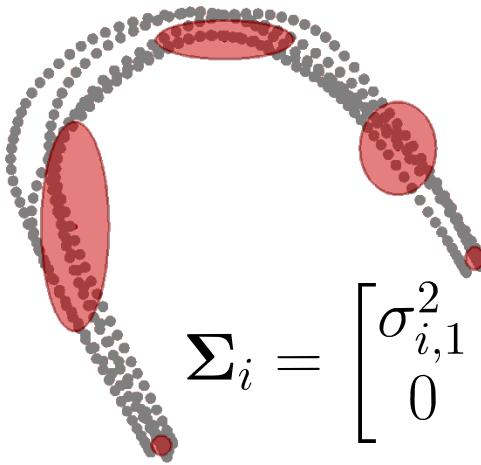


Covariance structures in GMM

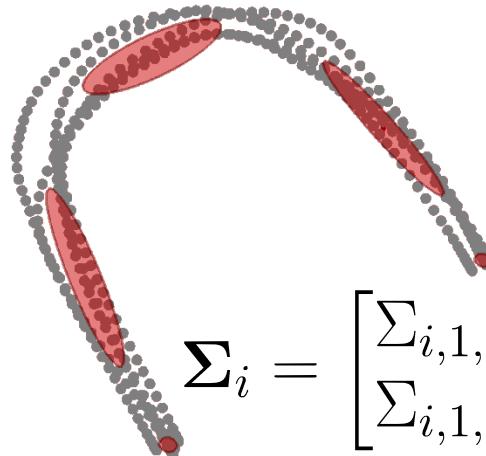
Isotropic



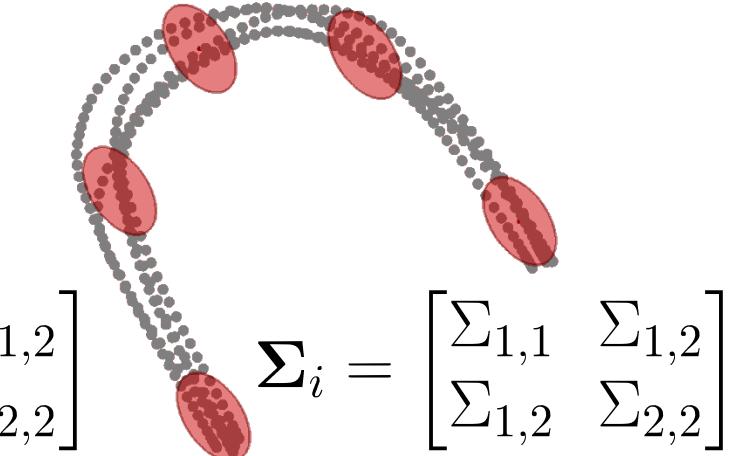
Diagonal



Full

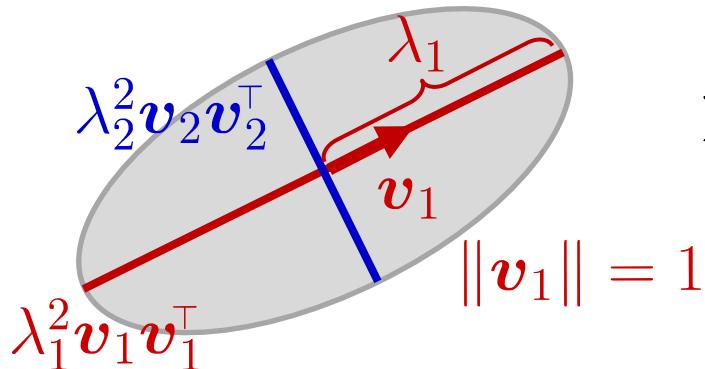


Tied



Multivariate normal distribution - Stochastic sampling

The eigendecomposition of Σ is expressed in a matrix form as



$$\Sigma = V D V^\top = \sum_{j=1}^D \lambda_j^2 \mathbf{v}_j \mathbf{v}_j^\top$$

with $V = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_D]$

$$D = \begin{bmatrix} \lambda_1^2 & 0 & \cdots & 0 \\ 0 & \lambda_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_D^2 \end{bmatrix}$$

By using this notation, datapoints can be stochastically generated with



$$\xi \sim \mathcal{N}(\mu, \Sigma) \iff \xi \sim \mu + V D^{\frac{1}{2}} \mathcal{N}(0, I)$$

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 \mathbf{z}_t .
The goal is to maximize the log-likelihood of the observed data

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

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

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

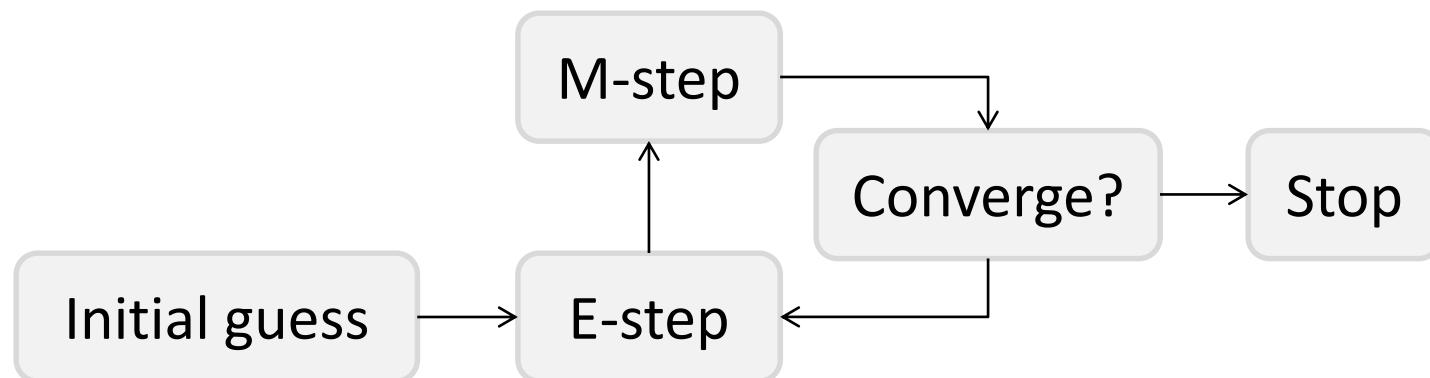
where $Q(\Theta, \Theta^{\text{old}})$ is called the auxiliary function.

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

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

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

$z_{t,i} = 1$ if ξ_i is part of cluster i .
 It is 0 otherwise.
e.g. $\prod_{i=1}^3 \pi_i^{z_i} = \pi_1^{z_1} \cdot \pi_2^{z_2} \cdot \pi_3^{z_3}$
 $= \pi_1^0 \cdot \pi_2^0 \cdot \pi_3^1$
 $= 1 \cdot 1 \cdot \pi_3$

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

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

EM for GMM

Setting

$$\frac{\partial \mathcal{Q}(\Theta, \Theta^{\text{old}})}{\partial \pi_i} = 0 \quad \frac{\partial \mathcal{Q}(\Theta, \Theta^{\text{old}})}{\partial \boldsymbol{\mu}_i} = 0 \quad \frac{\partial \mathcal{Q}(\Theta, \Theta^{\text{old}})}{\partial \Sigma_i} = 0$$

and solving for π_i , $\boldsymbol{\mu}_i$ and Σ_i results in an EM procedure to compute the maximum likelihood estimate of the parameters.

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

By using the linear algebra relations

($= 2Ax$ if A symmetric)

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

and the derivation chain rule, we obtain

$$\frac{\partial \mathcal{Q}(\Theta, \Theta^{\text{old}})}{\partial \mu_i} = \frac{1}{2} \sum_{t=1}^N h_{t,i} 2\Sigma_i^{-1} (\xi_t - \mu_i) = \Sigma_i^{-1} \sum_{t=1}^N h_{t,i} (\xi_t - \mu_i) = 0$$

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

$$\frac{\partial \mathcal{Q}(\Theta, \Theta^{\text{old}})}{\partial \Sigma_i} = \frac{1}{2} \Sigma_i \sum_{t=1}^N h_{t,i} - \frac{1}{2} \sum_{t=1}^N h_{t,i} (\xi_t - \mu_i)(\xi_t - \mu_i)^\top = 0$$

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

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

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

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

The sum over K of the above relation provides

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

which can be reintroduced in the equation to find

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

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:

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

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

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

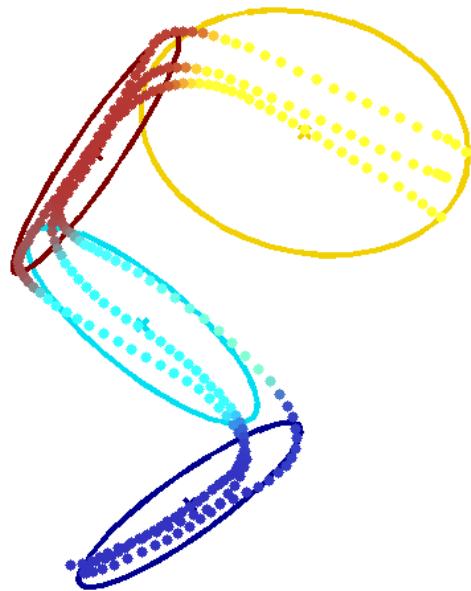
These results can be intuitively interpreted in terms of normalized counts.

EM provides a systematic approach to derive such procedure.

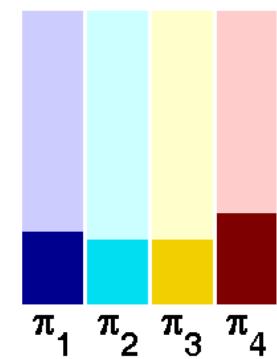
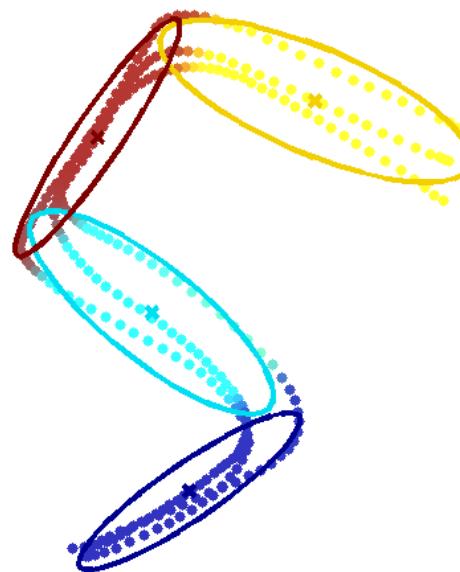
→ **Weighted averages taking into account the responsibility of each datapoint in each cluster.**

EM for GMM

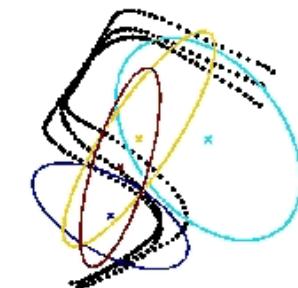
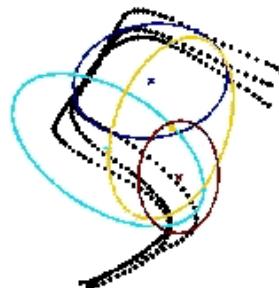
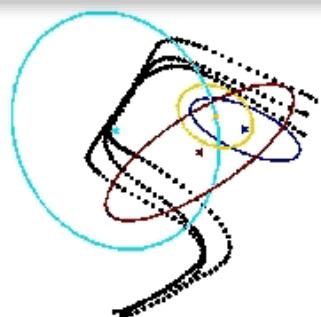
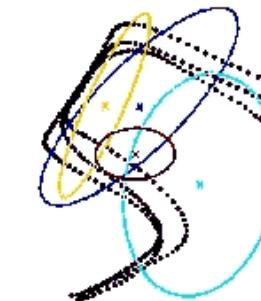
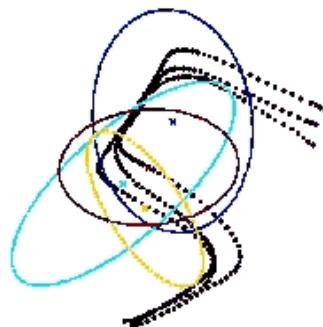
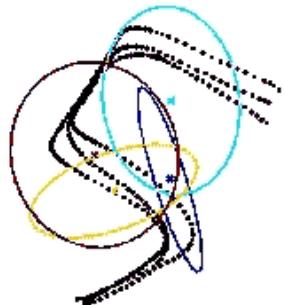
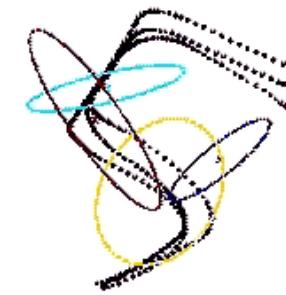
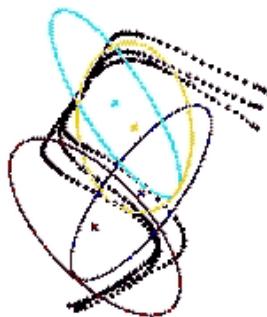
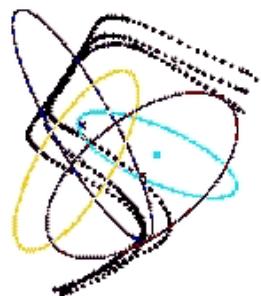
E-step



M-step

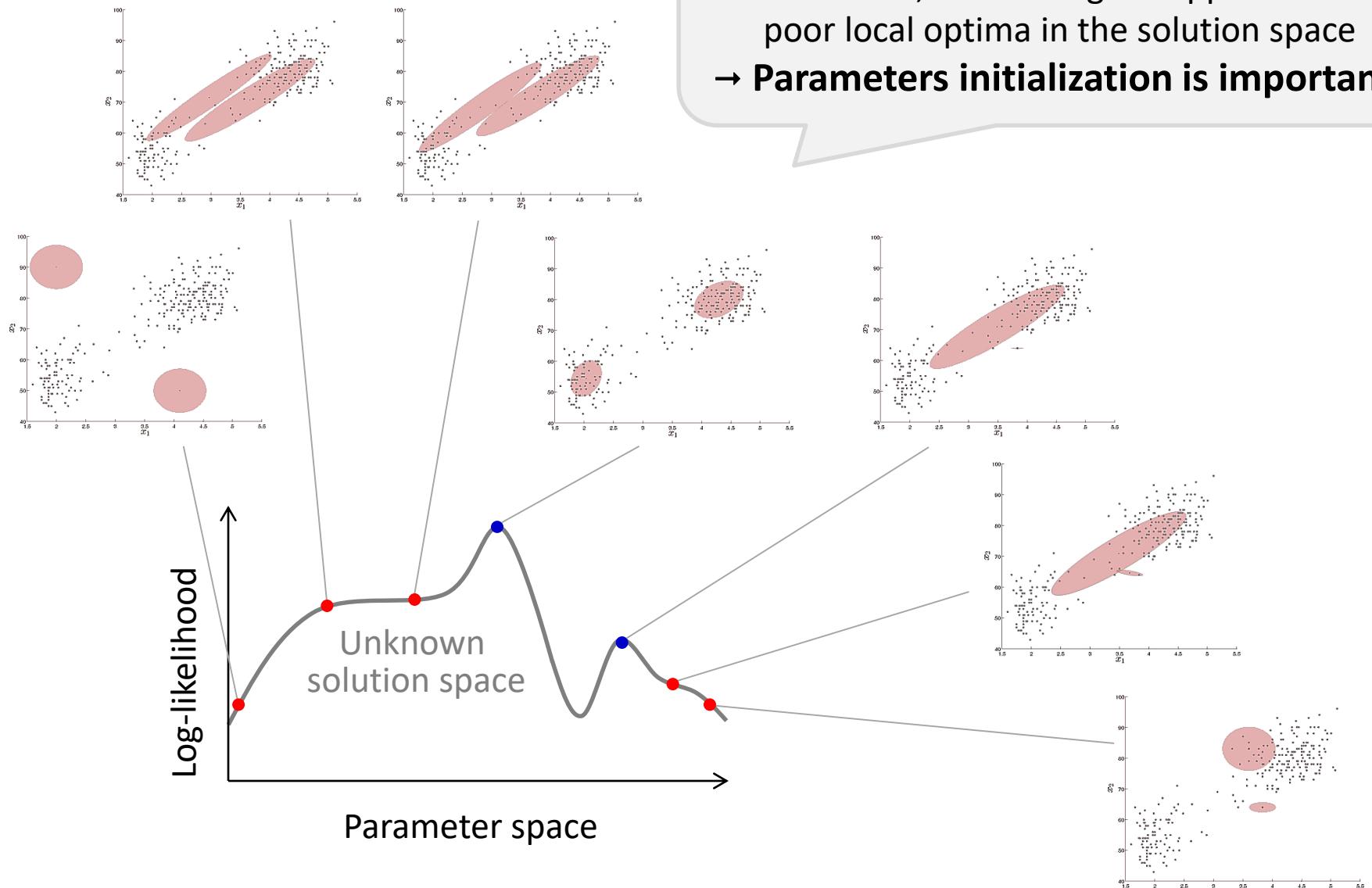


EM for GMM: Local optima issue



20

Local optima in EM

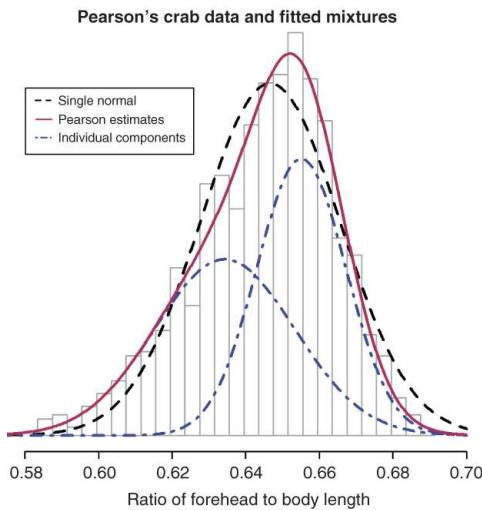


EM will improve the likelihood at each iteration, but it can get trapped into poor local optima in the solution space

→ Parameters initialization is important!

Parameters estimation in GMM... in 1893

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



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.

...—On the Dissection of Asymmetrical Frequency-Curves. General Theory, §§ 1–8. Example: Professor WELDON's measurements of the "Forehead" of Crabs. §§ 9–10	Page. 71–85
II.—On the Dissection of Symmetrical Frequency-Curves. General Theory, §§ 11–12 Application. Crabs "No. 1" &c. 19–15	85–90
	90–100

00–106

106

107

07–110

... which does not mean that moment-based approaches are old-fashioned!

They are actually today popular again with new developments related to spectral decomposition.

High-dimensional data clustering (HDDC)

Matlab code: demo_HDDC01.m

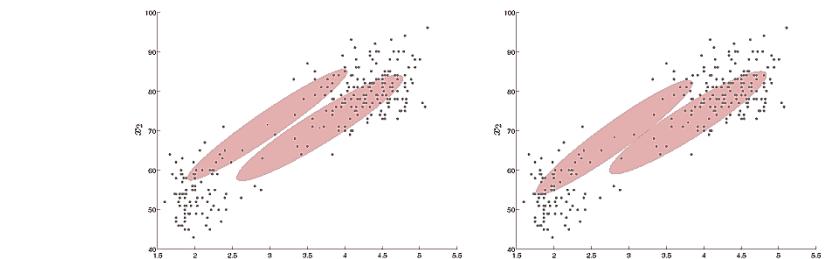
[C. Bouveyron and C. Brunet. Model-based clustering of high-dimensional data: A review. Computational Statistics and Data Analysis, 71:52–78, March 2014]

Curse of dimensionality

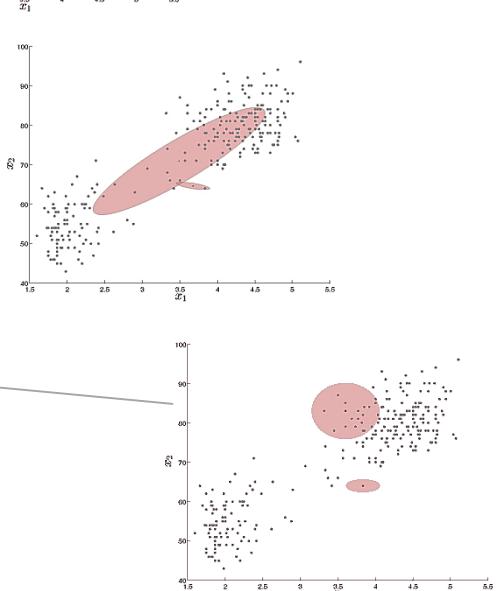
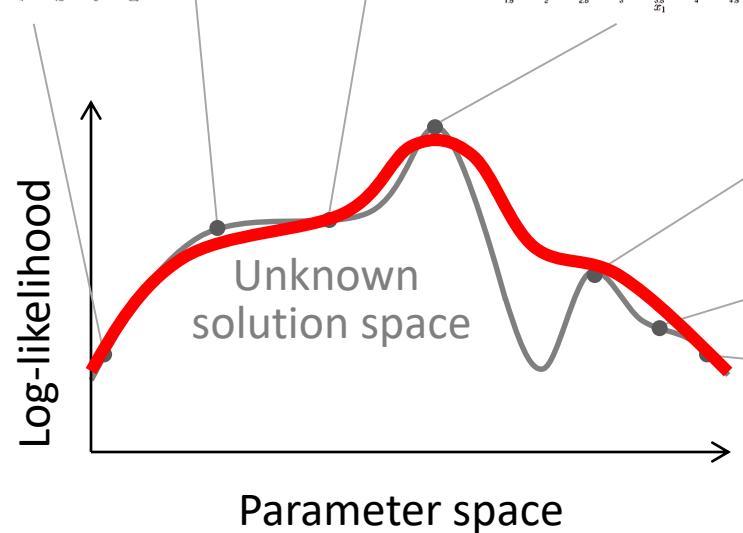
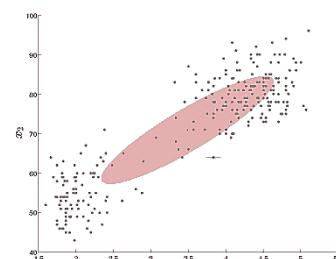
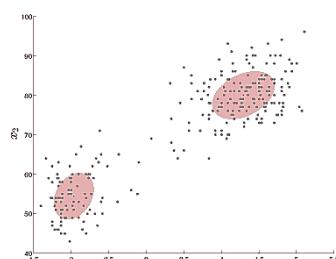
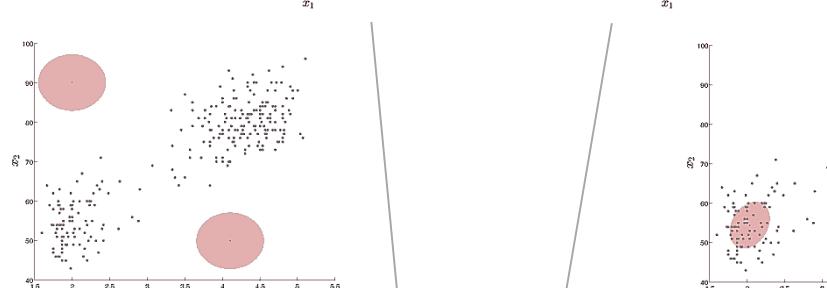
Bouveyron and Brunet (2014, COMPUT STAT DATA AN) reviewed various ways of viewing the problem and coping with high-dimensional data in clustering problems:

1. Since D is too large wrt N , a global dimensionality reduction should be applied as a pre-processing step to reduce D .
2. Since D is too large wrt N , the solution space contains many poor local optima; the solution space should be smoothed by introducing ridge or lasso regularization in the estimation of the covariance (avoiding numerical problem and singular solutions when inverting the covariances). A simple form of regularization can be achieved after the maximization step of each EM loop.
3. Since D is too large wrt N , the model is probably over-parametrized, and a more parsimonious model should be used (thus estimating a fewer number of parameters).

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:

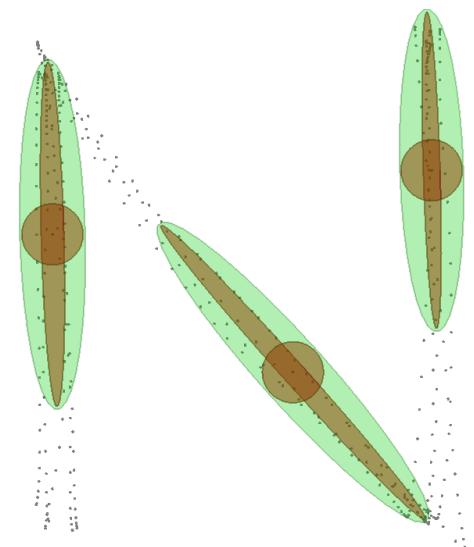
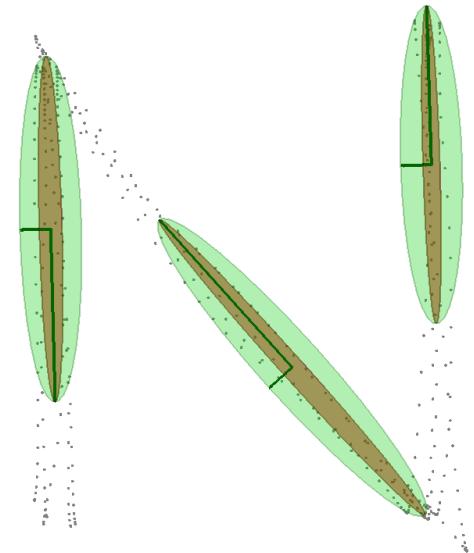
$$\Sigma_i \leftarrow V_i \tilde{D}_i V_i^\top$$

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

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

Tikhonov regularization with diagonal isotropic covariance:

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



High-dimensional data clustering (HDDC)

The HDDC approach (Bouveyron, 2007, COMPUT STAT DATA AN) addresses both subspace clustering and regularization.

One implementation considers that the subspace of each cluster i is generated by the first d_i eigenvectors associated with the first d_i eigenvalues $\lambda_{i,k}$, and that outside of this subspace, the variance is isotropic, modeled by

$$\bar{\lambda}_i = \frac{1}{D-d_i} \sum_{k=d_i+1}^D \lambda_{i,k} = \frac{1}{D-d_i} \left(\text{tr}(\Sigma_i) - \sum_{k=1}^{d_i} \lambda_{i,k} \right)$$

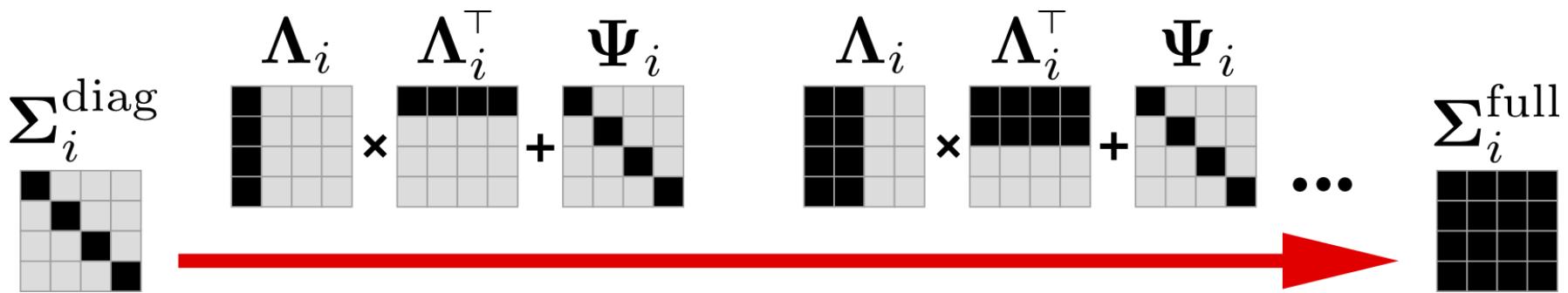
which is used to reconstruct a full covariance matrix by replacing the last $D - d_i$ eigenvalues with $\bar{\lambda}_i$.

Mixture of factor analyzers (MFA)

Matlab code: demo_MFA01.m

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

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 $\Psi_i = I\sigma_i^2$.

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

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

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



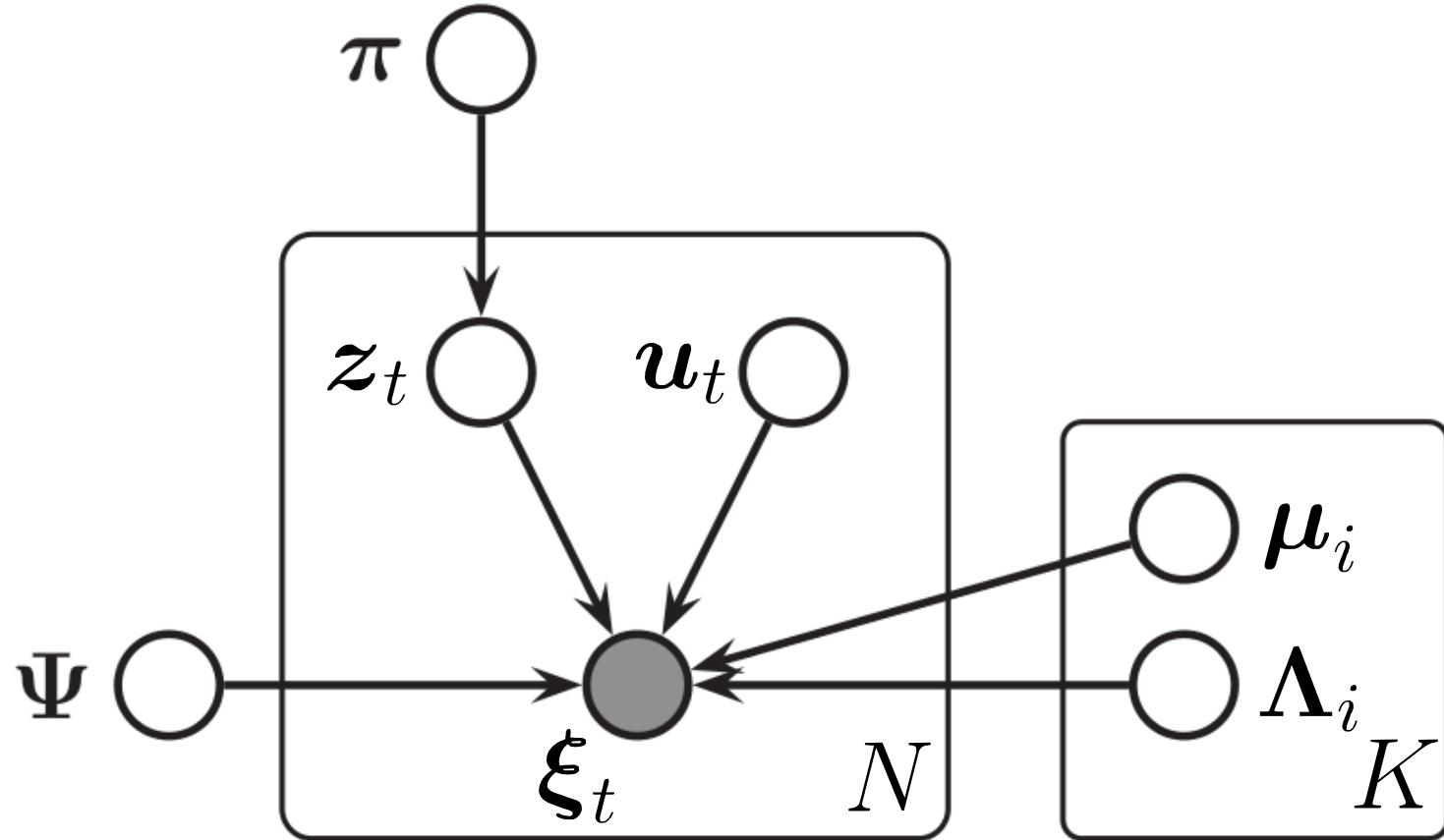
In MFA, the generative model for the i -th mixture component assumes that a D -dimensional random vector ξ is modeled using a d -dimensional vector of latent (unobserved) factors \mathbf{u}

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

where $\mu_i \in \mathbb{R}^D$ is the mean vector of the i -th factor analyzer, $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, I)$ (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) $\xi = \Lambda_i u + \mu_i + \epsilon_i$

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 u is

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

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

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

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

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

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

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

Estimation of parameters in MFA

In the case of MFA, it is considered that each datapoint ξ_t is associated with hidden variables \mathbf{z}_t and \mathbf{u}_t , and the goal is to maximize

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

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

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

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

Alternating Expectation Conditional Maximization (AECM)

In AECM, each iteration consists of the two cycles:

Cycle 1

Estimate $\boldsymbol{\mu}_i$ and π_i with missing variables \mathbf{z}_t based on auxiliary function $Q_1(\boldsymbol{\Theta}, \boldsymbol{\Theta}^{\text{old}})$.

Cycle 2

Estimate $\boldsymbol{\Lambda}_i$ and $\boldsymbol{\Psi}_i$ with missing variables \mathbf{z}_t and \mathbf{u}_t based on auxiliary function $Q_2(\boldsymbol{\Theta}, \boldsymbol{\Theta}^{\text{old}})$.

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

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

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

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

The auxiliary function $\mathcal{Q}_2(\Theta, \Theta^{\text{old}})$ to estimate Λ_i and Ψ_i becomes (see *McNicholas and Murphy (2008)* for details of computation)

$$\begin{aligned} \mathcal{Q}_2(\Theta, \Theta^{\text{old}}) = & \sum_{t=1}^N \sum_{i=1}^K h_{t,i} \left(\frac{1}{2} \log |\Psi_i^{-1}| - \text{tr}(\Psi_i^{-1} S_i) \right. \\ & \left. + \text{tr}(\Psi_i^{-1} \Lambda_i B_i S_i) - \frac{1}{2} \text{tr}(\Lambda_i^\top \Psi_i^{-1} \Lambda_i \theta_i) \right) + C \end{aligned}$$

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

$$\text{with } S_i = \frac{\sum_{t=1}^N h_{t,i} (\xi_t - \mu_i)(\xi_t - \mu_i)^\top}{\sum_{t=1}^N h_{t,i}}, \quad B_i = \Lambda_i^\top (\Lambda_i \Lambda_i^\top + \Psi_i)^{-1}$$

covariance as in GMM

$$\text{and } \theta_i = I - B_i \Lambda_i + B_i S_i B_i^\top$$

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

E-step:

$$\boldsymbol{\Sigma}_i = \boldsymbol{\Lambda}_i \boldsymbol{\Lambda}_i^\top + \boldsymbol{\Psi}_i$$

$$h_{t,i} = \frac{\pi_i \mathcal{N}(\boldsymbol{\xi}_t \mid \boldsymbol{\mu}_i, \boldsymbol{\Lambda}_i \boldsymbol{\Lambda}_i^\top + \boldsymbol{\Psi}_i)}{\sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{\xi}_t \mid \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k \boldsymbol{\Lambda}_k^\top + \boldsymbol{\Psi}_k)}$$

CM-step:

Same as standard GMM

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

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

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

$$\boldsymbol{\Psi}_i \leftarrow \text{diag}\{\boldsymbol{S}_i - \boldsymbol{\Lambda}_i \boldsymbol{B}_i \boldsymbol{S}_i\}$$

computed with the help of the intermediary variables

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

$$\boldsymbol{B}_i = \boldsymbol{\Lambda}_i^\top (\boldsymbol{\Lambda}_i \boldsymbol{\Lambda}_i^\top + \boldsymbol{\Psi}_i)^{-1}$$

covariance as in GMM

Mixture of probabilistic PCA (MPPCA)

Matlab code: [demo_MPPCA01.m](#)

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

Mixture of probabilistic PCA (MPPCA)

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

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

$$\tilde{\Lambda}_i \leftarrow S_i \Lambda_i (I \sigma_i^2 + M_i^{-1} \Lambda_i^\top S_i \Lambda_i)^{-1}$$

$$\Psi_i \leftarrow I \sigma_i^2$$

computed with the help of the intermediary variables

covariance as in GMM

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

$$M_i = \Lambda_i^\top \Lambda_i + I \sigma_i^2$$

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

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

A taxonomy of parsimonious GMMs

D in the slides of this lecture

Model name	Cov. structure	Nb. of parameters	$K = 4, d = 3$
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

$p = 100$

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

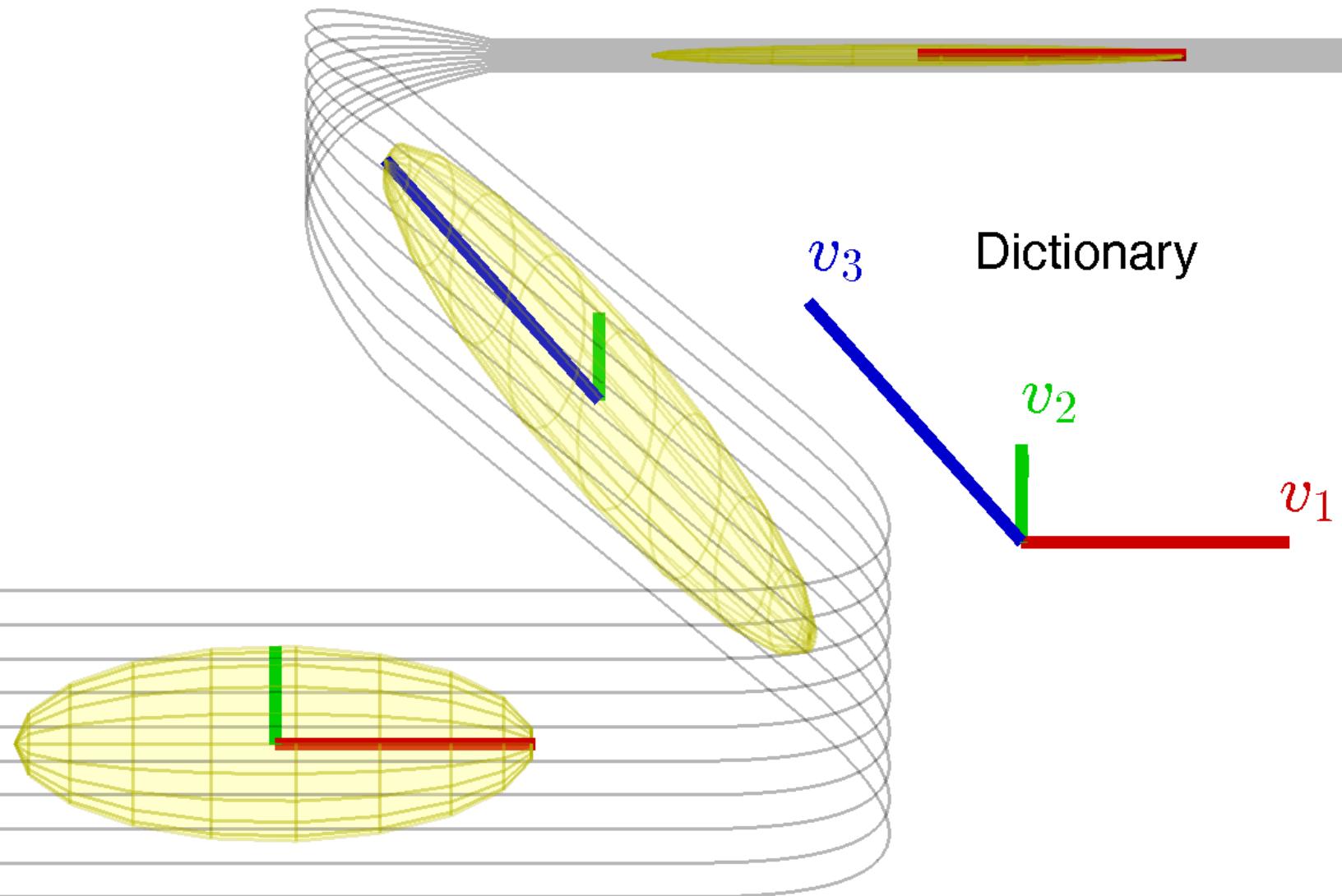
[C. Bouveyron and C. Brunet. Model-based clustering of high-dimensional data: A review. Computational Statistics and Data Analysis, 71:52–78, March 2014]

GMM with semi-tied covariance matrices

Matlab code: demo_semitiedGMM01.m

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

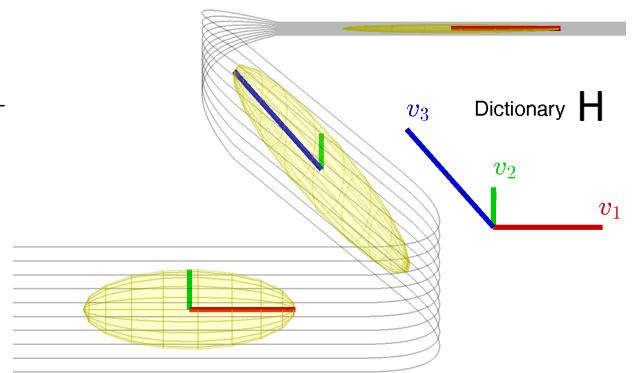
Sharing of parameters in mixture models



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

$$\boldsymbol{\Sigma}_i = \mathbf{H} \boldsymbol{\Sigma}_i^{\text{diag}} \mathbf{H}^\top$$

The parameters of a GMM with semi-tied covariances are $\boldsymbol{\Theta}^{\text{tiedGMM}} = \{\mathbf{H}, \{\pi_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i^{\text{diag}}\}_{i=1}^K\}$. By setting $\mathbf{B} = \mathbf{H}^{-1}$, we have

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

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

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

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

GMM with semi-tied covariance matrices

$$\boldsymbol{\Sigma}_i = \mathbf{B}^{-1} \boldsymbol{\Sigma}_i^{\text{diag}} \mathbf{B}^{-\top}$$

Setting $\frac{\partial \mathcal{Q}(\boldsymbol{\Theta}, \boldsymbol{\Theta}^{\text{old}})}{\partial \mathbf{B}}$ and $\frac{\partial \mathcal{Q}(\boldsymbol{\Theta}, \boldsymbol{\Theta}^{\text{old}})}{\partial \boldsymbol{\Sigma}_i^{\text{diag}}}$ equal to 0, and solving for \mathbf{B} and $\boldsymbol{\Sigma}_i^{\text{diag}}$ 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}

$$\begin{aligned} \mathbf{B}^{-1} &= \frac{\text{cof}(\mathbf{B})^\top}{|\mathbf{B}|} \\ \iff \left[\begin{array}{c} \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_D \end{array} \right] &= |\mathbf{B}| (\mathbf{B}^\top)^{-1} \end{aligned}$$

$$\mathbf{b}_d = \mathbf{c}_d \mathbf{G}_d^{-1} \sqrt{\frac{\sum_{t=1}^T \sum_{i=1}^K h_{t,i}}{\mathbf{c}_d \mathbf{G}_d^{-1} \mathbf{c}_d^\top}}$$

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

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

GMM with semi-tied covariance matrices

$$\boldsymbol{\Sigma}_i = \mathbf{B}^{-1} \boldsymbol{\Sigma}_i^{\text{diag}} \mathbf{B}^{-\top}$$

$\Sigma_{i,d}^{\text{diag}}$ 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 $\boldsymbol{\Sigma}_i^{\text{diag}}$ is computed as

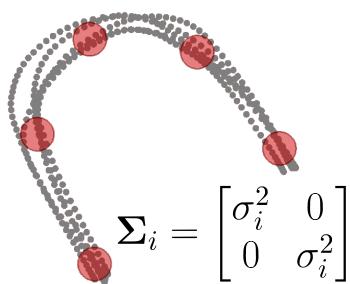
$$\boldsymbol{\Sigma}_i^{\text{diag}} = \text{diag} \{ \mathbf{B} \mathbf{S}_i \mathbf{B}^\top \}$$

Note the variational nature of optimisation where the current estimate of $\boldsymbol{\Sigma}_i^{\text{diag}}$ is dependent on \mathbf{B} and vice versa.

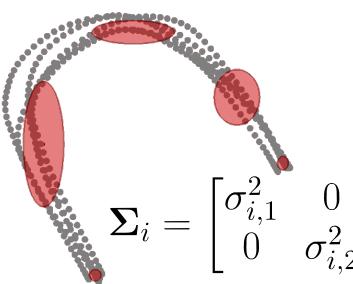
Both \mathbf{B} and $\boldsymbol{\Sigma}_i^{\text{diag}}$ are iteratively improved in each EM step and the likelihood is guaranteed to increase at each step until convergence.

Summary of relevant covariance structures

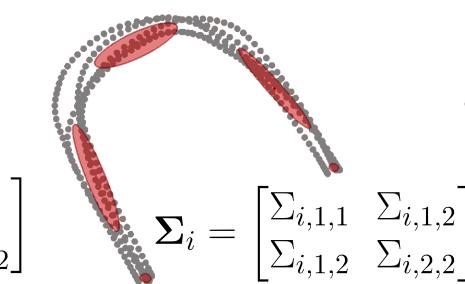
Isotropic



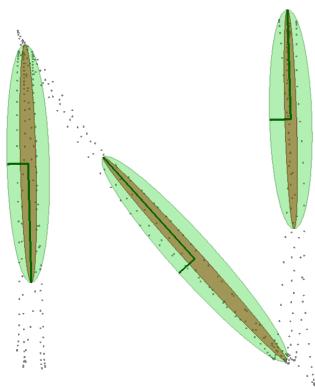
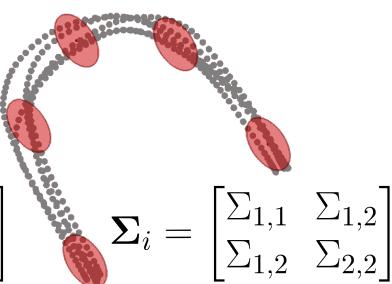
Diagonal



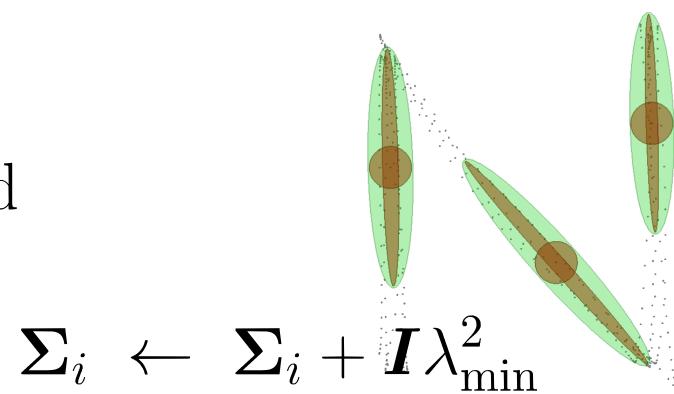
Full



Tied



$\Sigma_i \leftarrow \mathbf{V}_i \tilde{\mathbf{D}}_i \mathbf{V}_i^\top$ with
 $\tilde{\mathbf{D}}_i = \text{diag}(\tilde{\lambda}_{i,1}^2, \dots, \tilde{\lambda}_{i,D}^2)$ and
 $\tilde{\lambda}_{i,j}^2 = \max(\lambda_{i,j}^2, \lambda_{\min}^2)$



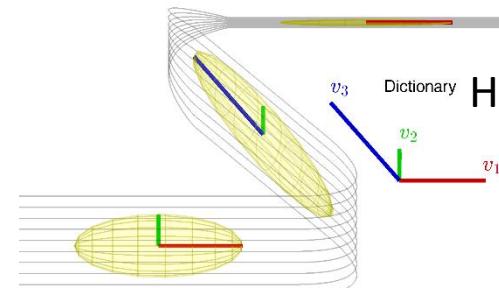
MFA: $\Sigma_i = \Lambda_i \Lambda_i^\top + \Psi_i$

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

Σ_i^{diag} Λ_i Λ_i^\top Ψ_i Λ_i Λ_i^\top Ψ_i ... Σ_i^{full}

$\Sigma_i^{\text{diag}} \xrightarrow{\quad} \Lambda_i \times \Lambda_i^\top + \Psi_i \xrightarrow{\quad} \Sigma_i^{\text{full}}$

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



Main 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

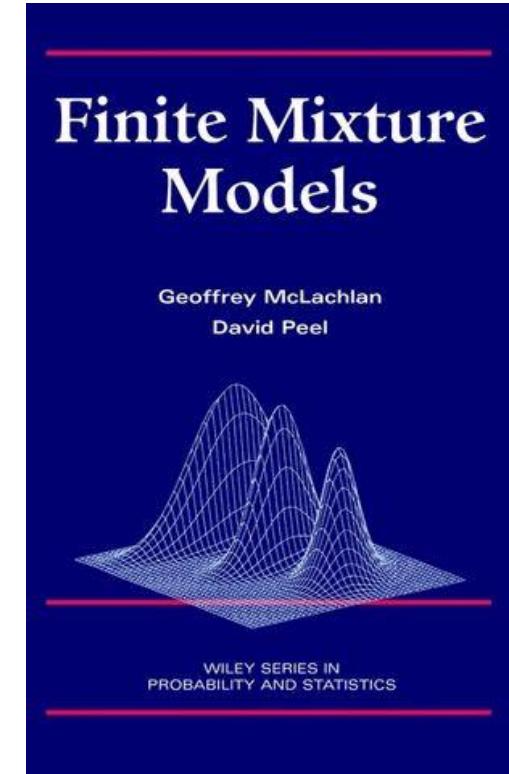
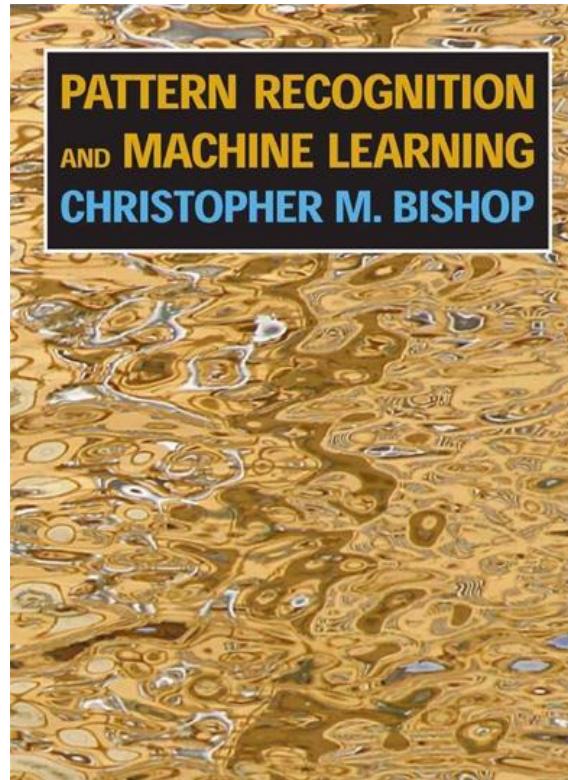
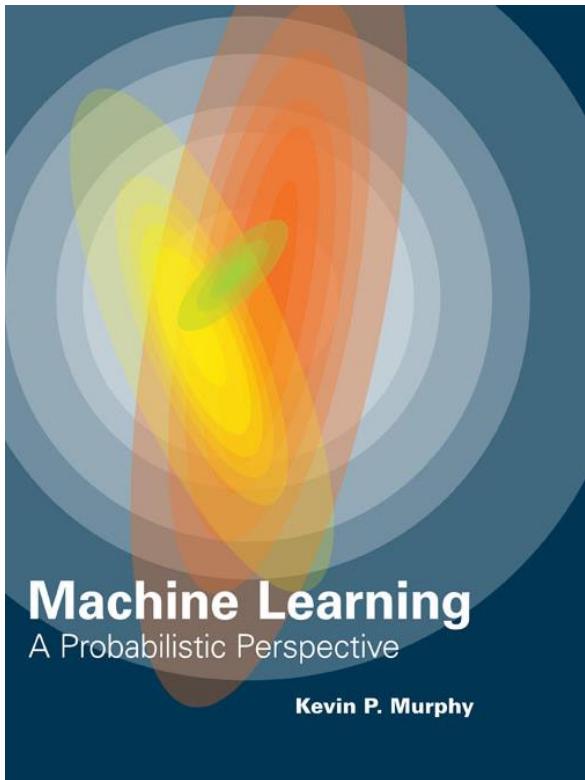
MPPCA

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

General textbooks



Advanced related research topics

(not covered in the course)

Coordinated MFA by using common factor loadings

J. Baek, G. J. McLachlan, and L. K. Flack. Mixtures of factor analyzers with common factor loadings: Applications to the clustering and visualization of high-dimensional data. *IEEE Trans. Pattern Anal. Mach. Intell.*, 32(7):1298–1309, 2010

J. Verbeek. Learning nonlinear image manifolds by global alignment of local linear models. *IEEE Trans. on Pattern Analysis & Machine Intelligence*, 28(8):1236–1250, August 2006

Estimation of K and d_k in MFA with Bayesian nonparametrics

Y. Wang and J. Zhu. DP-space: Bayesian nonparametric subspace clustering with small-variance asymptotics. In Proc. Intl Conf. on Machine Learning (ICML), pages 1–9, Lille, France, 2015

Online parameters estimation in MPPCA

A. Bellas, C. Bouveyron, M. Cottrell, and J. Lacaille. Model-based clustering of high-dimensional data streams with online mixture of probabilistic PCA. *Advances in Data Analysis and Classification*, 7(3):281–300, 2013

Advanced related research topics

(not covered in the course)

Sparse subspace clustering with L₁ regularization

H. Zou, T. Hastie, and R. Tibshirani. Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, 15(2):265–286, 2006

Y. Guan and J. G. Dy. Sparse probabilistic principal component analysis. In *Intl Conf. on Artificial Intelligence and Statistics*, pages 185–192, 2009

Deep MFA

Y. Tang, R. Salakhutdinov, and G. Hinton. Deep mixtures of factor analysers. In *Proc. Intl Conf. on Machine Learning (ICML)*, Edinburgh, Scotland, 2012

Mixture of tensor analyzers (MTA)

Y. Tang, R. Salakhutdinov, and G. Hinton. Tensor analyzers. In *Proc. Intl Conf. on Machine Learning (ICML)*, Atlanta, USA, 2013