

	<p style="text-align: center;">PCA and Mixtures of PCA: Improving the robustness to outliers</p> <p style="text-align: center;">Joint work with - Nicolas Delannay, UCL machine learning group - Cédric Archambeau, University College London</p> <p style="text-align: right;">26 October 2007</p>
M. Verleysen UCL 1	

	<p>Overview</p>
	<ul style="list-style-type: none">■ Principal Component Analysis: a reminder■ Probabilistic PCA■ Robust probabilistic PCA■ Mixtures of (robust) probabilistic PCA■ Experiments
M. Verleysen UCL 2	

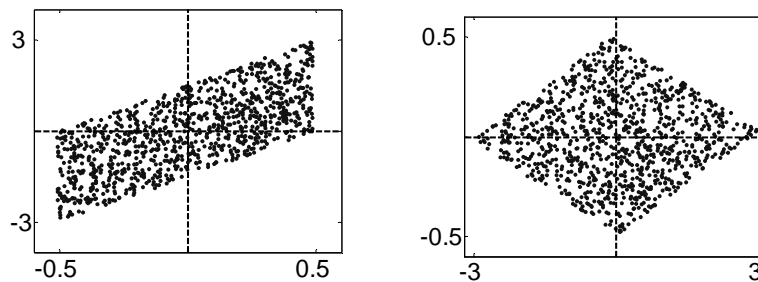
Overview

- Principal Component Analysis: a reminder
- Probabilistic PCA
- Robust probabilistic PCA
- Mixtures of (robust) probabilistic PCA
- Experiments

M. Verleysen
UCL
3

Maximum variance direction [1/2]

- We look for an orthogonal coordinate system such that the elements of \mathbf{X} in the new system are uncorrelated
- =
- We look for axes that maximize variance after projection



M. Verleysen
UCL
4

Maximum variance direction [2/2]

- We look for axes which maximise the variance after projection

i.e. along axis \mathbf{u} with $\|\mathbf{u}\| = \sqrt{\mathbf{u}\mathbf{u}^T} = 1$

- In the new coordinate system : $x_1 = \mathbf{u}\mathbf{Y}$

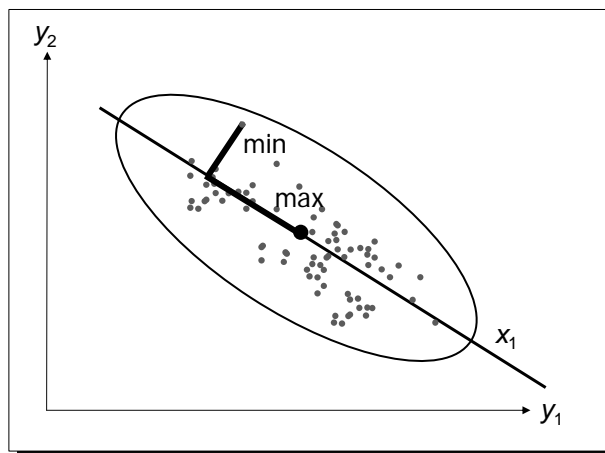
- Therefore :

$$\text{Var}(X_1) = \sigma_{X_1}^2 = E[X_1 X_1] = \mathbf{u}^T E[\mathbf{Y}\mathbf{Y}^T] \mathbf{u}$$

- Maximum variance = min. distortion (LMS)

M. Verleysen
UCL
5

Maximum variance = minimum distortion



M. Verleysen
UCL
6

Choice of direction [1/3]

- Choice of \mathbf{u} ?

$$\begin{aligned} \mathbf{e} &= \underset{\mathbf{u}}{\operatorname{argmax}} \mathbf{u}^T E[\mathbf{Y}\mathbf{Y}^T] \mathbf{u} \\ &= \underset{\mathbf{u}}{\operatorname{argmax}} \mathbf{u}^T \mathbf{C}_{\mathbf{Y}\mathbf{Y}} \mathbf{u} \\ &= \underset{\mathbf{u}}{\operatorname{argmax}} \mathbf{u}^T \underbrace{\boldsymbol{\Theta}\boldsymbol{\Lambda}\boldsymbol{\Theta}^T}_{\operatorname{EVD}(\mathbf{C}_{\mathbf{Y}\mathbf{Y}})} \mathbf{u} \end{aligned}$$

$$\sigma_{X_1, \max}^2 = \mathbf{e}^T \underbrace{\boldsymbol{\Theta}\boldsymbol{\Lambda}\boldsymbol{\Theta}^T}_{\operatorname{EVD}(\mathbf{C}_{\mathbf{Y}\mathbf{Y}})} \mathbf{e}$$

M. Verleysen
UCL
7

Choice of direction [2/3]

- Choice of \mathbf{u} ?

$$\mathbf{e} = \underset{\mathbf{u}}{\operatorname{argmax}} \mathbf{u}^T E[\mathbf{Y}\mathbf{Y}^T] \mathbf{u}$$

$$\mathbf{e} = \boldsymbol{\Theta}_1$$

where $\boldsymbol{\Theta}_1$ is the eigenvector corresponding to the largest eigenvalue of $\mathbf{C}_{\mathbf{Y}\mathbf{Y}}$

$$\sigma_{X_1, \max}^2 = \boldsymbol{\Theta}_1^T \underbrace{\boldsymbol{\Theta}\boldsymbol{\Lambda}\boldsymbol{\Theta}^T}_{\operatorname{eig}(\mathbf{C}_{\mathbf{Y}\mathbf{Y}})} \boldsymbol{\Theta}_1 = [1 \ 0 \ \dots \ 0] \cdot \begin{bmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \lambda_n \end{bmatrix} \cdot [1 \ 0 \ \dots \ 0]^T$$

$$\sigma_{X_1, \max}^2 = \lambda_1$$

where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$

M. Verleysen
UCL
8

Choice of direction [3/3]

- Classical result :

Best choice for \mathbf{u}_1 is the eigenvector Θ_1 associated to the largest eigenvalue λ_1 of matrix \mathbf{C}_{YY} .

- In the space orthogonal to \mathbf{u}_1 :

Best choice for \mathbf{u}_2 is the eigenvector Θ_2 associated to the largest eigenvalue λ_2 of matrix \mathbf{C}_{YY} .

- And so on ...

M. Verleysen
UCL
9

Overview

- Principal Component Analysis: a reminder
- Probabilistic PCA
- Robust probabilistic PCA
- Mixtures of (robust) probabilistic PCA
- Experiments

M. Verleysen
UCL
10

	<h2 style="text-align: center;">Probabilistic PCA</h2>
<p>M. Verleysen UCL 11</p>	<ul style="list-style-type: none"> ■ Idea #1: generative probabilistic model with latent variables \mathbf{x} ■ Idea #2: linear dependency between y_n and x_n ■ Idea #3: Gaussian distribution of latent variables ■ Idea #4: Gaussian additive noise $\{y_n\}_{n=1}^N \subset \mathbb{R}^D, \{x_n\}_{n=1}^N \subset \mathbb{R}^J, J < N$ <div style="border: 1px solid black; padding: 5px; width: fit-content; margin: 10px auto;"> $P(x) = \mathcal{N}(x 0, I_J)$ $P(y x) = \mathcal{N}(y Wx + \mu, \tau^{-1}I_D)$ </div>

	<h2 style="text-align: center;">Probabilistic PCA</h2>
<p>M. Verleysen UCL 12</p>	<div style="display: flex; align-items: center; justify-content: space-between;"> <div style="text-align: center;"> $P(x) = \mathcal{N}(x 0, I_J)$ $P(y x) = \mathcal{N}(y Wx + \mu, \tau^{-1}I_D)$ </div> <div style="text-align: right;"> <p>Isotropic: natural (what else?)</p> <p>Arbitrary center and variance: no problem because compensated by μ and W</p> </div> </div> <p style="margin-left: 40px;">↕</p> <p>Gaussians: natural, and mathematical convenience!</p> <ul style="list-style-type: none"> ■ Marginal distribution: <div style="border: 1px solid black; padding: 5px; width: fit-content; margin: 10px auto;"> $P(y) = \int_{\mathbf{x}} P(y x)P(x)dx = \mathcal{N}(y \mu, \Sigma)$ </div> <p style="margin-left: 40px;">where $\Sigma = WW^T + \tau^{-1}I_D$</p>

Probabilistic PCA training

- Finding parameters $\theta = \{W, \mu, \tau\}$ in

$$P(x) = N(x|O, I_J)$$

$$P(y|x) = N(y|Wx + \mu, \tau^{-1}I_D)$$

- How ? By finding the parameters that lead to maximum likelihood of the observations y_n

$$\theta_{\text{ML}} = \arg \max_{\theta} \left(\prod_n P(y_n) \right) = \arg \max_{\theta} \left(\log \left(\prod_n P(y_n) \right) \right)$$

$$= \arg \min_{\theta} \left(- \sum_n \log(P(y_n)) \right)$$

$$= \arg \min_{\theta} \left(\sum_n \frac{\tau}{2} (y_n - \mu)^T \Sigma^{-1} (y_n - \mu) + \frac{N}{2} \log 2\pi\tau^{-1} \right)$$

M. Verleysen
UCL
13

Probabilistic PCA training

- How to find θ_{ML} ?

$$\theta_{\text{ML}} = \arg \min_{\theta} \left(\sum_n \frac{\tau}{2} (y_n - \mu)^T \Sigma^{-1} (y_n - \mu) + \frac{N}{2} \log 2\pi\tau^{-1} \right)$$

- \equiv fitting a multivariate Gaussian distribution to the observations, with a constrained covariance matrix

$$\Sigma = WW^T + \tau^{-1}I_D$$

- How: EM algorithm

M. Verleysen
UCL
14

	Probabilistic PCA
	<ul style="list-style-type: none">■ Is it useful? (does the probabilistic formulation add something to the deterministic one?)

M. Verleysen
UCL
15

	Probabilistic PCA
	<ul style="list-style-type: none">■ Is it useful? (does the probabilistic formulation add something to the deterministic one?)■ Answer: no ☹️ Tipping and Bishop showed that the axes found by PPCA span the same subspace as those found by PCA (standard equivalence between ML + Gaussian noise hypothesis, and LMS criterion)

M. Verleysen
UCL
16

	Probabilistic PCA
	<ul style="list-style-type: none">■ Is it useful? (does the probabilistic formulation add something to the deterministic one?)■ Answer: yes 😊 The probabilistic formulation makes it possible to introduce new, more realistic hypotheses

M. Verleysen
UCL
17

	Probabilistic PCA
	<ul style="list-style-type: none">■ Is it useful? (does the probabilistic formulation add something to the deterministic one?)■ Answer: yes 😊 The probabilistic formulation makes it possible to introduce new, more realistic hypotheses■ Ex: LMS criterion (\equiv Gaussian hyp.) does not correspond to natural goals!!!

M. Verleysen
UCL
18

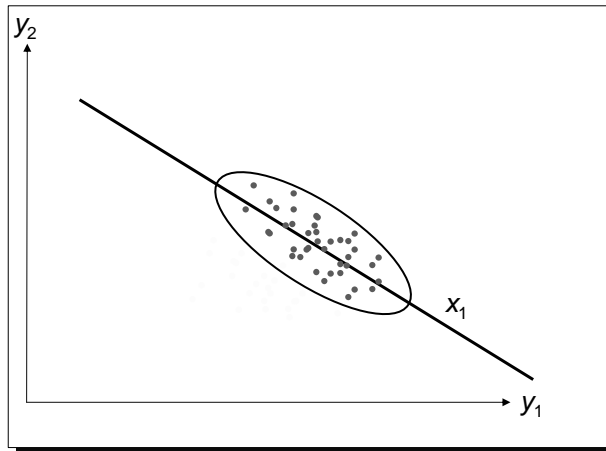
Overview

- Principal Component Analysis: a reminder
- Probabilistic PCA
- Robust probabilistic PCA
- Mixtures of (robust) probabilistic PCA
- Experiments

M. Verleysen
UCL
19

PCA and outliers

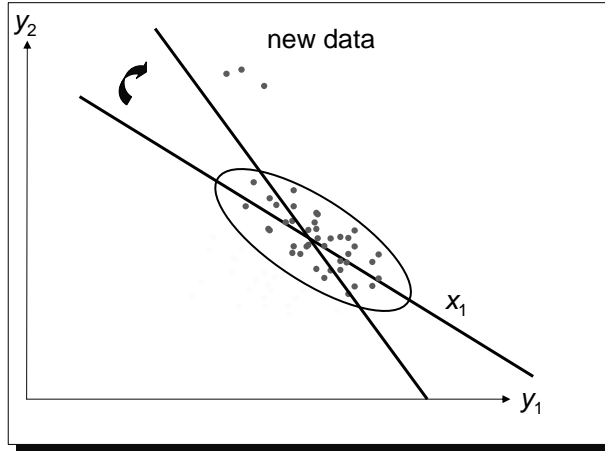
- With "easy" data



M. Verleysen
UCL
20

PCA and outliers

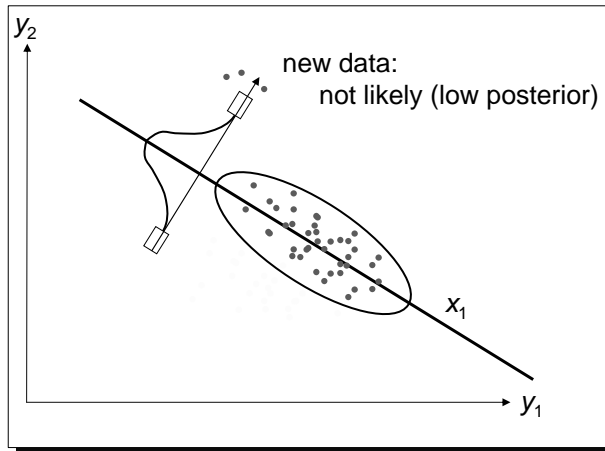
- With a few outliers



M. Verleysen
UCL
21

PCA and outliers

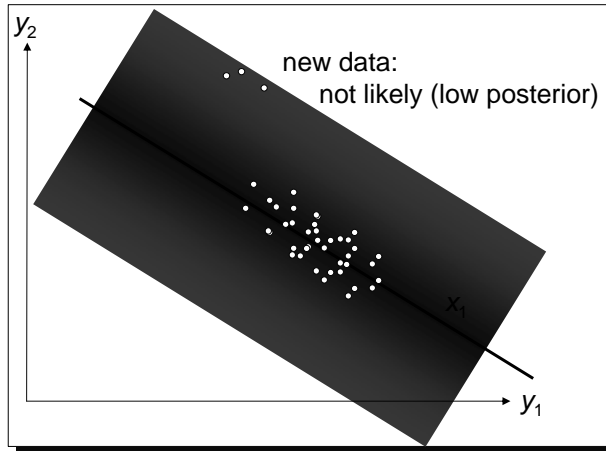
- Why is it the case?



M. Verleysen
UCL
22

PCA and outliers

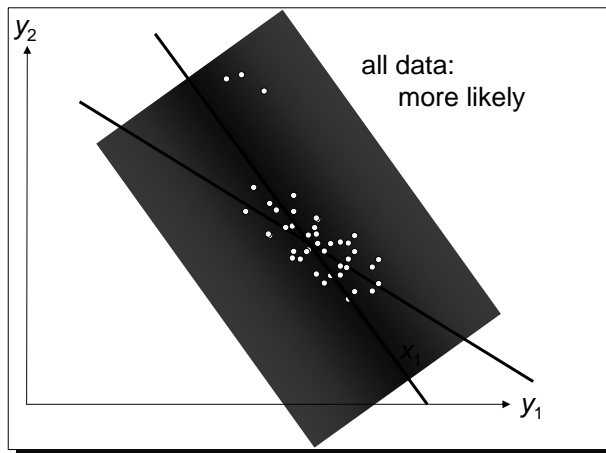
- Why is it the case?



M. Verleysen
UCL
23

PCA and outliers

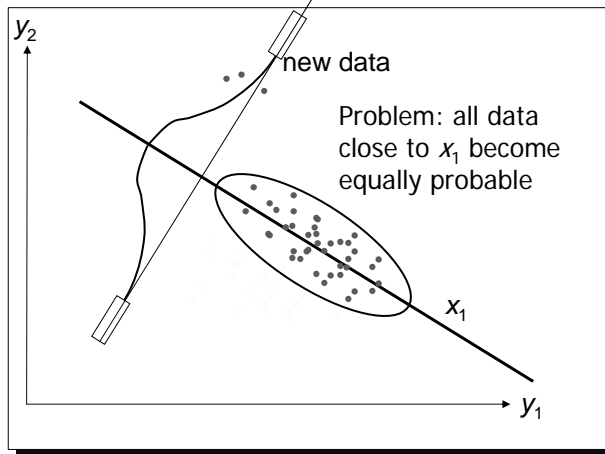
- Why is it the case?



M. Verleysen
UCL
24

PCA robust to outliers

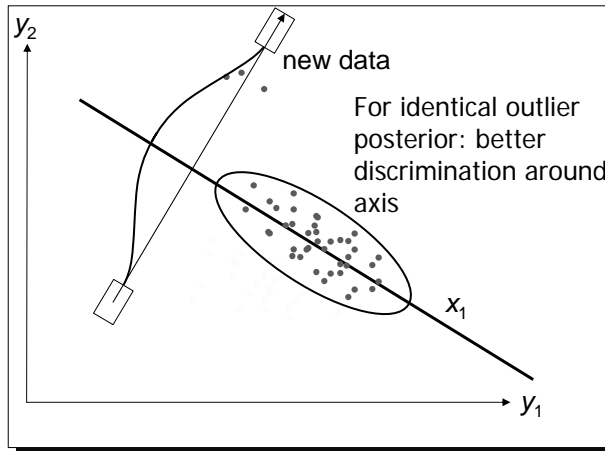
- Increasing the probability of outliers: Gaussian



M. Verleysen
UCL
25

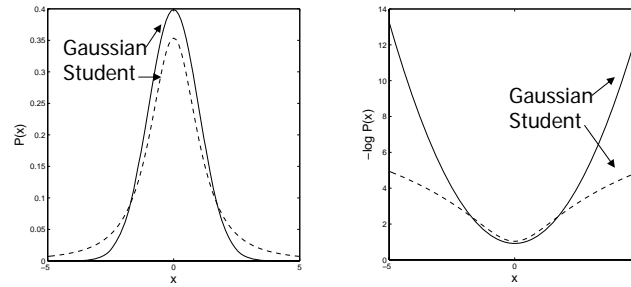
PCA robust to outliers

- Increasing the probability of outliers: Student



M. Verleysen
UCL
26

Student- t distribution



M. Verleysen
UCL
27

Robust PPCA

- Probabilistic PCA

$$P(x) = N(x|O, I_J)$$

$$P(y|x) = N(y|Wx + \mu, \tau^{-1}I_D)$$

- Robust Probabilistic PCA

$$P(x) = \text{St}(x|O, I_J, \nu)$$

$$P(y|x) = \text{St}(y|Wx + \mu, \tau^{-1}I_D, \nu)$$

Identical for simplicity !

M. Verleysen
UCL
28

Robust PPCA

- Model
$$\begin{cases} P(x) = \text{St}(x|O, I_J, \nu) \\ P(y|x) = \text{St}(y|Wx + \mu, \tau^{-1}I_D, \nu) \end{cases}$$
- But Student- t = infinite mixture of Gaussians:

$$\text{St}(y|\mu, \Sigma, \nu) = \int_0^\infty \mathcal{N}\left(y|\mu, \frac{1}{u}\Sigma\right) \text{Ga}\left(u\left|\frac{\nu}{2}, \frac{\nu}{2}\right.\right) du, \nu > 0$$
 where $\text{Ga}(u|..)$ is a Gamma distribution

- Therefore the generative model is

$$\begin{aligned} P(u) &= \text{Ga}\left(u\left|\frac{\nu}{2}, \frac{\nu}{2}\right.\right) \\ P(x|u) &= \mathcal{N}\left(x\left|O, \frac{1}{u}I_J\right.\right) \\ P(y|x, u) &= \mathcal{N}\left(y\left|Wx + \mu, \frac{1}{u\tau}I_D\right.\right) \end{aligned}$$

M. Verleysen
UCL
29

Robust PPCA

- Model
$$\begin{aligned} P(u) &= \text{Ga}\left(u\left|\frac{\nu}{2}, \frac{\nu}{2}\right.\right) \\ P(x|u) &= \mathcal{N}\left(x\left|O, \frac{1}{u}I_J\right.\right) \\ P(y|x, u) &= \mathcal{N}\left(y\left|Wx + \mu, \frac{1}{u\tau}I_D\right.\right) \end{aligned}$$
- Good news: the posterior is tractable

$$P(y) = \int_0^\infty \int_X P(y|x, u)P(x|u)P(u)dx = \text{St}(y|\mu, \Sigma, \nu)$$
 where again $\Sigma = WW^T + \tau^{-1}I_D$

M. Verleysen
UCL
30

Robust PPCA training

- Finding parameters $\theta = \{W, \mu, \tau, \nu\}$ in

$$P(u) = \text{Ga}\left(u \mid \frac{\nu}{2}, \frac{\nu}{2}\right)$$

$$P(x|u) = \text{N}\left(x \mid O, \frac{1}{u} I_J\right)$$

$$P(y|x, u) = \text{N}\left(y \mid Wx + \mu, \frac{1}{u\tau} I_D\right)$$

- How ? By finding the parameters that lead to maximum likelihood of the observations y_n

$$\theta_{\text{ML}} = \underset{\theta}{\text{argmin}} \left(- \sum_n \log(P(y_n)) \right)$$

M. Verleysen
UCL
31

Robust PPCA advantages

1. Only one parameter to fix in advance (the dimension of the latent –projection- space)
2. Natural framework for an extension to mixtures

M. Verleysen
UCL
32

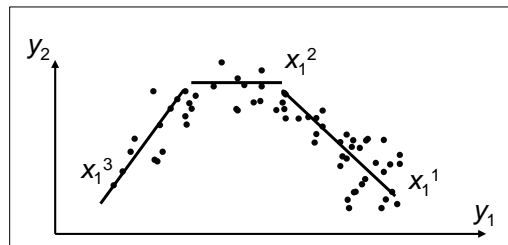
Overview

- Principal Component Analysis: a reminder
- Probabilistic PCA
- Robust probabilistic PCA
- Mixtures of (robust) probabilistic PCA
- Experiments

M. Verleysen
UCL
33

Mixtures of (robust) PPCA

- (P)PCA: linear dependencies in data only
- Mixtures of (P)PCA: nonlinear dependencies, through mixtures of linear manifolds

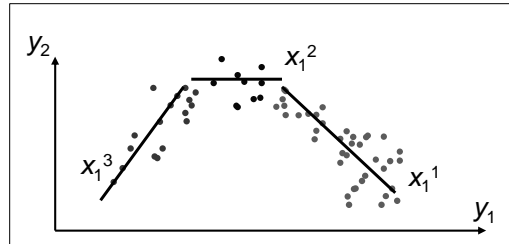


- Latent variable model, *no* common projection space

M. Verleysen
UCL
34

Mixtures of (robust) PPCA

- Concept of (hidden) cluster membership



$$P(y) = \sum_k \pi_k P_k(y)$$

Single PCA model

Mixture proportions $\pi_k > 0$

$$\sum_k \pi_k = 1$$

M. Verleysen
UCL
35

Mixtures of (robust) PPCA

- $z = [z_1, z_2, \dots, z_k]$ with
 $z_k = 1$ if y_n was generated by component k
 $z_k = 0$ otherwise

- Latent variable model

$$P(z) = \prod_k \pi_k^{z_k}$$

$$P(u|z) = \prod_k \text{Ga}\left(u_k \mid \frac{v_k}{2}, \frac{v_k}{2}\right)^{z_k}$$

$$P(x|u, z) = \prod_k \text{N}\left(x \mid 0, \frac{1}{u_k} I_J\right)^{z_k}$$

Possible different dimensionalities

$$P(y|x, u, z) = \prod_k \text{N}\left(y \mid W_k x + \mu_k, \frac{1}{u_k \tau_k} I_D\right)^{z_k}$$

M. Verleysen
UCL
36

Mixtures of (robust) PPCA training	
M. Verleysen UCL 37	<ul style="list-style-type: none"> ■ Finding parameters $\theta = \{(W_k, \mu_k, \tau_k, \nu_k, \pi_k)\}_{k=1\dots K}$ in $P(z) = \prod_k \pi_k^{z_k}$ $P(u z) = \prod_k \text{Ga}\left(u_k \mid \frac{\nu_k}{2}, \frac{\nu_k}{2}\right)^{z_k}$ $P(x u, z) = \prod_k \text{N}\left(x \mid 0, \frac{1}{u_k} I_J\right)^{z_k}$ $P(y x, u, z) = \prod_k \text{N}\left(y \mid W_k x + \mu_k, \frac{1}{u_k \tau_k} I_D\right)^{z_k}$ ■ How ? By finding the parameters that lead to maximum likelihood of the observations y_n $\theta_{\text{ML}} = \underset{\theta}{\text{argmin}} \left(- \sum_n \log(P(y_n)) \right)$

Mixtures of (robust) PPCA training	
M. Verleysen UCL 38	<ul style="list-style-type: none"> ■ In practice: EM algorithm ■ E-step: fix the model parameters, and compute <i>expectations</i> (over an approximate distribution) of latent variables ■ M-step: update the model parameters to <i>maximize</i> the likelihood ■ Two hyper-parameters: <ul style="list-style-type: none"> – the number of components – the dimensionality of the latent representations

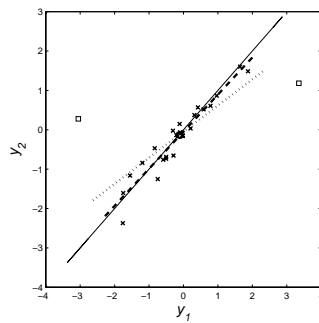
Overview

- Principal Component Analysis: a reminder
- Probabilistic PCA
- Robust probabilistic PCA
- Mixtures of (robust) probabilistic PCA
- Experiments

M. Verleysen
UCL
39

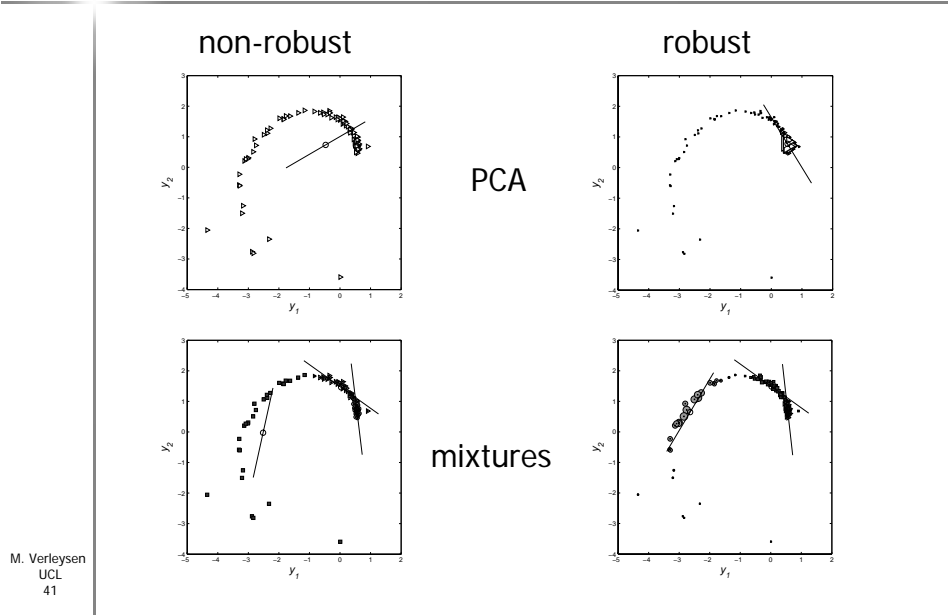
Experiments

- PCA and robust PCA



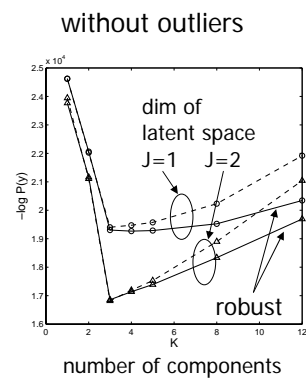
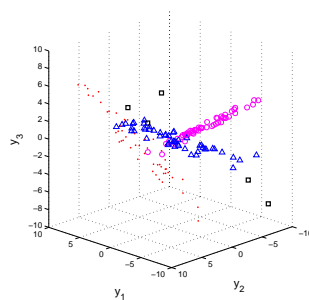
M. Verleysen
UCL
40

Experiments



Experiments

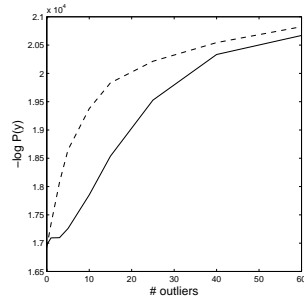
- Three 3-D Gaussian clusters
 - Diagonal covariance matrix: $\text{diag}(5, 1, 0.2)$ before rotation
 - Intrinsic dimensionality: 2



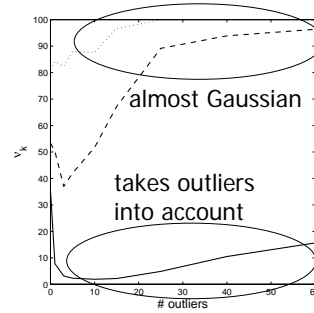
Experiments

- True model ($K = 3, J = 2$)
- Outliers

Performances on test set



Degree of freedom parameters



M. Verleysen
UCL
43

Experiments

- USPS handwritten digit dataset
 - around 700 images of "2"
 - around 700 images of "3"
 - around 100 images of "0" (as outliers)
- Experiment: two 1-D components (two clusters)
- Illustrations: images close from the 1-D subspaces

standard mixtures



robust mixtures



M. Verleysen
UCL
44

Conclusions

- Probabilistic formulation of PCA:
 - identical to PCA
 - but possible to introduce other hypotheses
- Robust PCA:
 - replacing Gaussians by Student- t
 - makes outliers more likely \rightarrow lower influence on the model
- Mixtures of PCA
 - latent variable model
 - better than mixtures of Gaussians through regularization (dimensionality of latent space)
- Mixtures of Robust PCA
 - replacing Gaussians by Student- t