
Principal Component Analysis

*Prof. Nicholas Zabaras
University of Notre Dame
Notre Dame, IN, USA*

*Email: nzabaras@nd.edu
URL: <https://www.zabaras.com/>*

November 7, 2017



Contents

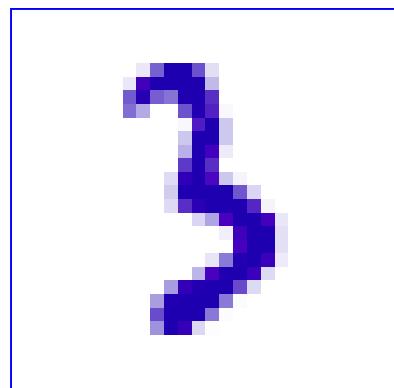
- Continuous Latent Variable Model, Low-Dimensional Manifold of a Data Set, Generative Point of View, Unidentifiability
- Principal Component Analysis, Maximum variance formulation, Minimum-error formulation, PCA Reconstruction, PCA and SVD, Canonical Correlation Analysis
- Applications of PCA (off-line Digit Images), Whitening of the data with PCA, PCA for Visualization, PCA for High-Dimensional Data
- Probabilistic PCA, Problem setup, Maximum likelihood PCA, EM algorithm for PCA, EM Missing Values, EM for $\sigma^2 \rightarrow 0$, Model Selection

Following closely Chris Bishop's PRML book (Chapter 12)



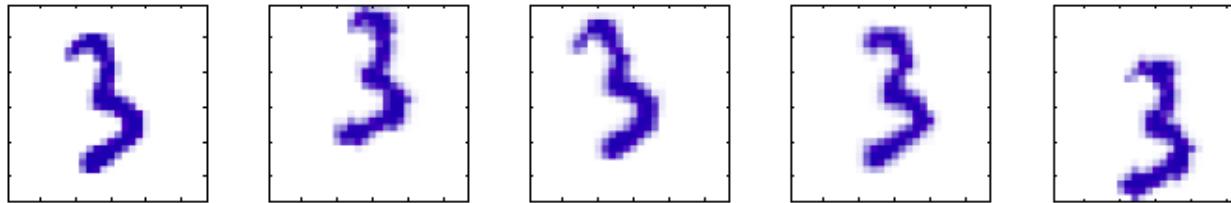
Continuous Latent Variables

- In many data sets, the data points lie close to a manifold of much lower dimensionality than that of the original data space.
- Consider a data set constructed by taking one of the off-line digits, represented by a 64×64 pixel grey-level image.
- Embed this in a larger image of size 100×100 by padding with pixels having the value zero (white pixels).
- Each image is represented by a point in 10,000 dimensional space.



Continuous Latent Variables

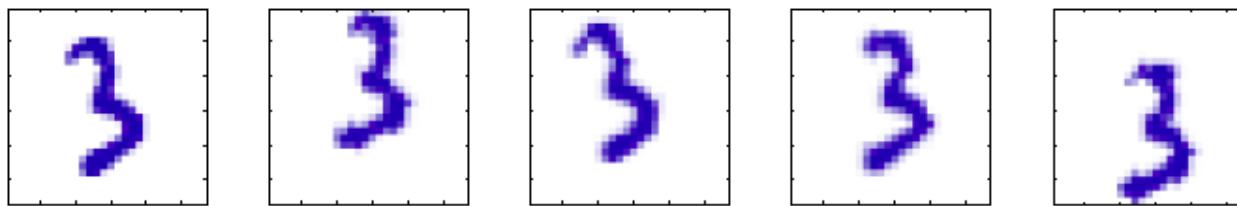
- We create multiple copies in which the location and orientation of the digit is varied at random.



- In this data set, there are 3 DOF of variability:
 - Vertical displacement
 - Horizontal displacement and
 - Rotation

Intrinsic Dimensionality of a Data Set

- The intrinsic dimensionality of the data set is three (vertical and horizontal displacement and rotation)



- The manifold is non-linear: when translating the digit past a particular pixel, that pixel value goes from 0 (white) to 1 (black) and back to zero again. This is a nonlinear function of the digit position.
- The translation and rotation parameters are latent variables: we observe only the image vectors without knowing the translation or rotation variables used to create them.

Other Latent Variables

- For real digit image data, additional DOF arise from
 - Scaling
 - Complex deformations due to the variability of an individual's writing style,
 - Etc.
- The number of such degrees of freedom will still be small compared to the dimensionality of the data set.
- For data compression, we are interested to explore the manifold structure.

Generative Point of View

- The data points are often not confined precisely to a smooth low-dimensional manifold.
- The departure of data points from the manifold is interpreted as ‘noise’.
- **Generative view:**
 - (a) Select a point within the manifold according to some latent variable distribution $p(z)$ and
 - (b) Generate an observed data point by adding noise, drawn from some conditional distribution of the data variables given the latent variables $p(x|z)$.



Generative Point of View

- Example – Linear Gaussian latent variable model:
 - (a) assume Gaussian distributions for the latent variables
 $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ and the observed variables
 $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|W\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$
 - (b) make use of a linear-Gaussian dependence of the observed variables on the state of the latent variables,
 $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|W\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$
- This leads to a probabilistic formulation of principal component analysis (PCA).

Generative Point of View

□ Example – Factor Analysis:

- (a) assume Gaussian distributions for the latent variables

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

and the observed variables

$$p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i | \mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}), \\ \mathbf{x}_i \in \mathbb{R}^D, \mathbf{z}_i \in \mathbb{R}^M, \mathbf{W} \in \mathbb{R}^{D \times M}, \boldsymbol{\Psi} \in \mathbb{R}^{D \times D}$$

- (b) We take $\boldsymbol{\Psi}$ to be diagonal. This overall model is called **factor analysis** or FA.
- The special case in which $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$ is called **probabilistic principal components analysis** or PPCA.

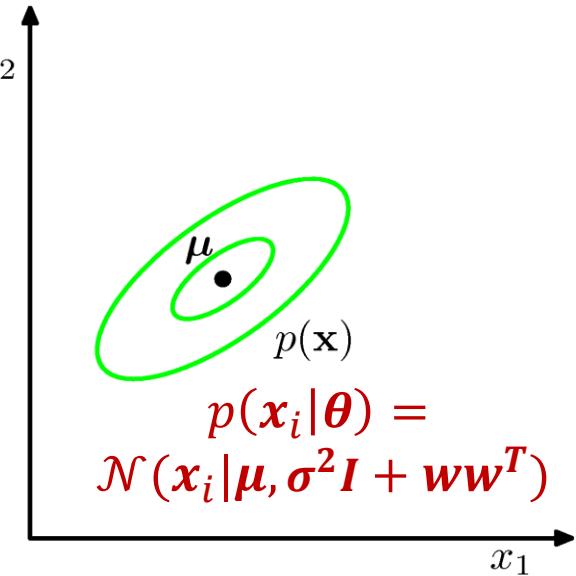
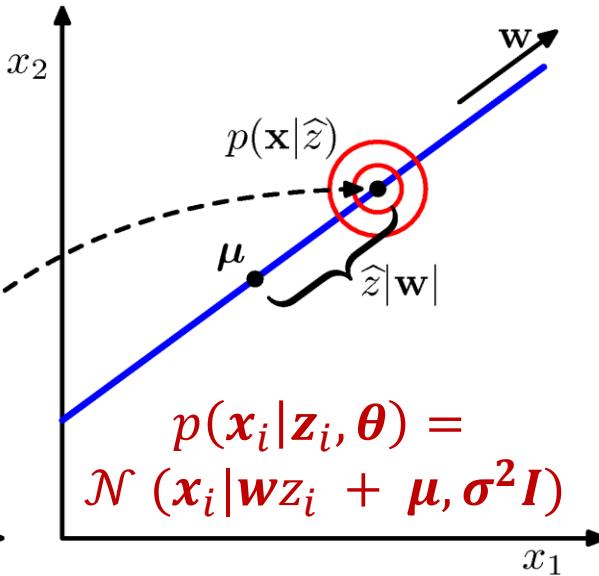
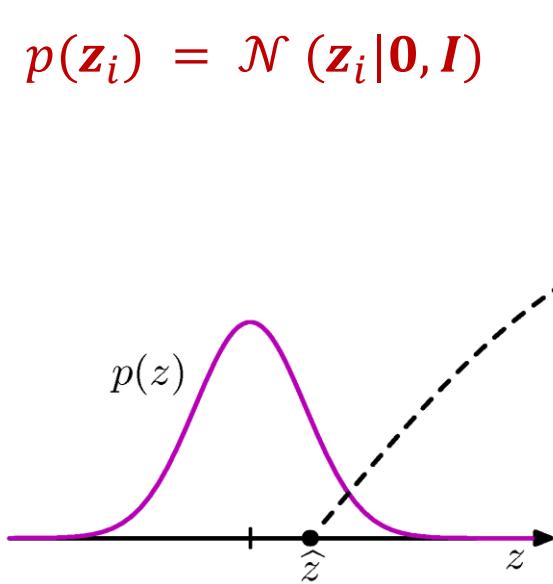


Generative Point of View

- PPCA is illustrated for $D = 2, M = 1$ (latent dimension).

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \int \mathcal{N}(\mathbf{x}_i|W\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}) \mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) d\mathbf{z}_i = \mathcal{N}(\mathbf{x}_i|W\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + W\boldsymbol{\Sigma}_0 W^T)$$

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|\mathbf{0}, \mathbf{I})$$



$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

$$p(\mathbf{x}_i|\mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|W\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}, \boldsymbol{\Psi} + W\boldsymbol{\Sigma}_0 W^T)$$

- Based on the obtained marginal, we see that without sacrificing generality, we can simplify as:

$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|\mathbf{0}, \mathbf{I})$$

$$p(\mathbf{x}_i|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}, \boldsymbol{\Psi} + W\boldsymbol{\Sigma}_0 W^T)$$

Inference on the Latent Factors

- Using $p(\mathbf{x}_i|\mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i|\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi})$, $p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$, Bayes rule and standard eqs. for linear Gaussian systems, we can derive the posterior of the latent variables:

$$\begin{aligned} p(\mathbf{z}_i|\mathbf{x}_i, \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{z}_i|\mathbf{m}_i, \boldsymbol{\Sigma}_i), \\ \mathbf{m}_i &= \boldsymbol{\Sigma}(\mathbf{W}^T \boldsymbol{\Psi}^{-1}(\mathbf{x}_i - \boldsymbol{\mu}) + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0) \\ \boldsymbol{\Sigma} &= (\boldsymbol{\Sigma}_0^{-1} + \mathbf{W}^{-1} \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} \end{aligned}$$

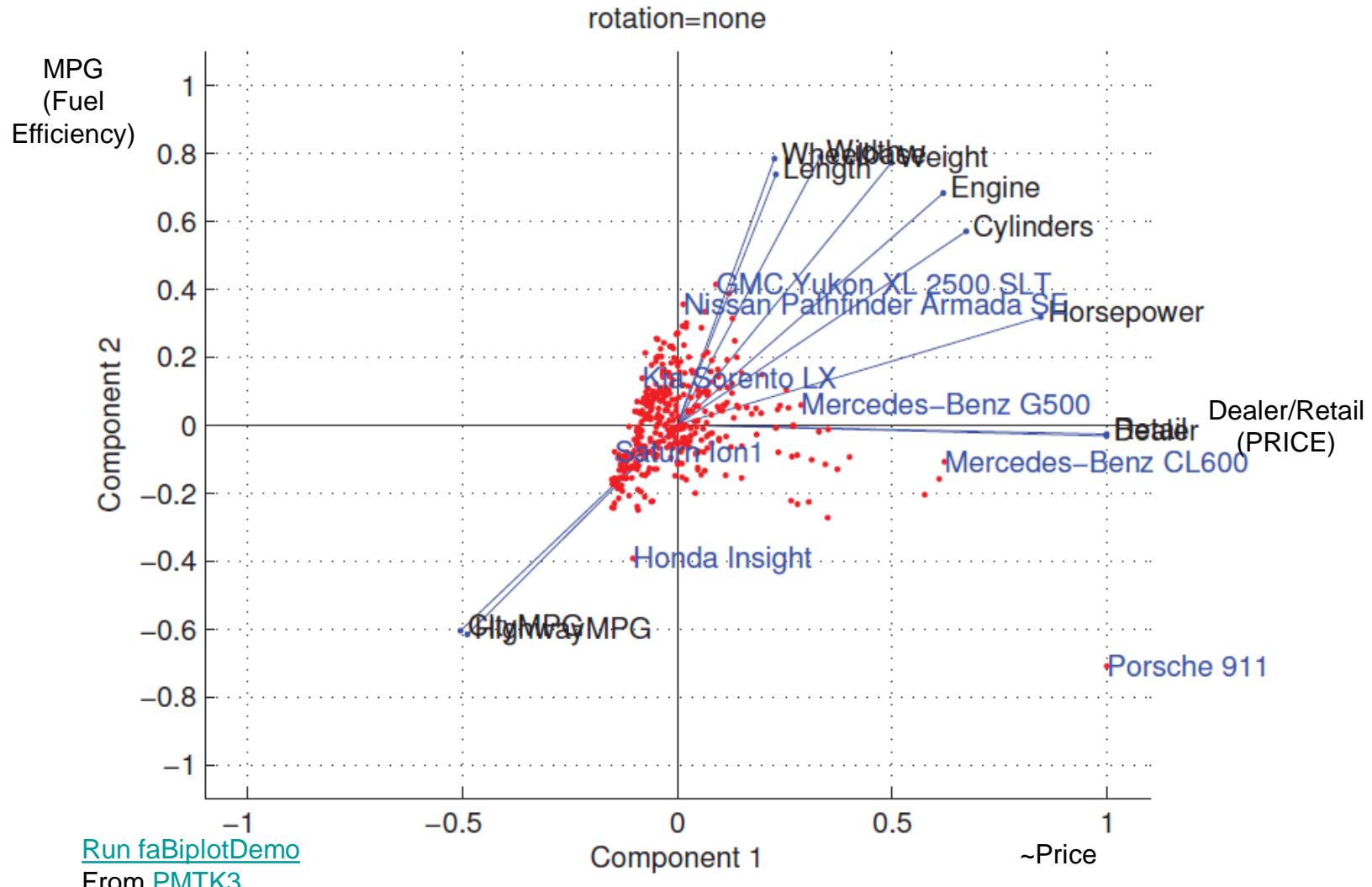
- Computing $\boldsymbol{\Sigma}$ takes $\mathcal{O}(M^3 + M^2D)$ time, and computing each $\mathbf{m}_i = \mathbb{E}[\mathbf{z}_i|\mathbf{x}_i, \boldsymbol{\theta}]$ takes $\mathcal{O}(M^2 + MD)$.
- The \mathbf{m}_i are sometimes called **the latent scores**, or **latent factors**.

Example Based on Factor Analysis

- Consider a dataset of $D = 11$ variables and $N = 387$ cases describing aspects of cars (engine size, # of cylinders, MPG, etc.).
- We fit a $M = 2$ dimensional model. We plot the $\mathbf{m}_i = \mathbb{E}[\mathbf{z}_i | \mathbf{x}_i, \theta]$ scores as points in \mathbb{R}^2 , to visualize the data.
- We also project unit vectors corresponding to each of the feature dimensions, $\mathbf{e}_1 = (1, 0, \dots, 0)$, $\mathbf{e}_2 = (0, 1, 0, \dots, 0)$, etc. into the low dimensional space (blue lines).
- We see that the horizontal axis represents price corresponding to the features labeled “dealer” and “retail” (expensive cars on the right).
- The vertical axis represents fuel efficiency (in terms of MPG) versus size: heavy vehicles are less efficient and are higher up, whereas light vehicles are lower down.
- Unfortunately, interpreting latent variable models is fraught with difficulties.



2D Projection of Car Data



Unidentifiability

- Consider an arbitrary rotation such that: $\tilde{\mathbf{W}} = \mathbf{W}\mathbf{R}$.
- We can immediately see that the likelihood function $p(\mathbf{x}_i|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T)$ remains the same:

$$\tilde{\mathbf{W}}\tilde{\mathbf{W}}^T = \mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T = \mathbf{W}\mathbf{W}^T$$

- We need to remove $M(M - 1)/2$ DOF since that is the number of orthonormal matrices of size $M \times M$.
- Then the FA model has $D + MD - M(M - 1)/2$ free parameters (excluding the mean), where the first term arises from $\boldsymbol{\Psi}$. This needs to be $\leq D(D + 1)/2$, i.e. the number of parameters in an unconstrained symmetric covariance matrix.
- This gives us an upper bound on M : $M_{max} = \lfloor D + 0.5(1 -$

Unidentifiability

To address the unidentifiability the following options are available:

- Making \mathbf{W} orthonormal: This is the case with PCA
 - columns arranged in order of decreasing variance
 - Making \mathbf{W} lower triangular: The first feature is then generated only by the 1st latent variable, the second by the first two, etc.
 - The number of parameters is equal to M_{max} .
 - One needs to properly select the first M visible variables as they affect the latent factors.
- [Lopes, H. and M. West \(2004\). Bayesian model assessment in factor analysis. *Statistica Sinica* 14, 41– 67.](#)



Unidentifiability

- Sparsity promoting priors on the weights: Use sparse factor analysis that forces some entries in W to zero. They may include ℓ_1 regularization, automatic relevance determination or spike and slab priors. But note that this does not ensure a unique MAP estimate.
- Choosing an informative rotation matrix R : Choosing appropriate R to increase interpretable sparse W .
- Using non-Gaussian priors for the latent factors in $p(\mathbf{z}_i)$: This leads to Independent Component Analysis (ICA).

- Zou, H. (2006). [The adaptive Lasso and its oracle properties](#). *J. of the Am. Stat. Assoc.*, 1418–1429.
- Bishop, C. (1999). [Bayesian PCA](#). In *NIPS*.
- Archambeau, C. and F. Bach (2008). [Sparse probabilistic projections](#). In *NIPS*.
- Kaiser, H. (1958). <https://www.dropbox.com/s/o4hejvydu1ntsk6/Kaiser1958.pdf?dl=0>. *Psychometrika* 23(3).

Latent Models

- PCA arises as the Maximum Likelihood solution to a linear Gaussian Latent Variable Model.
- Probabilistic formulation allows:
 - Use of EM for parameter estimation
 - Using mixtures of PCA models
 - A Bayesian formalism in addition allows computing the number of principal components directly from the data.
- Non Gaussian latent variable models lead to Independent Component Analysis.
- Models can also be considered with a non-linear relation between latent and observed variables.

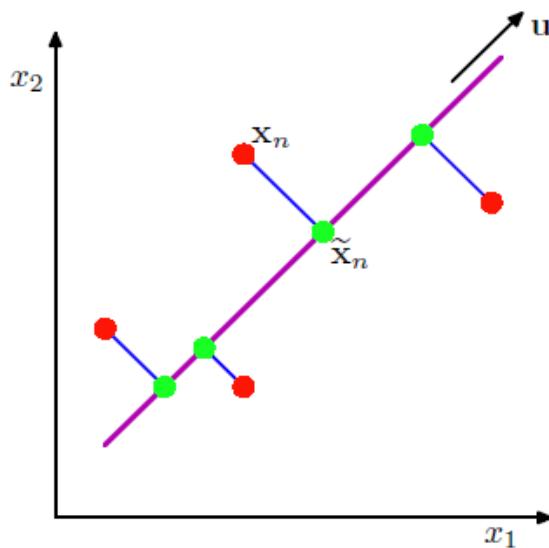
Principal Component Analysis

- Principal component analysis (PCA) – also known as Karhunen-Loeve expansion - is used widely for (Jolliffe, 2002)
 - Dimensionality reduction
 - Lossy data compression
 - Feature selection and
 - Data Visualization
 - Principal component analysis seeks a lower dimensionality linear space (principal subspace) such that
 - The orthogonal projection of the data points onto this subspace **maximizes the variance of the projected points** (Hotelling, 1933)
 - An alternative definition of PCA is based on **minimizing the sum-of-squares of the projection errors** (Pearson, 1901)
-
- Jolliffe, I. T. (2002). Principal Component Analysis (Second ed.). Springer.
 - Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology* **24**, 417–441.
 - Pearson, K. (1901). On lines and planes of closest fit to systems of points in space. *The London, Edinburgh and Dublin Philosophical Magazine and Journal of Science, Sixth Series* **2**, 559–572.



Principal Component Analysis

- PCA seeks a space of lower dimensionality (magenta line) such that:
 - (1) the orthogonal projection of the data points (red dots) onto this subspace **maximizes the variance of the projected points** (green dots).
 - (2) minimizing the sum-of-squares of the projection errors (blue lines)



- [J. Shlens \(2005\). *A Tutorial on Principal Component Analysis*](#)

PCA: Maximum Variance Formulation

- Consider a data set of observations $\{x_n\}, n = 1, \dots, N$, where x_n has dimensionality D .
- Our goal is to project the data onto a space having dimensionality $M < D$ (principal subspace) while maximizing the variance of the projected data (*M here is fixed*).
- Let $\{u_i\}, i = 1, \dots, M$ the basis vectors ($(D \times 1)$ vectors) of the principal subspace.
- We define the sample mean ($(D \times 1)$ vector) by:
$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$
 and the sample covariance matrix ($(D \times D)$ matrix) by:

$$S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T$$



PCA with 1D Principal Subspace

- Consider first the projection onto 1D space ($M = 1$).
- We define the direction of this subspace using vector \mathbf{u}_1 , which we choose to be a unit vector: $\mathbf{u}_1^T \mathbf{u}_1 = 1$
- Each data point \mathbf{x}_n is then projected onto the scalar $\mathbf{u}_1^T \mathbf{x}_n$
- The mean of the projected data is given as: $\mathbf{u}_1^T \bar{\mathbf{x}}$
- Similarly, the variance of the projected data is:

$$\frac{1}{N} \sum_{n=1}^N \left\{ \mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}} \right\}^2 = \frac{1}{N} \sum_{n=1}^N \left\{ \mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}} \right\} \left\{ \mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}} \right\}^T = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

- Key idea of PCA: Maximize the projected variance $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ with respect to \mathbf{u}_1 under the constraint $\mathbf{u}_1^T \mathbf{u}_1 = 1$



PCA with 1D Principal Subspace

- Introduce the Lagrange multiplier λ_1 and define the unconstrained maximization of

$$\max_{\mathbf{u}_1, \lambda_1} \left\{ \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1) \right\}$$

- We can see immediately that the solution satisfies:

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

- \mathbf{u}_1 must be an eigenvector of \mathbf{S} with eigenvalue λ_1
- From the eigen-problem note that the variance of the projected data is $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1$ so λ_1 needs to be the largest eigenvalue of \mathbf{S} .
- \mathbf{u}_1 is called the first principal component.

PCA with 1D Principal Subspace

- Additional Principal Components: maximize the projected variance amongst all possible directions orthogonal to those already considered.
- Using induction, you can show: For an M -dimensional projection space, the optimal linear projection for which the variance of the projected data is maximized is defined by the M eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_M$ of the data covariance matrix S corresponding to the largest eigenvalues $\lambda_1, \dots, \lambda_M$.
- The computational cost of finding the M principal components (i.e. finding the first M eigenvalues and eigenvectors of S) is $\mathcal{O}(MD^2)$.

Total Variance

- ❑ Essentially PCA becomes the calculation of the eigenvectors of the data covariance matrix corresponding to the largest M eigenvalues.
- ❑ $\sum_{i=1}^D \lambda_i$ is the **total variance**.
- ❑ The contribution to the total variance by one component u_i is measured as:

$$\frac{\lambda_i}{\sum_{i=1}^D \lambda_i}$$



Proof by Induction: M-Dimensional Subspace

- We have shown our result for $M = 1$. Assume it is valid for projection space of dimensionality M . We will show that it is then valid for projection space of dimensionality $M + 1$.
- We need to maximize the following:

$$\max_{\mathbf{u}_{M+1}} \left\{ \mathbf{u}_{M+1}^T S \mathbf{u}_{M+1} + \lambda_{M+1} \left(1 - \mathbf{u}_{M+1}^T \mathbf{u}_{M+1} \right) + \sum_{i=1}^M \eta_i \mathbf{u}_{M+1}^T \mathbf{u}_i \right\}$$

- The stationary points satisfy:

$$0 = 2S\mathbf{u}_{M+1} - 2\lambda_{M+1}\mathbf{u}_{M+1} + \sum_{i=1}^M \eta_i \mathbf{u}_i$$

- Taking dot product with $\mathbf{u}_j^T, j = 1, \dots, M$ gives $\eta_j = 0$. Thus $S\mathbf{u}_{M+1} = \lambda_{M+1}\mathbf{u}_{M+1} \Rightarrow \mathbf{u}_{M+1}$ eigenvector of S with the variance

in this direction $\mathbf{u}_{M+1}^T S \mathbf{u}_{M+1} = \lambda_{M+1}$

- λ_{M+1} needs to be the next largest eigenvalue after the first M ones \rightarrow The argument by induction holds for any $M \leq D$.

PCA: Minimum Error Formulation

- This is based on projection error minimization.
- Consider a D -dimensional basis $\{\mathbf{u}_i\}$ where $i = 1, \dots, D$ satisfying

$$\mathbf{u}_1^T \mathbf{u}_j = \delta_{ij}$$

- Each data point \mathbf{x}_n can be represented by:

$$\mathbf{x}_n = \sum_{i=1}^D a_{ni} \mathbf{u}_i, \text{ where: } a_{ni} = \mathbf{x}_n^T \mathbf{u}_i$$

- We approximate \mathbf{x}_n using a representation with $M < D$ variables z_{ni} via a projection onto a lower-dim subspace:

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$

where the $\{z_{ni}\}$ depend on the particular data point, whereas the $\{b_i\}$ are the same constants for all data points.

- We are free to choose the $\{\mathbf{u}_i\}$, $\{z_{ni}\}$, and $\{b_i\}$

PCA: Minimum Error Formulation

- Key idea of PCA: Minimize the distortion J introduced by the dimensionality reduction:

$$\min_{z_{nj}, b_j, \mathbf{u}_i} J = \frac{1}{N} \sum_{i=1}^N \| \mathbf{x}_n - \tilde{\mathbf{x}}_n \|^2, \quad \tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$

- Setting the derivatives wrt $\{z_{ni}\}$, and $\{b_i\}$ to zero and using $\mathbf{u}_1^T \mathbf{u}_j = \delta_{ij}$ gives:

$$z_{nj} = \mathbf{x}_n^T \mathbf{u}_j, \quad j = 1, \dots, M$$

$$b_j = \bar{\mathbf{x}}^T \mathbf{u}_j, \quad j = M+1, \dots, D, \text{ where } \bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

- If we now substitute these expressions for $\{z_{ni}\}$, and $\{b_i\}$

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \sum_{i=1}^D (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i - \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i - \sum_{i=M+1}^D (\bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i = \sum_{i=M+1}^D \{(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_i\} \mathbf{u}_i$$

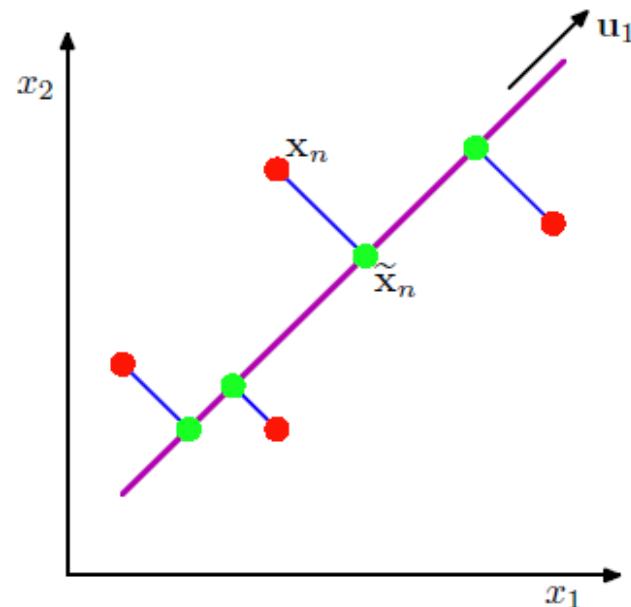
PCA: Minimum Error Formulation

$$x_n - \tilde{x}_n = \sum_{i=M+1}^D \{(x_n - \bar{x})^T u_i\} u_i$$

- We see that the displacement vector $x_n - \tilde{x}_n$ lies in the space orthogonal to the principal subspace, i.e. is a linear combination of $\{u_i\}$ for $i = M + 1, \dots, D$.

- We can also see this geometrically:

- The projected points \tilde{x}_n must lie within the principal subspace, but can move them freely within that subspace.
- The minimum error is then obtained by the orthogonal projection.



PCA: Minimum Error Formulation

- We obtain an expression for the distortion measure J as a function purely of the $\{\mathbf{u}_i\}$ in the form

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \sum_{i=M+1}^D \left\{ (\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_i \right\} \mathbf{u}_i$$
$$J = \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D \left(\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i \right)^2 = \sum_{i=M+1}^D \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i$$
$$J = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$$

- The minimum of J is obtained when $\{\mathbf{u}_i\}, i = M + 1, \dots, D$ are the eigenvectors of \mathbf{S} associated to the smallest eigenvalues.
- Consider $D = 2$ and $M = 1$. Let $\lambda_1 > \lambda_2$. The principal subspace is aligned with the eigenvector with the larger eigenvalue λ_1 , and the min value of $J = \lambda_2$ is obtained by choosing \mathbf{u}_2 corresponding to λ_2 .

PCA: Minimum Error Formulation

- The distortion measure is given as

$$J = \sum_{i=M+1}^D \lambda_i$$

- Obtain the minimum value of J by selecting these eigenvectors corresponding to the $D - M$ smallest eigenvalues. The eigenvectors defining the principal subspace are those corresponding to the M largest eigenvalues.
- The two approaches discussed
 - Maximum Variance and
 - Minimum Errorboth lead to the same main result of the PCA.

Compression of the Original Data Set

- Using the earlier equations, we can derive:

$$\begin{aligned}\tilde{\mathbf{x}}_n &= \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i \\ z_{nj} &= \mathbf{x}_n^T \mathbf{u}_j, j = 1, \dots, M & \Rightarrow \tilde{\mathbf{x}}_n &= \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i + \sum_{i=M+1}^D (\bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i \Rightarrow \\ b_j &= \bar{\mathbf{x}}^T \mathbf{u}_j, j = M+1, \dots, D\end{aligned}$$

Data Compression: $\tilde{\mathbf{x}}_n = \bar{\mathbf{x}} + \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i$

- In deriving the last equation, we used the completeness of \mathbf{u}_i , i.e.
- This clearly shows the compression of the data set: From the D -dimensional \mathbf{x}_n to the M -dimensional vector with components $(\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i)$.

PCA Reconstruction

- To reconstruct the data in the original D -dimensional space from a representation in the M -dimensional principal subspace, we simply use:

$$\tilde{\mathbf{x}}_n = \bar{\mathbf{x}} + \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i$$

- If $M = D$, there is no dimensionality reduction but a rotation to align with the principal components e.g from $\{\mathbf{x}_{n_1}, \dots, \mathbf{x}_{n_D}\}$ to $\{\mathbf{a}_{n_1}, \dots, \mathbf{a}_{n_D}\}$ where:

$$\mathbf{x}_n = \sum_{i=1}^{M=D} a_{ni} \mathbf{u}_i$$

PCA and Singular Value Decomposition

- We have defined the solution to PCA in terms of eigenvectors of the empirical covariance matrix (here denoted as $\widehat{\Sigma}$). The same solution can be obtained based on SVD on the **standardized data matrix X** .

U, V : orthonormal

$\sigma_1, \sigma_2, \dots, \sigma_D$: singular values (≥ 0)

$$\begin{matrix} & D \\ N & = \end{matrix} \begin{matrix} D & N - D \\ \text{---} & \text{---} \\ \sigma_1 & \dots & \sigma_D \\ \text{---} & & 0 \end{matrix} \begin{matrix} D \\ \text{---} \\ D \end{matrix} \quad X = U S V^T$$

- The shaded entries in S and the off diagonal terms are zero. So we don't need to compute the shaded entries of U and S .
- It can be shown that the eigenvectors of $X^T X$ are equal to V (right singular vectors of X) and the eigenvalues of $X^T X$ equal to $D = S^2$ (squared singular values) $X^T X = V S U^T U S V^T = V S^2 V^T \rightarrow (X^T X)V = VD$
- Similarly: $XX^T = USV^T VSU^T = US^2 U^T \rightarrow (XX^T)U = UD$

PCA and Singular Value Decomposition

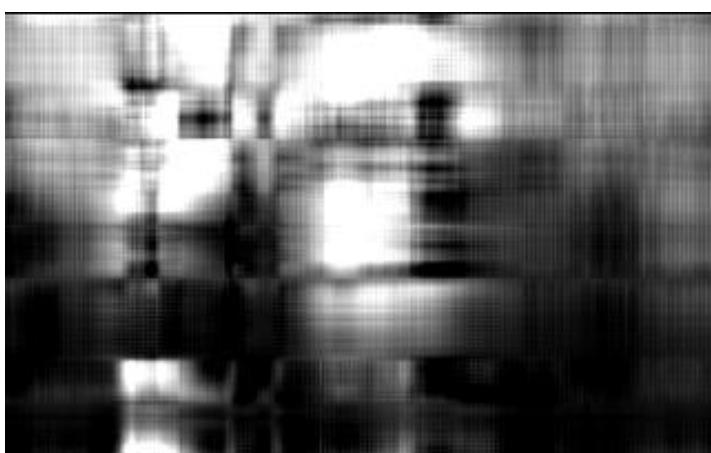
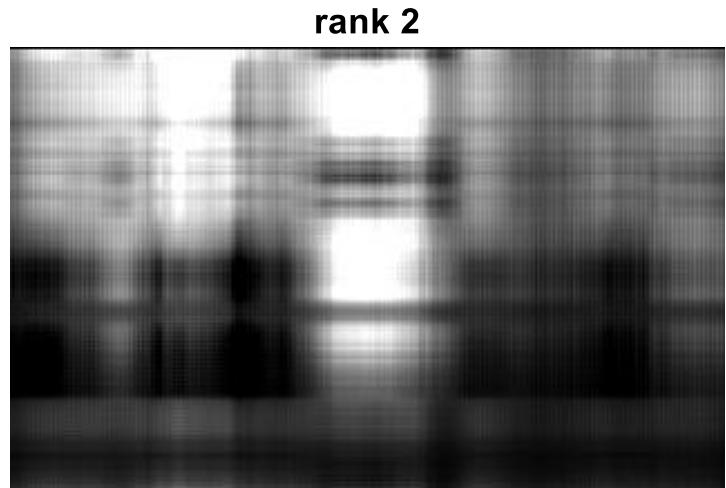
- If the singular values die off quickly, we can produce a rank M approximation of \mathbf{X} as follows:

$$\begin{array}{c} D \\ \boxed{\phantom{\text{D}}} \\ N \end{array} \quad \approx \quad \begin{array}{c} M \\ \boxed{\phantom{\text{M}}} \\ \sigma_1 \\ \vdots \\ \sigma_M \end{array} \quad \begin{array}{c} D \\ \boxed{\phantom{\text{D}}} \\ M \end{array}$$
$$\mathbf{X} \quad \approx \quad \mathbf{U}_M \quad \mathbf{S}_M \quad \mathbf{V}_M^T$$

- The number of parameters for this truncated SVD are: $NM + M + MD$

PCA and Singular Value Decomposition

- Consider the image shown of size 200×320 (rank 200) and the SVD approximations of rank 2, 5, and 20.



Run [svdImageDemo](#)
From [PMTK3](#)



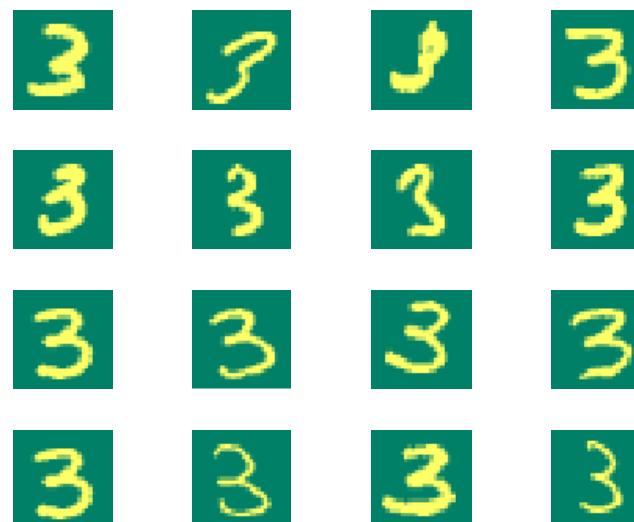
Canonical Correlation Analysis (CCA)

- CCA is a closely related linear dimensionality reduction technique.
- PCA works with a single random variable, CCA considers two (or more) variables.
- CCA finds a corresponding pair of linear subspaces that have high cross-correlation, so that each component within one of the subspaces is correlated with a single component from the other subspace.
- The CCA solution is expressed in terms of a generalized eigenvector problem.
 - Bach, F. R. and M. I. Jordan (2002). [Kernel independent component analysis](#). *Journal of Machine Learning Research* **3**, 1–48.
 - [Hotelling, H.](#) (1936). [Relations between two sets of variables](#). *Biometrika* **28**, 321–377.



Off-line Digit Images: An Example of PCA

- Let us apply PCA in the digit images discussed earlier.
- Consider a synthetic data set obtained by taking one of the digit images and creating multiple copies ($N = 10,000$) in each of which the digit has undergone a random displacement and rotation.
- The resulting images each have $D = 28 \times 28 = 784$ pixels.



MatLab Code



Off-line Digit Images: PCA Algorithm

- Our dataset is $X \in D \times N$, where $D = 784$ (pixels) and $N = 10,000$ (images).
- Each column of the matrix X is a data point (image) x_n , ($n = 1, 2, \dots, 10,000$) of size 784.
- We define the centered covariance matrix as:

$$S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T$$

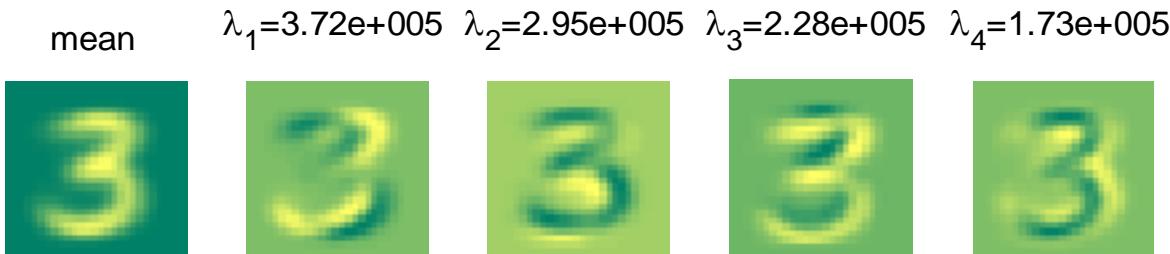
- To compute the principal components, we solve the following eigenvalue problem:

$$Su = \lambda u$$



Off-line Digit Images: PCA Results

- Using the algorithm stated before, we apply PCA on the off-line digits dataset.
- Since each eigenvector of S is a vector in the original D -dimensional space, we can represent the eigenvectors as images of the same size as the data points (see Fig. below)
- Below is the mean vector \bar{x} along with the first four PCA eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_4$ for the digit 3 from the off-line digits data set, together with the corresponding eigenvalues.



[MatLab Code](#)

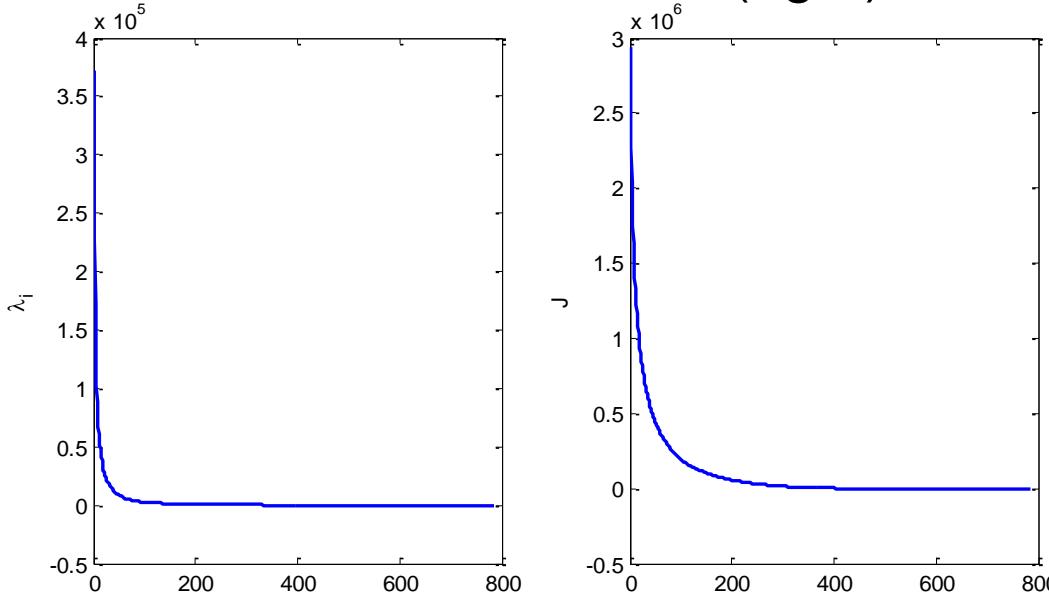


Off-line Digit Images: PCA Results

- The spectrum of eigenvalues, sorted into decreasing order, is shown below (left).
- The distortion measure J associated with choosing a particular value of M is given by

$$J = \sum_{i=M+1}^D \lambda_i$$

and is plotted for different values of M (right).



[MatLab Code](#)



Off-line Digit Images: PCA Reconstruction

- Below shows an original example of digit 3 from the off-line digits data set together with its PCA reconstructions obtained by retaining M principal components for various values of M .

$$\tilde{x}_n = \bar{x} + \sum_{i=1}^M (x_n^T u_i - \bar{x}^T u_i) u_i$$

- As M increases the reconstruction becomes more accurate and would become perfect when $M = D = 28 \times 28 = 784$.



[MatLab Code](#)



Off-line Digit Images: PCA Reconstruction



mean



principal basis 1



Using 5 bases



Using 10 bases



principal basis 2



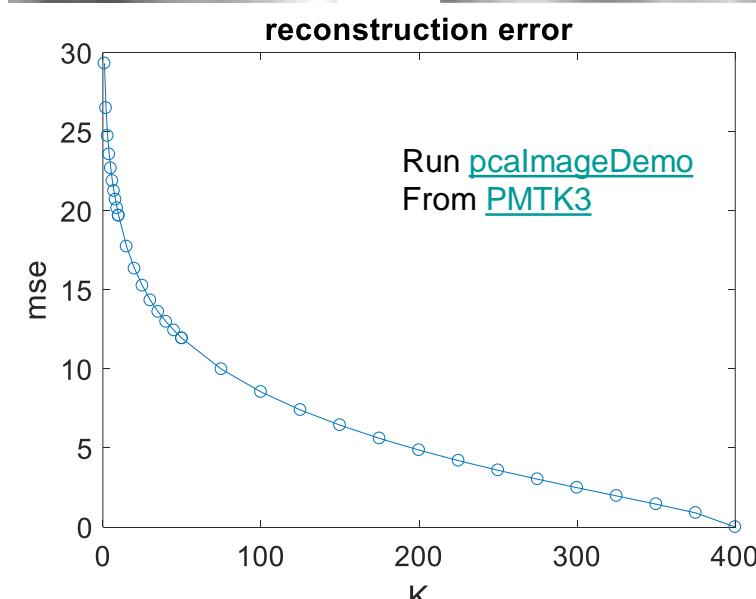
principal basis 3



Using 20 bases



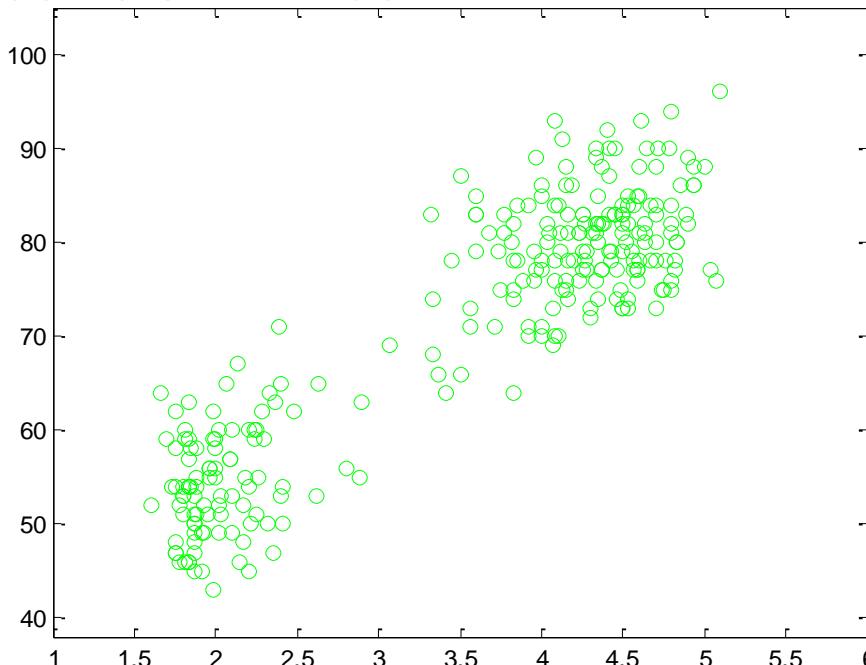
Using 400 bases



Run [pcalimageDemo](#)
From [PMTK3](#)

Old Faithful Data Set: Whitening

- ❑ Another example is the linear pre-processing of the old faithful dataset.
- ❑ Here the goal is not dimensionality reduction but rather the transformation of a data set in order to standardize the data (zero mean/unit variance).
- ❑ The original dataset is shown below:



MatLab Code



Old Faithful Data Set: Standardizing

- Standardizing the data is useful for successful application of pattern recognition techniques.
- It is essential when the original variables are measured in various different units or have different scaling.
- In the Old Faithful data set, the time between eruptions is typically an order of magnitude greater than the duration of an eruption.
- Before applying classification techniques (e.g. K-Means) we perform linear re-scaling of the individual variables such that **each variable has zero mean and unit variance**.
- This is known as **standardizing the data**.



Old Faithful Data Set: Standardizing

- The covariance matrix for the standardized data has components

$$\rho_{ij} = \frac{1}{N} \sum_{n=1}^N \frac{x_{ni} - \bar{x}_i}{\sigma_i} \frac{x_{nj} - \bar{x}_j}{\sigma_j}, \text{ where } \sigma_i = \text{std of } x_i$$

- This is the correlation matrix of the original data.
- If two components x_i and x_j of the data are perfectly correlated, then $\rho_{ij} = 1$, and if they are uncorrelated, then $\rho_{ij} = 0$.



Old Faithful Data Set: Whitening

- With PCA we can normalize the data to give them zero mean and **unit covariance (different variables become decorrelated)**.
- Consider the key eigenvalue problem in PCA in a matrix form:

$$\mathbf{S}\mathbf{U} = \mathbf{U}\mathbf{L}, \quad \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_D), \quad \mathbf{U} = [\mathbf{u}_1 \dots \mathbf{u}_D] \text{ (orthogonal)}$$

- For each data point \mathbf{x}_n , define a transformed value as:

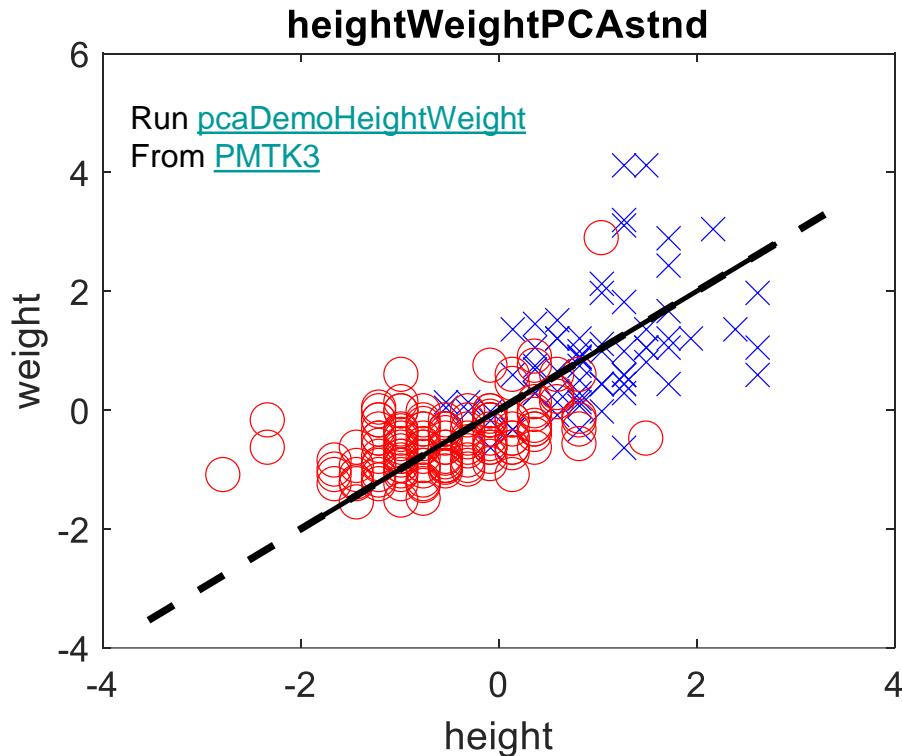
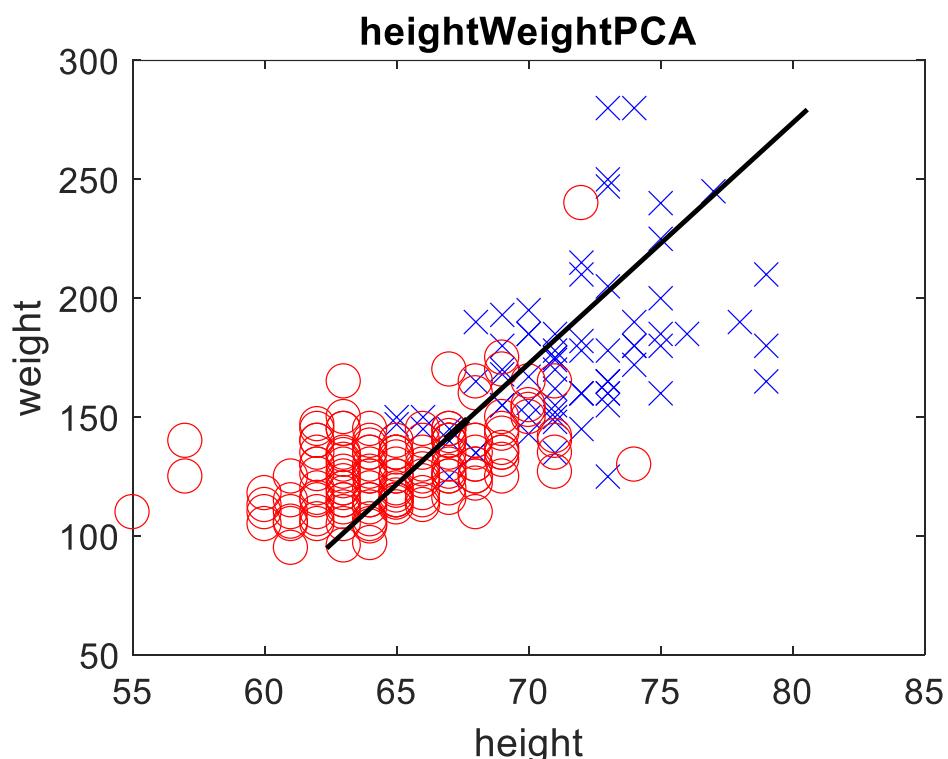
Whitening of the data: $\mathbf{y}_n = \mathbf{L}^{-1/2} \mathbf{U}^T (\mathbf{x}_n - \bar{\mathbf{x}})$

- The set \mathbf{y}_n has zero mean and its covariance is the identity:

$$\frac{1}{N} \sum_{n=1}^N \mathbf{y}_n \mathbf{y}_n^T = \frac{1}{N} \sum_{n=1}^N \mathbf{L}^{-1/2} \mathbf{U}^T (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{U} \mathbf{L}^{-1/2} = \mathbf{L}^{-1/2} \mathbf{U}^T \mathbf{S} \mathbf{U} \mathbf{L}^{-1/2} = \mathbf{I}$$

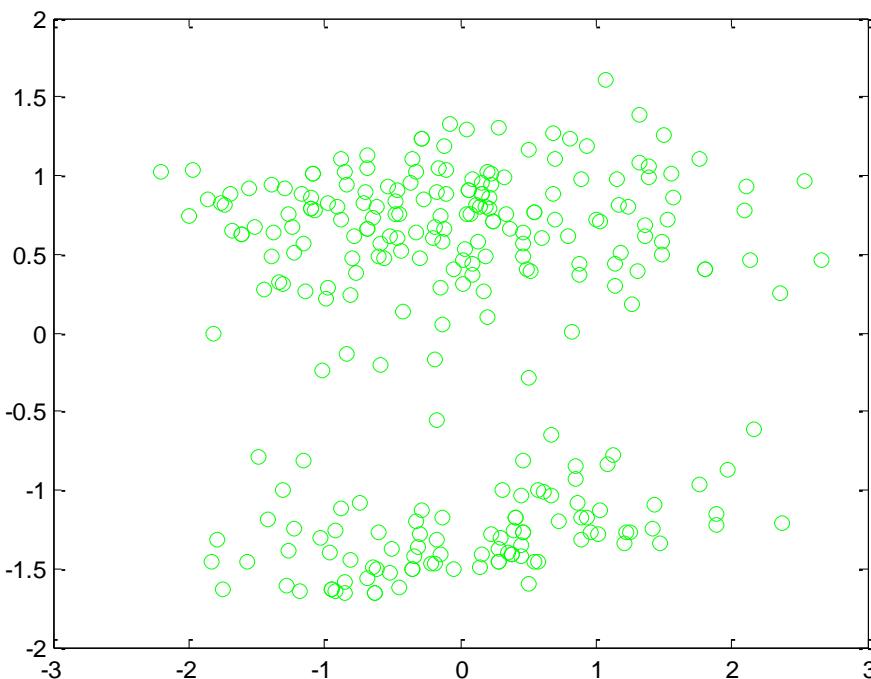
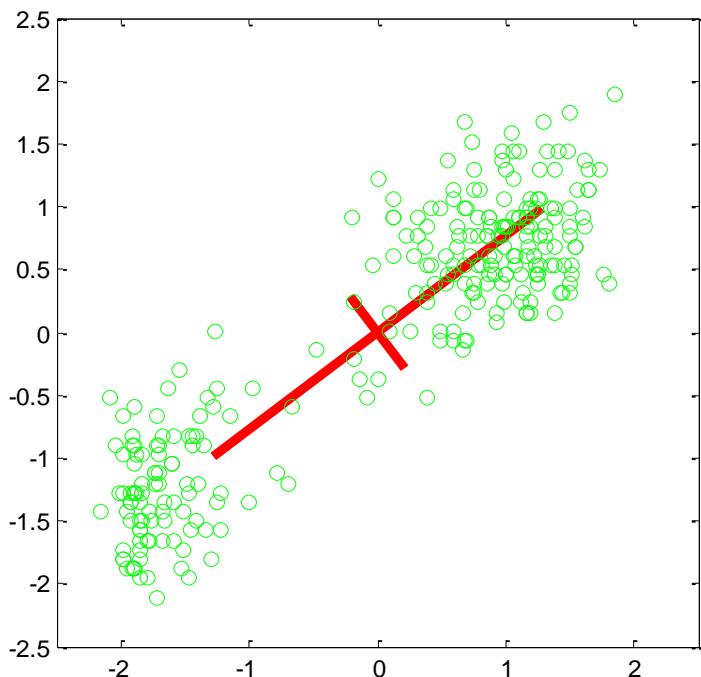
Old Faithful Data Set: PCA Example

- It should be standard practice to standardize the data first before performing PCA. This is equivalent to working with correlation matrices instead of covariance matrices.
- (Left) PCA for raw data (Right) PCA of standardized data.



Old Faithful Data Set: PCA Example

- Left: Standardizing individual variables to zero mean and unit variance. The principal axes of the normalized set are shown for the range $\pm\lambda_i^{1/2}$ (variables still correlated)
- Right: Whitening of the data (zero mean, unit covariance)

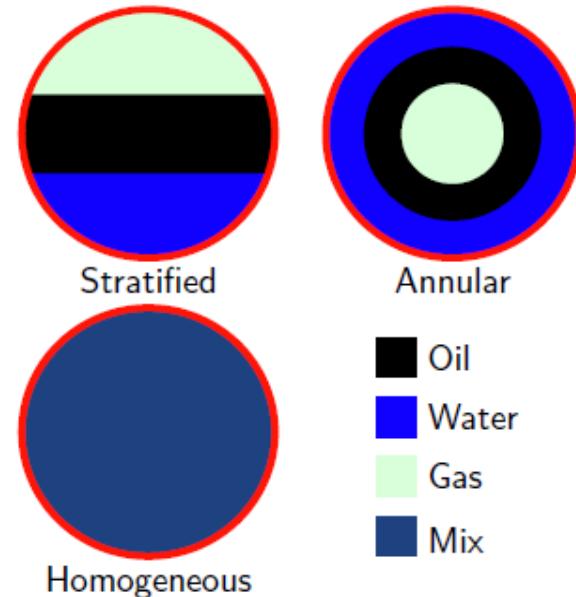


MatLab Code



Oil Flow Data Set

- Consider the oil flow data set.
- For a given geometrical configuration of the gas, water, and oil phases) there are only 2 DOF of variability -- **the fraction of oil in the pipe** and **the fraction of water** (the fraction of gas then being determined).
- The data space comprises 12 measurements.
- However, **the points in the data set lie close to a 2D manifold embedded within this 12th dimensional space.**

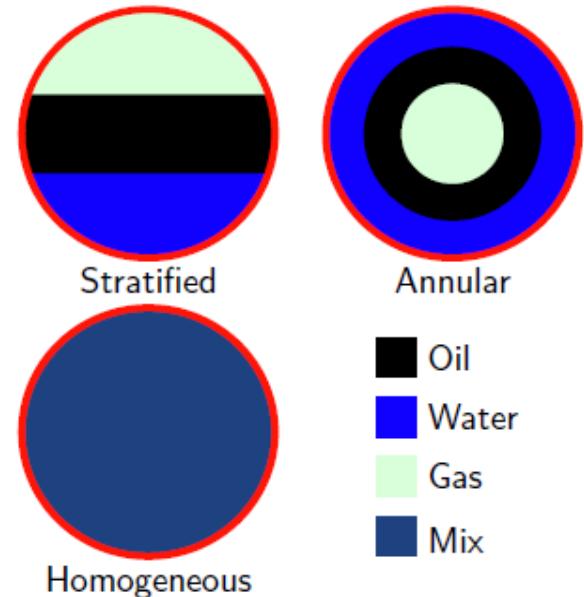


Flow Configurations

Red : 'laminar'
Blue: 'homogeneous'
Green: 'annular'

Oil Flow Data Set

- The manifold comprises several distinct segments corresponding to the different flow regimes, each such segment being a (noisy) continuous 2D manifold.
- For data compression, there is merit in exploiting this manifold structure.



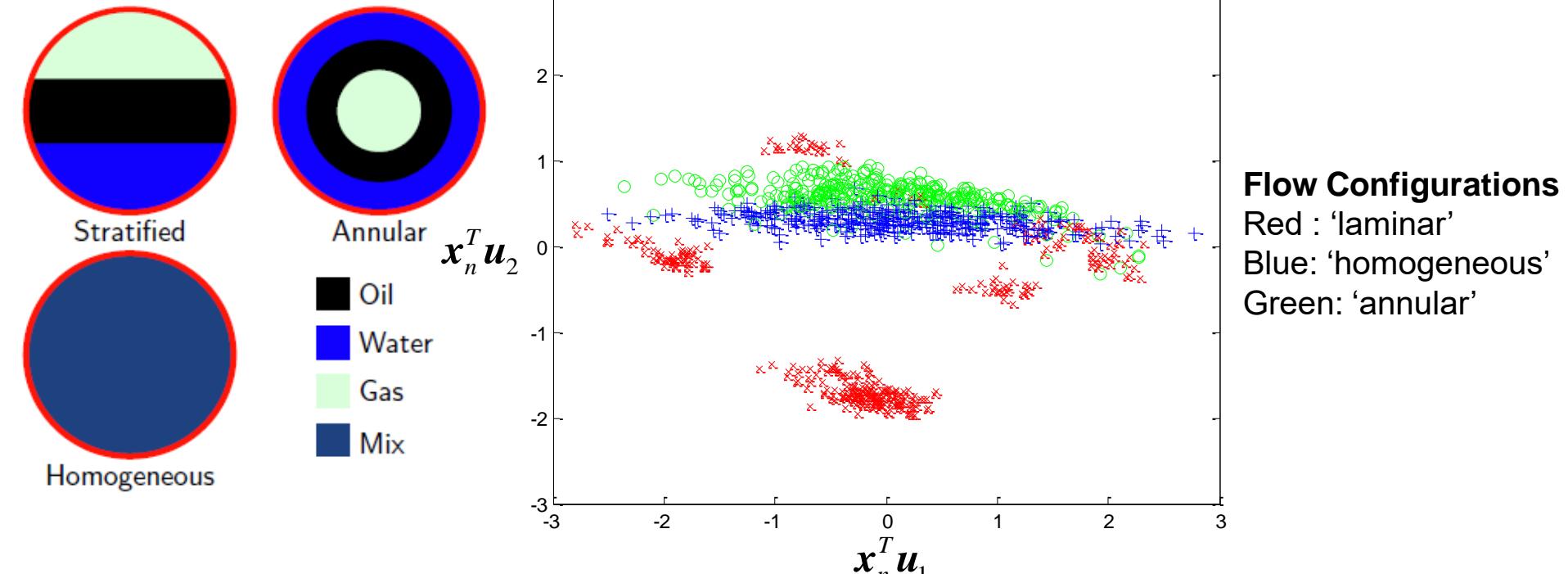
Flow Configurations

Red : 'laminar'
Blue: 'homogeneous'
Green: 'annular'



Oil Flow Data Set: PCA for Visualization

- Visualization of the oil flow data set obtained by projecting the data onto the first two principal components.

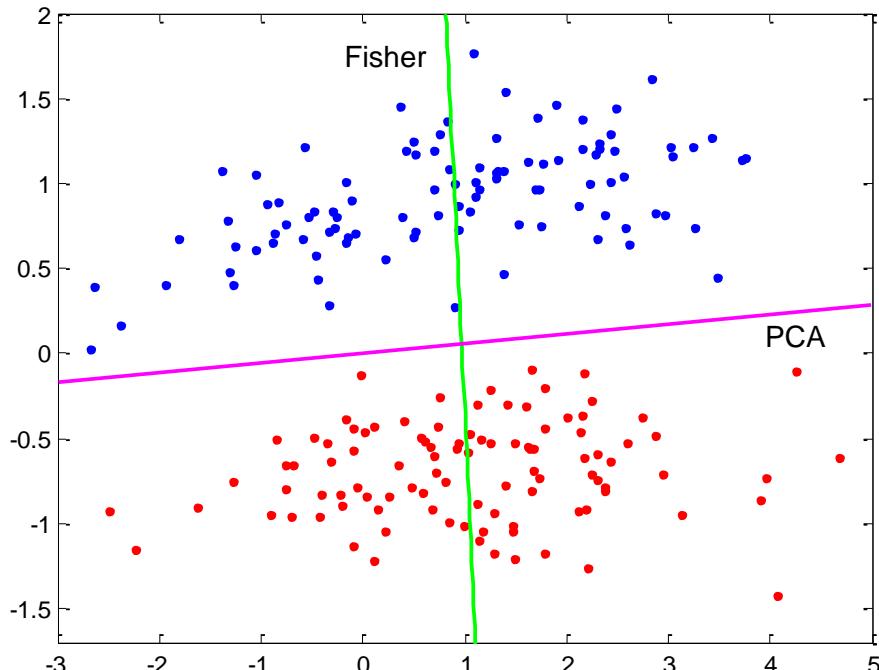


[MatLab Code](#)



PCA Vs. Fisher's Discriminant Method

- PCA and the Fisher linear discriminant method can be seen as linear dimensionality reduction techniques.
- PCA is unsupervised and depends only on x_n whereas Fisher linear discriminant also uses class-label information.



MatLab Code

- 2D data, 2 classes
- We are interested projecting the data in 1D.
- PCA chooses the direction of maximum variance (magenta curve) leading to strong class overlap.
- The Fisher linear discriminant projects (green curve) to maximize class separation.

PCA For High Dimensional Data

- In some applications of PCA, the number of data points (N) is smaller than the dimensionality (D) of the data space (e.g. 100 images each with 100,000 pixels).
- N points in a D -dimensional space, where $N \ll D$, defines a linear subspace whose dimensionality is at most $N - 1$.
 - There is little point in applying PCA for values of M that are greater than $N - 1$.
- If we perform PCA, we will find that **at least $D - N + 1$** of the eigenvalues are zero, corresponding to eigenvectors along whose directions the data set has zero variance.

PCA For High Dimensional Data

- Typical algorithms for finding the eigenvectors of a $D \times D$ matrix have a computational cost that scales like $\mathcal{O}(D^3)$.
- Direct application of PCA in e.g. the image example will be computationally infeasible.



PCA For High Dimensional Data

- To resolve the problem, let \mathbf{X} to be the $N \times D$ -dimensional centered data matrix, whose n^{th} row is given by $(\mathbf{x}_n - \bar{\mathbf{x}})^T$.

- The covariance $D \times D$ matrix can then be written as:

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T = \frac{1}{N} \mathbf{X}^T \mathbf{X}$$

- The corresponding eigenvector equation is:

$$\frac{1}{N} \mathbf{X}^T \mathbf{X} \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

- Pre-multiplying both sides from the left by \mathbf{X} gives:

$$\frac{1}{N} \mathbf{X} \mathbf{X}^T (\mathbf{X} \mathbf{u}_i) = \lambda_i (\mathbf{X} \mathbf{u}_i)$$



PCA For High Dimensional Data

$$\frac{1}{N} \mathbf{X} \mathbf{X}^T (\mathbf{X} \mathbf{u}_i) = \lambda_i (\mathbf{X} \mathbf{u}_i)$$

- Let us define:

$$\mathbf{v}_i = \mathbf{X} \mathbf{u}_i$$

- The eigenvalue problem becomes:

$$\frac{1}{N} \mathbf{X} \mathbf{X}^T \mathbf{v}_i = \lambda_i \mathbf{v}_i$$

- This is an eigenvector equ. for the $N \times N$ matrix $\mathbf{X} \mathbf{X}^T / N$
- This has the same $N - 1$ eigenvalues as the original covariance matrix (which has an additional $D - N + 1$ zero eigenvalues). Note we have $N - 1$ (and not N) eigenvalues because the data are centered ($\mathbf{X} \mathbf{X}^T$ has rank $N - 1$).
- Thus we can solve the eigenvector problem in spaces of lower dimensionality with computational cost $\mathcal{O}(N^3)$.



PCA For High Dimensional Data

- To determine the eigenvectors \mathbf{u}_i , we multiply by \mathbf{X}^T :

$$\frac{1}{N} (\mathbf{X}^T \mathbf{X}) \mathbf{X}^T \mathbf{v}_i = \lambda_i (\mathbf{X}^T \mathbf{v}_i)$$

- The covariance matrix recall has eigenvectors \mathbf{u}_i . Thus with proper rescaling (assuming \mathbf{v}_i is already normalized):

$$\mathbf{u}_i = \frac{1}{(N\lambda_i)^{1/2}} \mathbf{X}^T \mathbf{v}_i$$

- This approach is indeed simple:

- first evaluate $\mathbf{X}\mathbf{X}^T$ and then find its eigenvectors and eigenvalues and
- then compute the eigenvectors in the original data space from the equ. above.



Probabilistic PCA

PCA as the maximum likelihood solution of a probabilistic latent variable model (Tipping & Bishop [1997](#), [1999](#), Roweis, [1998](#))

- Constrained form of the Gaussian distribution: the number of free parameters is restricted while the model still captures dominant correlations in the data
- Derive an efficient EM algorithm for PCA
- Allows dealing with missing values in the dataset
- Mixture of probabilistic PCA models

- Tipping, M. E. and C. M. Bishop (1997). [Probabilistic principal component analysis. Technical Report NCRG/97/010, Neural Computing Research Group, Aston University.](#)
- Tipping, M. E. and C. M. Bishop (1999b). [Probabilistic principal component analysis. Journal of the Royal Statistical Society, Series B 21\(3\), 611–622.](#)
- Roweis, S. (1998). [EM algorithms for PCA and SPCA.](#) In M. I. Jordan, M. J. Kearns, and S. A. Solla (Eds.), [Advances in Neural Information Processing Systems, Volume 10, pp. 626–632. MIT Press.](#)



Probabilistic PCA

PCA as the maximum likelihood solution of a probabilistic latent variable model (Tipping & Bishop [1997](#), [1999](#), Roweis, [1998](#))

- Bayesian treatment of PCA -- the dimensionality of the principal subspace can be found automatically
- Probabilistic PCA can model class conditional densities and thus be used for classification
- Can be used generatively: samples from the distribution

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



Probabilistic PCA - Model

- Introduce a latent variable z corresponding to the principal-component subspace.
- Define a Gaussian prior distribution $p(z)$, together with a Gaussian conditional distribution $p(x|z)$ for the observed variable x conditioned on the value of z .
- The prior distribution over z is given by

$$p(z) = \mathcal{N}(z|\mathbf{0}, I)$$

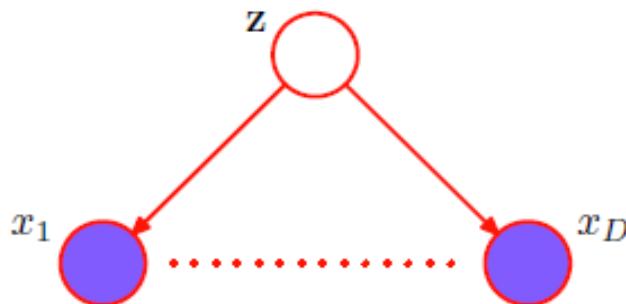
- The conditional distribution of the observed variable x , conditioned on the value of z , is again Gaussian:

$$p(x|z) = \mathcal{N}(x|Wz + \mu, \sigma^2 I)$$

The mean of x is a linear function of z governed by the $D \times M$ matrix W and the D -dimensional vector μ .

Probabilistic PCA: Naïve Bayes Model

- This model $p(x|z) = \mathcal{N}(x|Wz + \mu, \sigma^2 I)$ is an example of Naive Bayes' model - Conditioned on z , the components of the observed vector $x = (x_1, \dots, x_D)^T$ are assumed to be independent.



- We have 2 parameters in the model:
 - (a) W , its columns span a linear subspace within the data space that corresponds to the principal subspace;
 - (b) σ^2 – variance of the conditional distribution.
- The prior $p(z) = \mathcal{N}(z|\mathbf{0}, I)$ is without loss of generality.

Probabilistic PCA - Generative View Point

- A latent variable model seeks to relate a D -dimensional observation vector x to a corresponding M -dimensional Gaussian latent variable z

$$x = Wz + \mu + \varepsilon$$

where

- z is an M -dimensional Gaussian latent variable
- W is an $(D \times M)$ matrix (the latent space)
- ε is a D -dimensional Gaussian noise
- ε and z are independent
- μ is a parameter vector (non zero mean)

- Factor analysis:

$$\varepsilon \sim \mathcal{N}(\mathbf{0}, \psi)$$

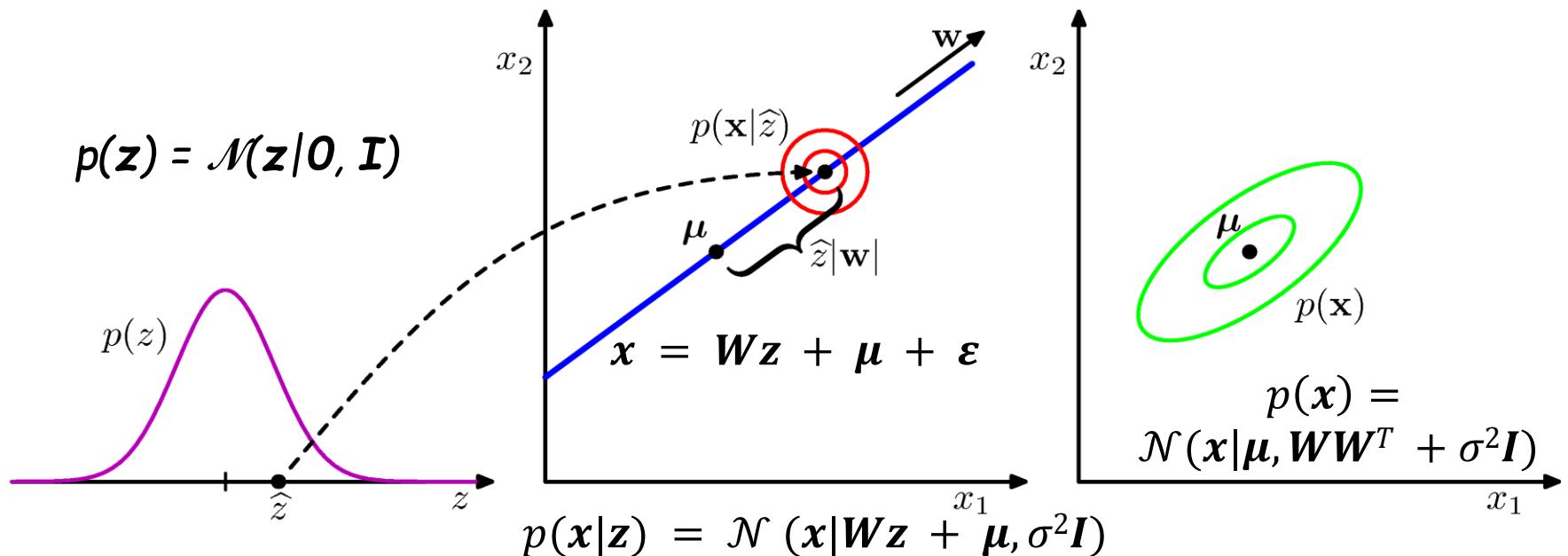
- Probabilistic PCA:

$$\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$$



Probabilistic PCA - Generative View Point

- Mapping from the latent space to the data space.
- Assume here 2D data and 1D latent space.
- An observed x is generated by drawing a value \hat{z} from $p(z)$ & then a value for x from an isotropic Gaussian distribution (red circles) having mean $w\hat{z} + \mu$ and covariance $\sigma^2 I$. The green ellipses are the density contours of $p(x)$.



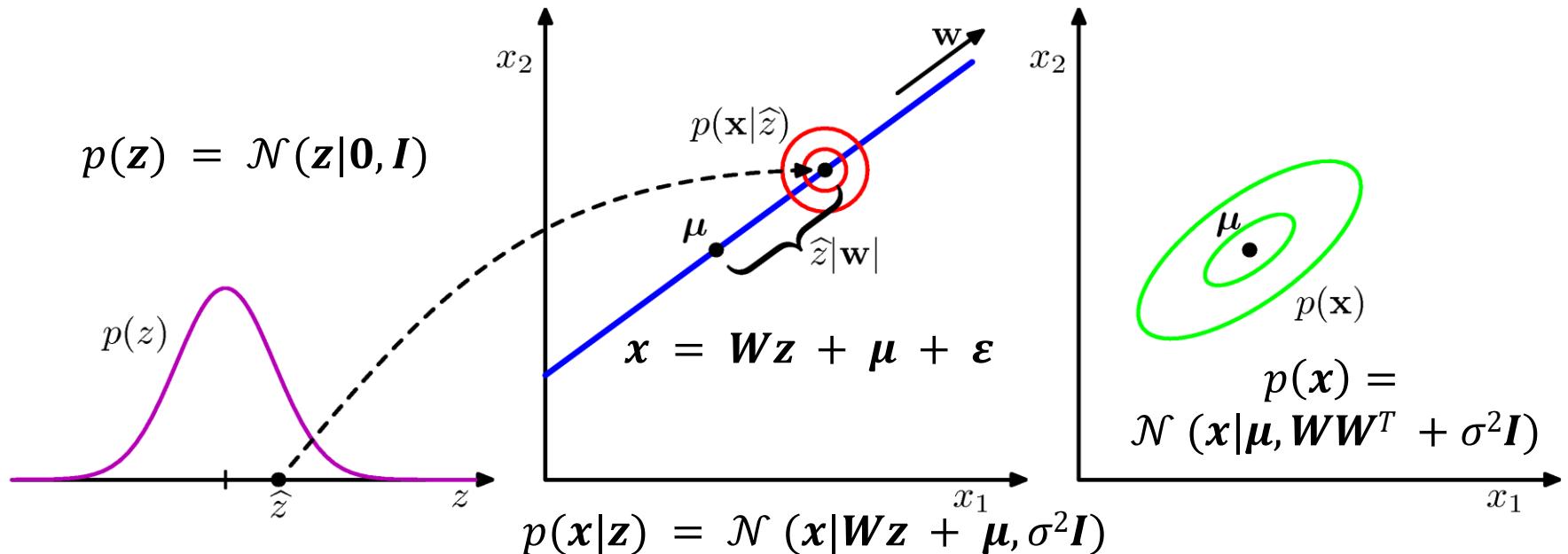
Probabilistic PCA - Generative View Point

- To compute the likelihood function, we need an expression for the marginal distribution $p(x)$ of the observed variable:

$$p(x) = \int p(x | z) p(z) dz$$

- Using Eqs. for linear Gaussian models, this is given as:

$$p(x) = \mathcal{N}(x | \mu, WW^T + \sigma^2 I)$$



Probabilistic PCA - Generative View Point

$$p(x) = \mathcal{N}(x | \mu, WW^T + \sigma^2 I)$$

- This result can be derived directly by noting that the predictive distribution will be Gaussian and then evaluating its mean and covariance.

$$\mathbb{E}[x] = \mathbb{E}[Wz + \mu + \varepsilon] = W\mathbb{E}[z] + \mu + \mathbb{E}[\varepsilon] = \mu$$

$$\text{var}[x] = \mathbb{E}\left[(Wz + \mu + \varepsilon)(Wz + \mu + \varepsilon)^T\right]$$

$$= WW^T \text{var}[z] + \mathbb{E}\left[\varepsilon\varepsilon^T\right] = WW^T + \sigma^2 I$$

- We used here that z and ε are independent, thus uncorrelated.

The Assumed Prior $p(z)$ is Not Restrictive

- The prior we assumed is not restrictive. Indeed consider a more general prior of the form:

$$p(z) = \mathcal{N}(z | \mathbf{m}, \Sigma)$$

- Then the marginal distribution is of the form:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{W}\mathbf{m} + \boldsymbol{\mu}, \mathbf{W}\boldsymbol{\Sigma}^{-1}\mathbf{W}^T + \sigma^2\mathbf{I})$$

- Thus the marginal distribution has an identical form as before:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \widetilde{\boldsymbol{\mu}}, \widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^T + \sigma^2\mathbf{I}), \quad \begin{aligned} \widetilde{\boldsymbol{\mu}} &= \mathbf{W}\mathbf{m} + \boldsymbol{\mu} \\ \widetilde{\mathbf{W}} &= \boldsymbol{\Sigma}^{-1/2}\mathbf{W} \end{aligned}$$

- We choose the simplest prior: $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$



Probabilistic PCA - Generative View Point

- There is redundancy in this parametrization.
- There is a whole class of \mathbf{W} 's differing by a rotation of the latent space coordinates that leads to the same predictive distribution.
Indeed:

- Let

$$\widetilde{\mathbf{W}} = \mathbf{W}\mathbf{R}, \text{ where } \mathbf{R} = \text{orthogonal matrix}$$

- Then

$$\widetilde{\mathbf{W}}\widetilde{\mathbf{W}}^T = \mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T = \mathbf{W}\mathbf{W}^T$$



Probabilistic PCA - Predictive Distribution

$$p(x) = \mathcal{N}(x | \mu, C) = \mathcal{N}(x | \mu, WW^T + \sigma^2 I)$$

- To compute the predictive distribution, we need to be able to invert C ($D \times D$ matrix). We use the matrix inversion Lemma:

$$C^{-1} = (WW^T + \sigma^2 I)^{-1} = \sigma^{-2} I - \sigma^{-2} W M^{-1} W^T$$

where

$$M = W^T W + \sigma^2 I \text{ } (M \times M \text{ matrix})$$

- The cost of this inversion is reduced from $\mathcal{O}(D^3)$ to $\mathcal{O}(M^3)$!

Probabilistic PCA - Predictive Distribution

$$\mathbf{C}^{-1} = (\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})^{-1} = \sigma^{-2} \mathbf{I} - \sigma^{-2} \mathbf{W} \mathbf{M}^{-1} \mathbf{W}^T$$

$$\mathbf{M} = \mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I} \text{ } (M \times M \text{ matrix})$$

- To derive this we used the [Woodbury identity](#):

$$(\mathbf{A} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}$$

- This identity (easy to show with direct substitution) is used often when

- \mathbf{A} is large and diagonal (so easy to invert),
- \mathbf{B} has many rows but few columns $\boxed{\mathbf{B}}$ (and conversely for \mathbf{C} $\boxed{\mathbf{C}}$) so that the rhs is much cheaper to evaluate than the lhs.



Probabilistic PCA - Posterior Distribution

- The posterior distribution can be derived directly from earlier results on linear Gaussian models.
 - We know $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$
 - The conditional: $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$
 - From these we conclude:

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}\left(\mathbf{z} | \mathbf{M}^{-1}\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}), \sigma^2\mathbf{M}^{-1}\right), \quad \mathbf{M} = \mathbf{W}^T\mathbf{W} + \sigma^2\mathbf{I}$$

Here we used results from an earlier lecture:

$$\begin{cases} p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ p(\mathbf{y} | \mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \end{cases} \Rightarrow$$

$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}\left(\mathbf{x} | \boldsymbol{\Sigma}(\boldsymbol{\Lambda}\boldsymbol{\mu} + \mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b})), \boldsymbol{\Sigma}\right),$$
$$\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1}$$



Maximum Likelihood PCA

- Consider determining the model parameters using maximum likelihood. Using

$$p(\mathbf{x}) = \mathcal{N}\left(\mathbf{x} \mid \boldsymbol{\mu}, \underbrace{\mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}}_C\right)$$

we derive:

$$\begin{aligned} \ln p(\mathbf{X} \mid \boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \sum_{n=1}^N \ln p(\mathbf{x}_n \mid \boldsymbol{\mu}, \mathbf{W}, \sigma^2) \\ &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \end{aligned}$$

- Setting the derivative wrt $\boldsymbol{\mu}$ equal to zero gives:

$$\boldsymbol{\mu} = \sum_{n=1}^N \mathbf{x}_n \Big/ N \equiv \bar{\mathbf{x}}$$

Maximum Likelihood PCA

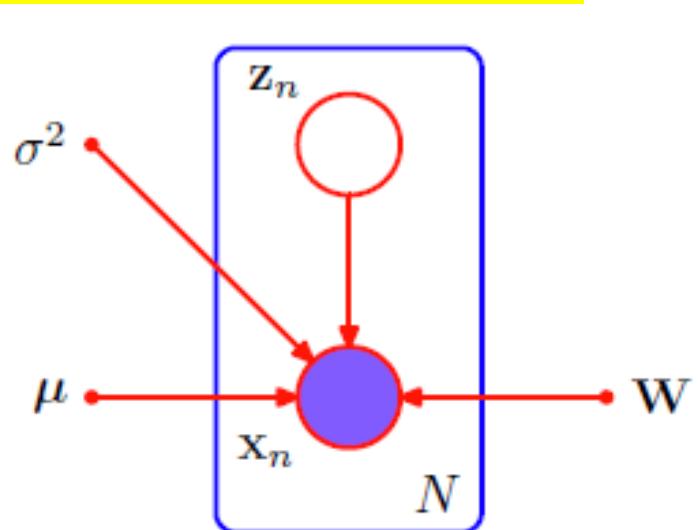
□ The log-likelihood is then simplified as:

$$\ln p(X | \mu, W, \sigma^2) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|C| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})^T C^{-1} (\mathbf{x}_n - \bar{\mathbf{x}})$$

or

$$\ln p(X | \mu, W, \sigma^2) = -\frac{N}{2} \left\{ D \ln(2\pi) + \ln|C| + \text{Tr}(C^{-1} S) \right\},$$

$$S = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$



□ Maximization wrt W and σ^2 can also be done analytically:

$$W_{ML} = U_M \left(L_M - \sigma^2 I \right)^{1/2} R$$

- Tipping, M. E. and C. M. Bishop (1999b). [Probabilistic principal component analysis. Journal of the Royal Statistical Society, Series B 21\(3\), 611–622.](#)
- Roweis, S. (1998). [EM algorithms for PCA and SPCA](#). In M. I. Jordan, M. J. Kearns, and S. A. Solla (Eds.), [Advances in Neural Information Processing Systems, Volume 10, pp. 626–632. MIT Press](#).

Maximum Likelihood PCA

$$\mathbf{W}_{ML} = \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{R}$$

- \mathbf{U}_M is a $D \times M$ matrix whose columns are given by any subset (of size M) of the eigenvectors of the data covariance matrix \mathbf{S} .
- \mathbf{L}_M is the $M \times M$ diagonal matrix with elements given by the corresponding eigenvalues λ_i of \mathbf{S} .
- \mathbf{R} is an arbitrary $M \times M$ orthogonal matrix.
- The max of the likelihood function is obtained when the M eigenvectors are chosen to be those whose eigenvalues are the M largest (all other solutions being saddle points).
- For eigenvectors arranged in order of decreasing λ_i values, the M principal eigenvectors are $\mathbf{u}_1, \dots, \mathbf{u}_M$. The columns of \mathbf{W} then define the principal subspace as in standard PCA.

Maximum Likelihood PCA

- The corresponding MLE solution for σ^2 is given as:

$$\sigma_{ML}^2 = \frac{1}{D-M} \sum_{i=M+1}^N \lambda_i$$

i.e. the average of the discarded eigenvalues.

- R in $W_{ML} = U_M (L_M - \sigma^2 I)^{1/2} R$ is a rotation matrix in the M dimensional latent space.
- Substituting this into the predictive variance $C = W_{ML} W_{ML}^T + \sigma_{ML}^2 I$ and using $R^T R = I$, we see that C is independent of R .
- The predictive density is unchanged by rotations in the latent space (statistical non-identifiability).



Maximum Likelihood PCA

$$\mathbf{W}_{ML} = \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{R}$$

- For $\mathbf{R} = \mathbf{I}$, the columns of \mathbf{W} are the principal component eigenvectors scaled by $(\lambda_i - \sigma^2)^{1/2}$.
- Interpretation: for the convolution of independent Gaussian distributions, here of \mathbf{z} and the noise $\boldsymbol{\varepsilon}$, the variances are additive.

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}$$

- The variance λ_i in the direction of an eigenvector \mathbf{u}_i is composed of the sum of a contribution $\lambda_i - \sigma^2$ from the projection of the unit-variance latent space $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ into data space through the corresponding column of \mathbf{W} , plus an isotropic contribution of variance σ^2 added in all directions by the noise model $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$.



Maximum Likelihood PCA

- Consider the variance of the predictive distribution $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$ along some direction defined by the unit vector $\boldsymbol{\nu}$ given by $\boldsymbol{\nu}^T \mathbf{C} \boldsymbol{\nu}$.
- Let $\boldsymbol{\nu}$ be orthogonal to the principal subspace \mathbf{U} , i.e. $\boldsymbol{\nu}$ is given by some linear combination of the discarded eigenvectors.
- Then $\boldsymbol{\nu}^T \mathbf{U} = 0$ and hence

$$\boldsymbol{\nu}^T \mathbf{C} \boldsymbol{\nu} = \boldsymbol{\nu}^T (\mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}) \boldsymbol{\nu} = \sigma^2$$

- Thus the model predicts a noise variance orthogonal to the principal subspace which it was shown to be the average of the discarded eigenvalues.

$$\sigma_{ML}^2 = \frac{1}{D-M} \sum_{i=M+1}^N \lambda_i$$



Maximum Likelihood PCA

- Consider the variance of the predictive distribution $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$ along some direction defined by the unit vector \mathbf{v} given by $\mathbf{v}^T \mathbf{C} \mathbf{v}$.
- Let now $\mathbf{v} = \mathbf{u}_i$ where \mathbf{u}_i is one of the retained eigenvectors defining the principal subspace.
- Then using $\mathbf{W}_{ML} = \mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}$, and $\mathbf{u}_i^T \mathbf{u}_j = 0, j = 1, \dots, M$ we see that:

$$\mathbf{v}^T \mathbf{C} \mathbf{v} = \mathbf{u}_i^T (\mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}) \mathbf{u}_i = (\lambda_i - \sigma^2) + \sigma^2 = \lambda_i$$

i.e. the model correctly captures the variance of the data along the principal axes.

- As shown earlier, the variance in all remaining directions is approximated with a single value of σ^2 . Variance is ‘lost’ in the projections.



Maximum Likelihood PCA

$$\mathbf{W}_{ML} = \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{R}$$

$$\sigma_{ML}^2 = \frac{1}{D-M} \sum_{i=M+1}^N \lambda_i$$

- We construct the MLE model by solving the underlying eigenvalue problem for the data covariance matrix and evaluate \mathbf{W} and σ^2 from the Eqs. above. We choose $\mathbf{R} = \mathbf{I}$.
- Note: if the MLE model is found using optimization methods or via the EM algorithm (see following slides), \mathbf{R} is arbitrary and the columns of \mathbf{W} not orthogonal.

In this case, orthogonality can be enforced

- as postprocessing or
 - by modifying the EM algorithm.
-
- Ahn, J. H. and J. H. Oh (2003). [A constrained EM algorithm for principal component analysis](#). *Neural Computation* 15(1), 57–65.



Maximum Likelihood PCA

- If we consider no dimensionality reduction ($M = D$), then

$$\mathbf{U}_M = \mathbf{U}, \quad \mathbf{L}_M = \mathbf{L},$$

and using $\mathbf{U}\mathbf{U}^T = \mathbf{I}$, $\mathbf{R}\mathbf{R}^T = \mathbf{I}$, we see that the covariance \mathbf{C} of the marginal distribution for \mathbf{x} becomes

$$\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I} = \mathbf{U}\left(\mathbf{L} - \sigma^2\mathbf{I}\right)^{1/2}\mathbf{R}\mathbf{R}^T\left(\mathbf{L} - \sigma^2\mathbf{I}\right)^{1/2}\mathbf{U}^T + \sigma^2\mathbf{I} = \mathbf{U}\mathbf{L}\mathbf{U}^T = \mathbf{S}$$

- This result is the standard MLE solution for an unconstrained Gaussian distribution in which the covariance matrix is given by the sample covariance.

Maximum Likelihood PCA

- Since \mathbf{C} is not full rank, we can use the matrix inversion lemma:

$$\mathbf{C}^{-1} = \frac{1}{\sigma^2} \left[\mathbf{I} - \mathbf{W} \left(\mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I} \right)^{-1} \mathbf{W}^T \right]$$

- Plugging the MLE estimates for \mathbf{W}, σ^2 gives: $\mathbf{W} = \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2}$

$$\begin{aligned}\mathbf{C}^{-1} &= \frac{1}{\sigma^2} \left[\mathbf{I} - \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \left(\mathbf{L}_M - \sigma^2 \mathbf{I} + \sigma^2 \mathbf{I} \right)^{-1} \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{U}_M^T \right] \\ &= \frac{1}{\sigma^2} \left[\mathbf{I} - \mathbf{U}_M \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{L}_M^{-1} \left(\mathbf{L}_M - \sigma^2 \mathbf{I} \right)^{1/2} \mathbf{U}_M^T \right] \\ &= \frac{1}{\sigma^2} \left[\mathbf{I} - \mathbf{U}_M \text{diag} \left(1 - \frac{\sigma^2}{\lambda_j} \right) \mathbf{U}_M^T \right]\end{aligned}$$

- Similarly: $\log |\mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}| = (D - M) \log \sigma^2 + \sum_{i=1}^M \log \lambda_i$

Use the [matrix inversion lemma](#) $\det(A + UWV^T) = \det(W^{-1} + V^T A^{-1} U) \det W \det A$

Maximum Likelihood PCA

- PCA is generally expressed as a projection of points from the D -dimensional dataspace onto an M -dimensional subspace.
- Probabilistic PCA is seen as the mapping $\mathbf{z} \rightarrow \mathbf{x}, \mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}$.
- The reverse mapping $\mathbf{x} \rightarrow \mathbf{z}$ is computed using the posterior:

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}\left(\mathbf{z} | \mathbf{M}^{-1}\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}), \sigma^2\mathbf{M}^{-1}\right), \quad \mathbf{M} = \mathbf{W}^T\mathbf{W} + \sigma^2\mathbf{I}$$

- Every point in data space is characterized by its posterior mean and covariance in latent space:

$$\mathbb{E}[z | x] = \mathbf{M}^{-1}\mathbf{W}_{ML}^T(x - \bar{x}), \quad \text{var}[z | x] = \sigma^2\mathbf{M}^{-1}$$

- This projects to a point in data space given by:

$$\mathbf{W}\mathbb{E}[z | x] + \boldsymbol{\mu}$$

- This is identical form as regularized linear regression.



Maximum Likelihood PCA vs Standard PCA

- In the limit $\sigma^2 \rightarrow 0$, the posterior mean becomes:

$$\mathbb{E}[z | x] = M^{-1} W_{ML}^T (x - \bar{x}), M = W_{ML}^T W_{ML} + \sigma^2 I \rightarrow W_{ML}^T W_{ML} \Rightarrow$$

$$\mathbb{E}[z | x] = (W_{ML}^T W_{ML})^{-1} W_{ML}^T (x - \bar{x})$$

- Now substitute $W_{ML} = U_M (L_M - \sigma^2 I)^{1/2} R = U_M L_M^{1/2}$ for $\sigma^2 \rightarrow 0$ with $R = I$ (for consistency with PCA). Then

$$\mathbb{E}[z | x] = (W_{ML}^T W_{ML})^{-1} W_{ML}^T (x - \bar{x}) = L_M^{-1/2} U_M^T (x - \bar{x})$$

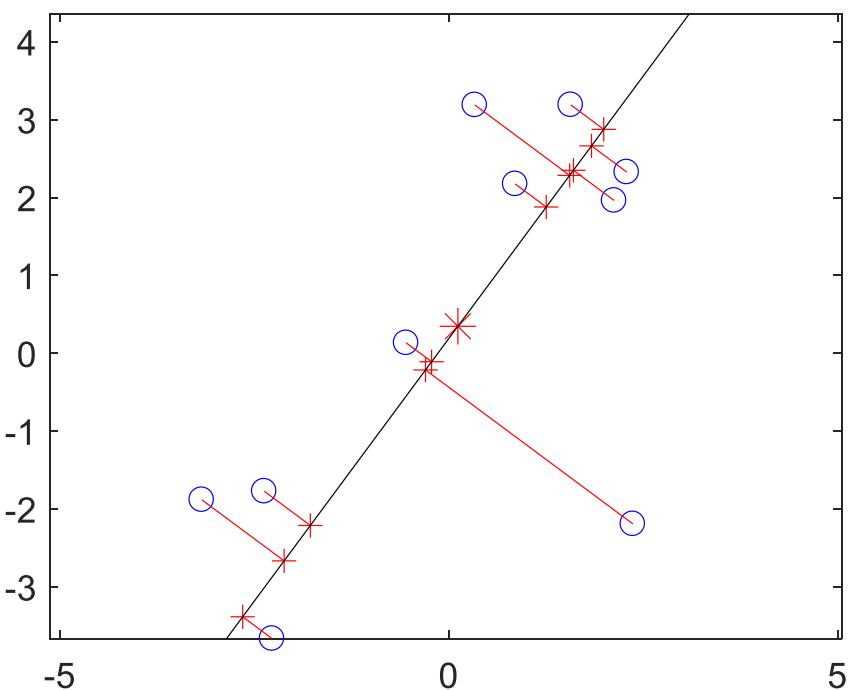
- This is an orthogonal projection of the data point onto the latent space, i.e. for the limit $\sigma^2 \rightarrow 0$, we recover the standard PCA model. The posterior covariance

$$\text{var}[z | x] = \sigma^2 M^{-1} \rightarrow 0 \text{ and the density becomes singular.}$$

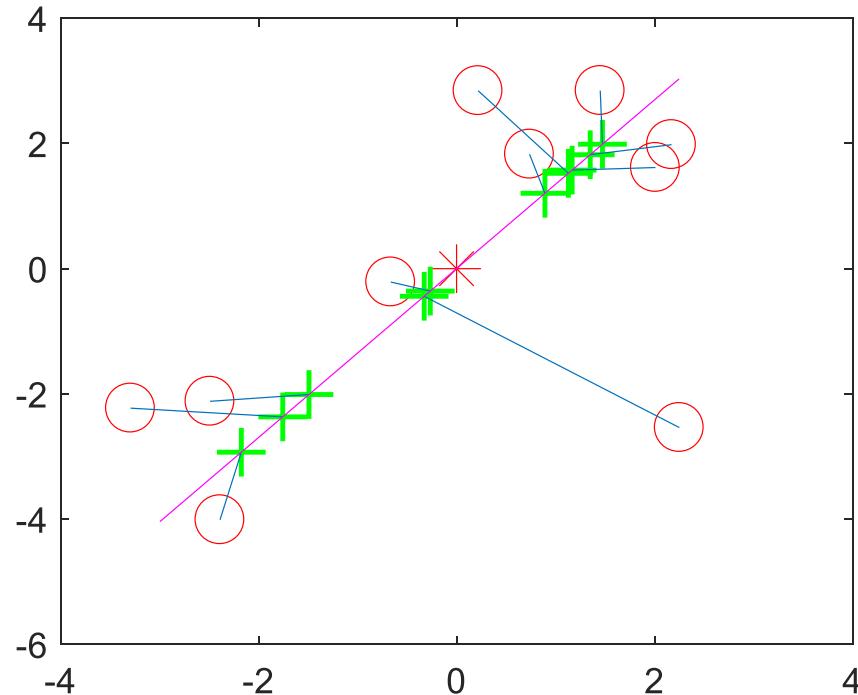
- For $\sigma^2 > 0$, the latent projection is shifted towards the origin, relative to the orthogonal projection.

PPCA Versus PCA

- Consider PCA and PPCA where $D = 2$ and $M = 1$. The red star is the data mean. In PCA the points are orthogonally projected onto the line. In PPCA the projection is no longer orthogonal and the reconstructions are shrunk towards the data mean (red star).



Run [pcaDemo2d](#)
From [PMTK3](#)



Run [ppcaDemo2d](#)
From [PMTK3](#)

Probabilistic PCA: Number of DOF

- PPCA defines a multivariate Gaussian distribution in which the number of DOF (independent parameters) can be controlled while still capturing the dominant correlations in the data.
- A general Gaussian distribution has $D(D + 1)/2$ independent parameters in its covariance matrix and D parameters in its mean.
 - The number of parameters scales as D^2 .
- For a diagonal covariance matrix we have D independent parameters and the number of parameters grows linearly. However, the variables are independent and hence this model cannot express any correlations between them.

Probabilistic PCA: Number of DOF

- In PPCA, the M most significant correlations are captured while the total number of parameters grows linearly with D .
- We can see this by evaluating the DOF in PPCA:
 - The covariance \mathbf{C} depends on \mathbf{W} ($D \times M$), and σ^2 : $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$ total parameters $DM + 1$.
 - We need to subtract the redundancy associated with rotations of the coordinate system in the latent space.
 - \mathbf{R} is $M \times M$. In the 1st column there are $M - 1$ independent parameters (must be normalized). In the 2nd column there are $M - 2$ independent parameters (normalized & orthogonal to the 1st column), etc. \mathbf{R} has a total of $M(M - 1)/2$ independent parameters.
- The number of degrees of freedom in \mathbf{C} grows linearly with D

$$D \times M + 1 - M \times (M - 1)/2$$



Probabilistic PCA: Number of DOF

- The number of DOF in \mathcal{C} grows linearly in D for a given M :

$$D \times M + 1 - M \times (M - 1)/2$$

- $M = D - 1$: recover the standard result for a full covariance Gaussian.
 - The variance along $D - 1$ linearly independent directions is controlled by the columns of \mathbf{W} , and
 - the variance along the remaining direction is given by σ^2 .
- $M = 0$: equivalent to the isotropic covariance case.

EM Algorithms for PCA

- The PPCA model involves a marginalization over a continuous latent space \mathbf{z} . For each \mathbf{x}_n there is a corresponding \mathbf{z}_n .
- Use EM to find the MLE of the model parameters.
- For high D , there are advantages using iteratively EM vs. working directly with the sample covariance.
- This approach can be extended to factor analysis for which there is no closed-form solution.
- Can also be used when values are missing & for mixture models.
- EM requires the complete-data log likelihood function:

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{n=1}^N \left\{ \ln p(\mathbf{x}_n | \mathbf{z}_n) + \ln p(\mathbf{z}_n) \right\}$$

- In the following, we substitute $\boldsymbol{\mu}$ with the sample mean $\bar{\mathbf{x}}$.

EM Algorithms for PCA

- We first take the expectation of the complete-data log likelihood with respect to the posterior distribution of the latent distribution evaluated using ‘old’ parameter values.
- Maximization of this expected complete data log likelihood then yields the ‘new’ parameter values.

$$\ln p(X, Z | \mu, W, \sigma^2) = \sum_{n=1}^N \{ \ln p(x_n | z_n) + \ln p(z_n) \}$$

- The n^{th} row of Z is given by z_n .
- Recall that $p(z) = \mathcal{N}(z|\mathbf{0}, I)$, $p(x|z) = \mathcal{N}(x|Wz + \mu, \sigma^2 I)$
- We can now write the expectation with respect to the posterior distribution over the latent variables.

EM Algorithms for PCA

$$\mathbb{E} \left[\ln p(X, Z | \mu, W, \sigma^2) \right] = -\sum_{n=1}^N \left\{ \begin{aligned} & \frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2} \text{Tr}(\mathbb{E}[z_n z_n^T]) \\ & + \frac{1}{2\sigma^2} \|x_n - \mu\|^2 - \frac{1}{\sigma^2} \mathbb{E}[z_n]^T W^T (x_n - \mu) \\ & + \frac{1}{2\sigma^2} \text{Tr}(\mathbb{E}[z_n z_n^T] W^T W) + \frac{M}{2 \ln(2\pi)} \end{aligned} \right\}$$

□ **E-Step:** We use the old parameters to evaluate:

$$\mathbb{E}[z_n] = M^{-1}W^T(x_n - \bar{x}) \quad \mathbb{E}[z_n z_n^T] = \sigma^2 M^{-1} + \mathbb{E}[z_n] \mathbb{E}[z_n]^T$$

□ This follows directly from

$$p(z | x) = \mathcal{N}(z | M^{-1}W^T(x - \mu), \sigma^2 M^{-1}), \quad M = W^T W + \sigma^2 I$$

together with the standard result

$$\mathbb{E}[z_n z_n^T] = \text{cov}[z_n] + \mathbb{E}[z_n] \mathbb{E}[z_n]^T$$



EM Algorithms for PCA

$$\mathbb{E} \left[\ln p(X, Z | \mu, W, \sigma^2) \right] = -\sum_{n=1}^N \left\{ \begin{aligned} & \frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2} \text{Tr}(\mathbb{E}[z_n z_n^T]) \\ & + \frac{1}{2\sigma^2} \|x_n - \mu\|^2 - \frac{1}{\sigma^2} \mathbb{E}[z_n]^T W^T (x_n - \mu) \\ & + \frac{1}{2\sigma^2} \text{Tr}(\mathbb{E}[z_n z_n^T] W^T W) + \frac{M}{2 \ln(2\pi)} \end{aligned} \right\}$$

□ **M-Step:** We maximize with respect to W and σ^2 keeping the posterior statistics fixed. We obtain:

$$\begin{aligned} W_{new} &= \left[\sum_{n=1}^N (x_n - \bar{x}) \mathbb{E}[z_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[z_n z_n^T] \right]^{-1} \\ \sigma_{new}^2 &= \frac{1}{ND} \sum_{n=1}^N \left\{ \begin{aligned} & \|x_n - \bar{x}\|^2 - 2 \mathbb{E}[z_n]^T W_{new} (x_n - \bar{x}) \\ & + \text{Tr}(\mathbb{E}[z_n z_n^T] W_{new}^T W_{new}) \end{aligned} \right\} \end{aligned}$$



Proof of the M-Step Equations

$$\mathbb{E} \left[\ln p(X, Z | \mu, W, \sigma^2) \right] = -\sum_{n=1}^N \left\{ \begin{aligned} & \frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2} \text{Tr} \left(\mathbb{E} [z_n z_n^T] \right) \\ & + \frac{1}{2\sigma^2} \|x_n - \mu\|^2 - \frac{1}{\sigma^2} \mathbb{E} [z_n]^T W^T (x_n - \mu) \\ & + \frac{1}{2\sigma^2} \text{Tr} \left(\mathbb{E} [z_n z_n^T] W^T W \right) + \frac{M}{2 \ln(2\pi)} \end{aligned} \right\}$$

- The two M-equations are derived by setting the derivatives wrt W and σ^2 equal to zero:

$$\frac{\partial \mathbb{E} \left[\ln p(X, Z | \mu, W, \sigma^2) \right]}{\partial W} = \sum_{n=1}^N \left\{ \frac{1}{\sigma^2} (x_n - \mu) \mathbb{E} [z_n]^T - \frac{1}{\sigma^2} W \mathbb{E} [z_n z_n^T] \right\} = 0$$

$$\frac{\partial \mathbb{E} \left[\ln p(X, Z | \mu, W, \sigma^2) \right]}{\partial \sigma^2} = \sum_{n=1}^N \left\{ \begin{aligned} & -\frac{D}{2\sigma^2} - \frac{1}{\sigma^4} \mathbb{E} [z_n]^T W^T (x_n - \mu) + \frac{1}{2\sigma^4} \|x_n - \mu\|^2 \\ & + \frac{1}{2\sigma^4} \text{Tr} \left(\mathbb{E} [z_n z_n^T] W^T W \right) \end{aligned} \right\} = 0$$

- Here we used $\frac{\partial}{\partial A} \text{Tr}(ABA^T) = A(B + B^T)$, $\frac{\partial}{\partial A} \text{Tr}(AB) = B^T$

EM Algorithms for PCA

- Initialize the parameters
- Compute the sufficient statistics of the latent space posterior distribution in the E-Step

$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x}_n - \bar{\mathbf{x}}) \quad \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] = \sigma^2 \mathbf{M}^{-1} + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n^T]$$

- Revise the parameter values in the M-Step.

$$\mathbf{W}_{new} = \left[\sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) \mathbb{E}[\mathbf{z}_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \right]^{-1}$$

$$\sigma_{new}^2 = \frac{1}{ND} \sum_{n=1}^N \left\{ \begin{aligned} & \|\mathbf{x}_n - \bar{\mathbf{x}}\|^2 - 2[\mathbf{z}_n]^T \mathbf{W}_{new}^T (\mathbf{x}_n - \bar{\mathbf{x}}) \\ & + Tr(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}_{new}^T \mathbf{W}_{new}) \end{aligned} \right\}$$

EM Algorithms for PCA: Computational Cost

- Each cycle of the EM algorithm can be computationally more efficient than conventional PCA in high dimensions.
- The eigen-decomposition of the covariance matrix requires $\mathcal{O}(D^3)$ computation. If interested only in the first M eigenvectors, we can use algorithms that are $\mathcal{O}(MD^2)$.
- However, the evaluation of the covariance matrix itself

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$

takes $\mathcal{O}(ND^2)$ computations.

- Algorithms such as the snapshot method ([Sirovich, 1987](#)), assume that the eigenvectors are linear combinations of the data vectors and avoid direct evaluation of the covariance matrix but are $\mathcal{O}(N^3)$ and hence unsuited to large data sets.

- Sirovich, L. (1987). [Turbulence and the dynamics of coherent structures](#). *Quarterly Applied Mathematics* **45**(3), 561–590.



EM Algorithms for PCA: Computational Cost

- The EM algorithm described here also does not construct the covariance matrix explicitly.
- Instead, the most computationally demanding steps are those involving sums over the data set that are $\mathcal{O}(NDM)$.

$$W_{new} = \left[\sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) \mathbb{E}[\mathbf{z}_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \right]^{-1}$$

- For large D , and $M \ll D$, $\mathcal{O}(NDM)$ is much better compared to $\mathcal{O}(ND^2)$ (that we show is needed to compute the sample covariance)

This offsets the iterative nature of the EM algorithm.



EM Algorithms for PCA: Online form

- The EM algorithm can be implemented in an on-line form: each D -dim data point is read/processed/then discarded before the next data point is considered.
- To see this, note that the quantities evaluated in the E-Step (an M -dimensional vector and an $M \times M$ matrix) can be computed for each data point separately.

$$\mathbb{E}[z_n] = M^{-1}W^T(x_n - \bar{x}) \quad \mathbb{E}[z_n z_n^T] = \sigma^2 M^{-1} + \mathbb{E}[z_n] \mathbb{E}[z_n]^T$$

- In the M-Step, we accumulate the sums over data points incrementally - advantageous if both N & D are large.

$$W_{new} = \left[\sum_{n=1}^N (x_n - \bar{x}) \mathbb{E}[z_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[z_n z_n^T] \right]^{-1}$$
$$\sigma_{new}^2 = \frac{1}{ND} \sum_{n=1}^N \left\{ \|x_n - \bar{x}\|^2 - 2\mathbb{E}[z_n]^T W_{new}^T (x_n - \bar{x}) \right\} \\ \left. + Tr(\mathbb{E}[z_n z_n^T] W_{new}^T W_{new}) \right\}$$



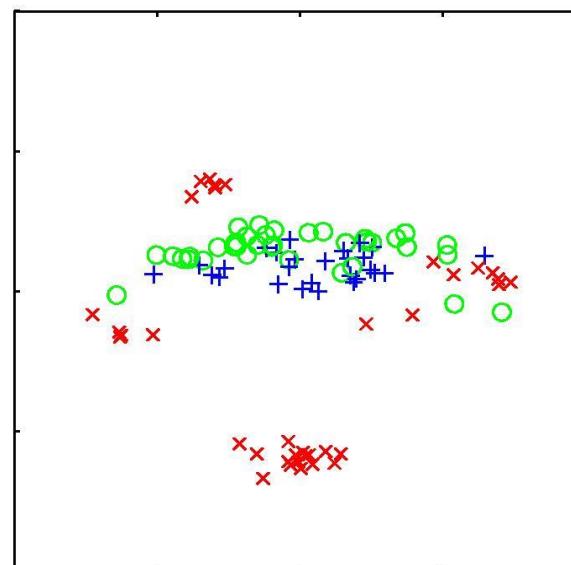
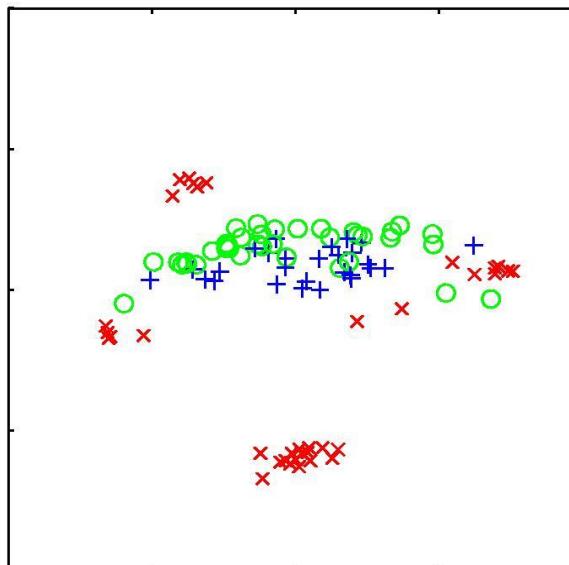
EM Algorithms for PCA: Missing Values

- Because we now have a fully probabilistic model for PCA, we can **deal with missing data**, provided that it is **missing at random**, by marginalizing over the distribution of the unobserved variables.
- Again these missing values can be treated using the EM algorithm.
- We give an example of the use of this approach for data visualization in the figure shown next.



EM Algorithms for PCA: Missing Values

- PPCA visualization for the first 100 data points of the oil flow data set.
- Left: the posterior mean projections of the data points on the principal subspace.
- Right: randomly omitting 30% of the variable values and using EM to handle the missing values. Even though each data point has at least one missing measurement, the plot is similar to the one obtained without missing values.



[Matlab Implementation](#)

EM Algorithm for PCA: Limit $\sigma^2 \rightarrow 0$

- When $\sigma^2 \rightarrow 0$, EM corresponds to standard PCA ([Roweis, 1998](#))
- Defining $\tilde{\mathbf{X}}$ a matrix of size $N \times D$ whose n^{th} row is given by $\mathbf{x}_n - \bar{\mathbf{x}}$.
- Defining a matrix Ω of size $M \times N$ whose n^{th} column is given by the vector $\mathbb{E}[\mathbf{z}_n]$.
- The E-Step for $\sigma^2 \rightarrow 0$ becomes

$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x}_n - \bar{\mathbf{x}}) \Rightarrow \Omega = (\mathbf{W}_{old}^T \mathbf{W}_{old})^{-1} \mathbf{W}_{old}^T \tilde{\mathbf{X}}^T$$

This is simply the orthogonal projection of the data points on the current estimate for the principal subspace.

- Roweis, S. (1998). [EM algorithms for PCA and SPCA](#). In M. I. Jordan, M. J. Kearns, and S. A. Solla (Eds.), [Advances in Neural Information Processing Systems, Volume 10](#), pp. 626–632. MIT Press.

EM Algorithm for PCA: Limit $\sigma^2 \rightarrow 0$

- When $\sigma^2 \rightarrow 0$, EM corresponds to standard PCA
- $\tilde{\mathbf{X}}$ an $N \times D$ matrix whose n^{th} row is given by $\mathbf{x}_n - \bar{\mathbf{x}}$.
- Ω a matrix $M \times N$ whose n^{th} column is given by $\mathbb{E}[\mathbf{z}_n]$.
- Noting that

$$\sigma^2 \rightarrow 0 \Rightarrow \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] = \sigma^2 \mathbf{M}^{-1} + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^T \rightarrow \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^T$$

the M-Step takes the form:

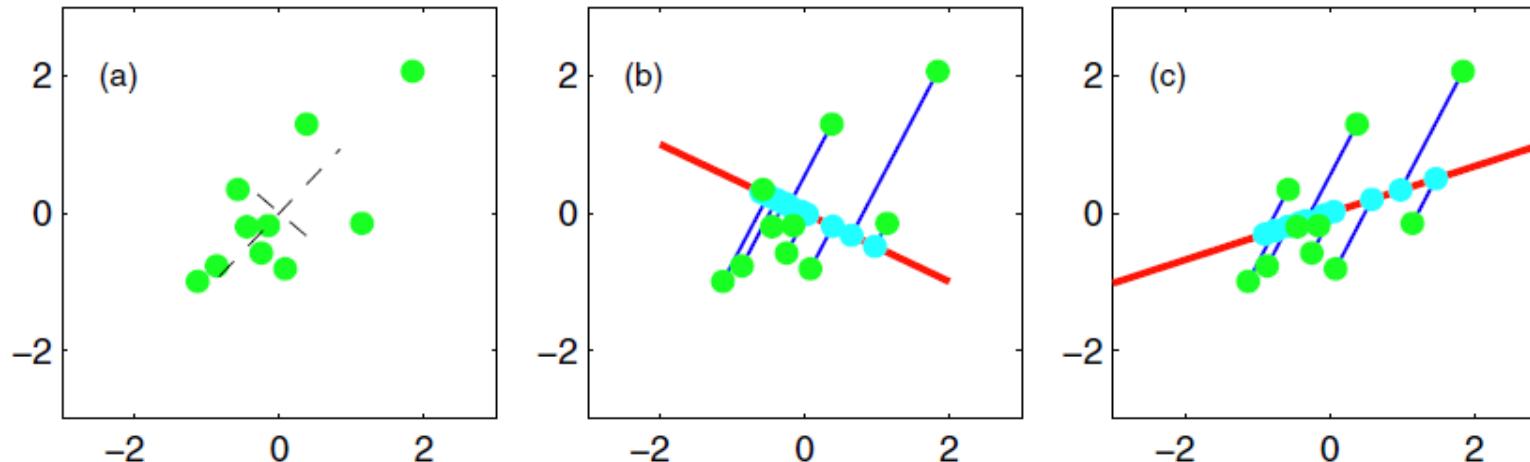
$$\mathbf{W}_{\text{new}} = \left[\sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) \mathbb{E}[\mathbf{z}_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \right]^{-1} \Rightarrow \mathbf{W}_{\text{new}} = \tilde{\mathbf{X}}^T \boldsymbol{\Omega}^T (\boldsymbol{\Omega} \boldsymbol{\Omega}^T)^{-1}$$

Re-estimation of the principal subspace minimizing the squared reconstruction errors in which the projections are fixed (see interpretation next and also [here](#))

EM Algorithm for PCA: Example, D=2, M=1

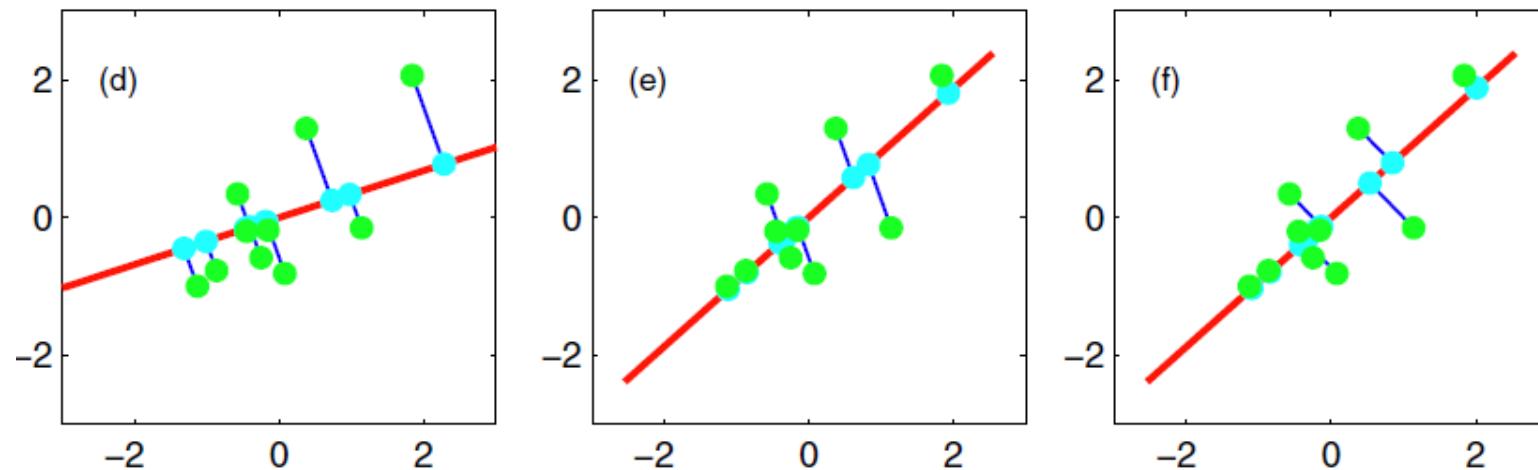
- Synthetic data illustrating the EM algorithm for PCA
 - (a) A data set X with the data points (green), together with the true principal components (eigenvectors scaled by the square roots of the eigenvalues).
 - (b) Initial configuration of the principal subspace defined by W (red) together with the projections of the latent points Z into the data space, given by ZW^T (cyan)
 - (c) After one M-Step, the latent space has been updated with Z held fixed.

[Matlab Implementation](#)



EM Algorithm for PCA: Example, D=2, M=1

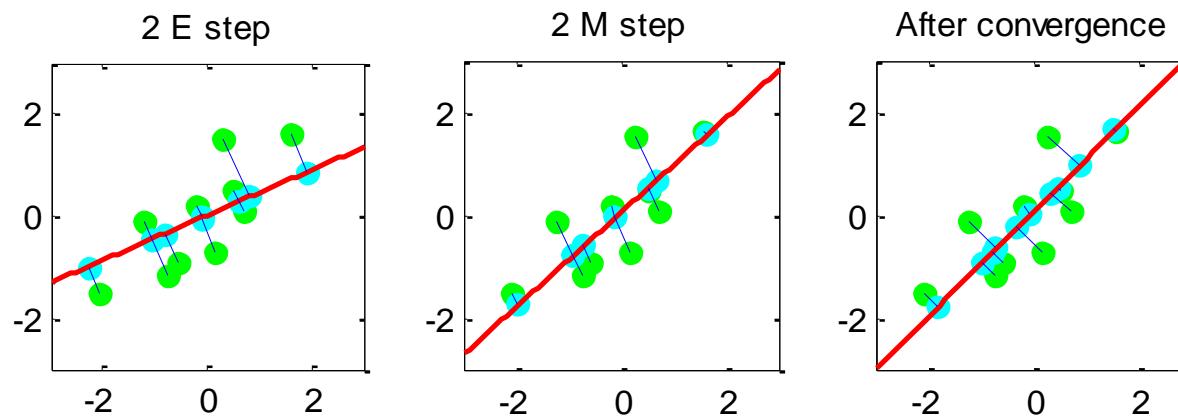
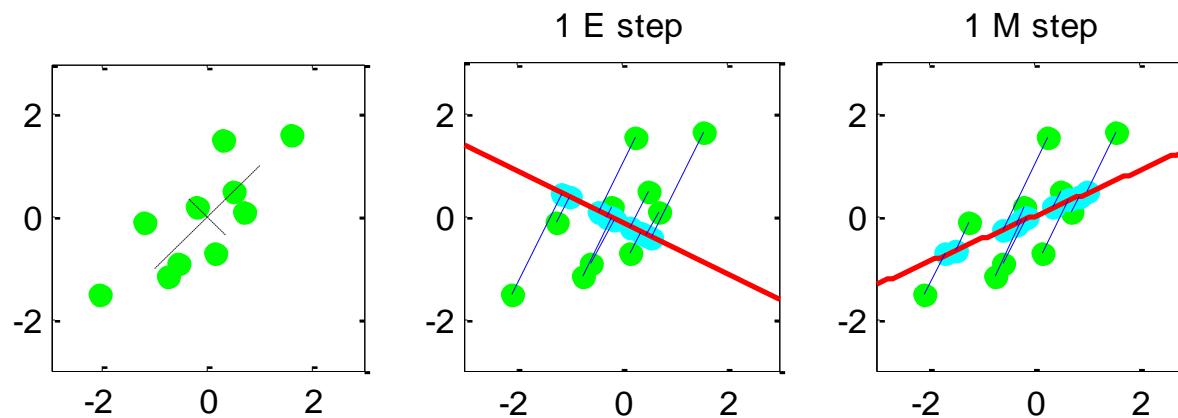
- Synthetic data illustrating the EM algorithm for PCA
 - (d) After the successive E-Step, the values of Z have been updated, giving orthogonal projections, with W held fixed.
 - (e) After the second M-Step.
 - (f) The converged solution.



$$\Omega = (\mathbf{W}_{old}^T \mathbf{W}_{old})^{-1} \mathbf{W}_{old}^T \tilde{\mathbf{X}}^T$$

$$\mathbf{W}_{new} = \tilde{\mathbf{X}}^T \Omega^T (\Omega \Omega^T)^{-1}$$

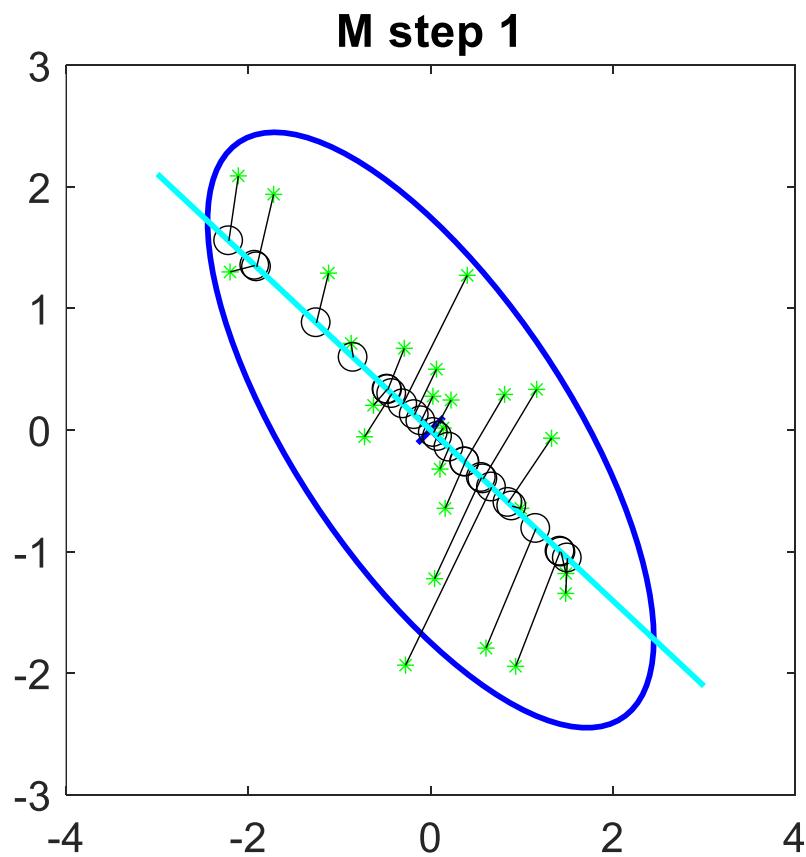
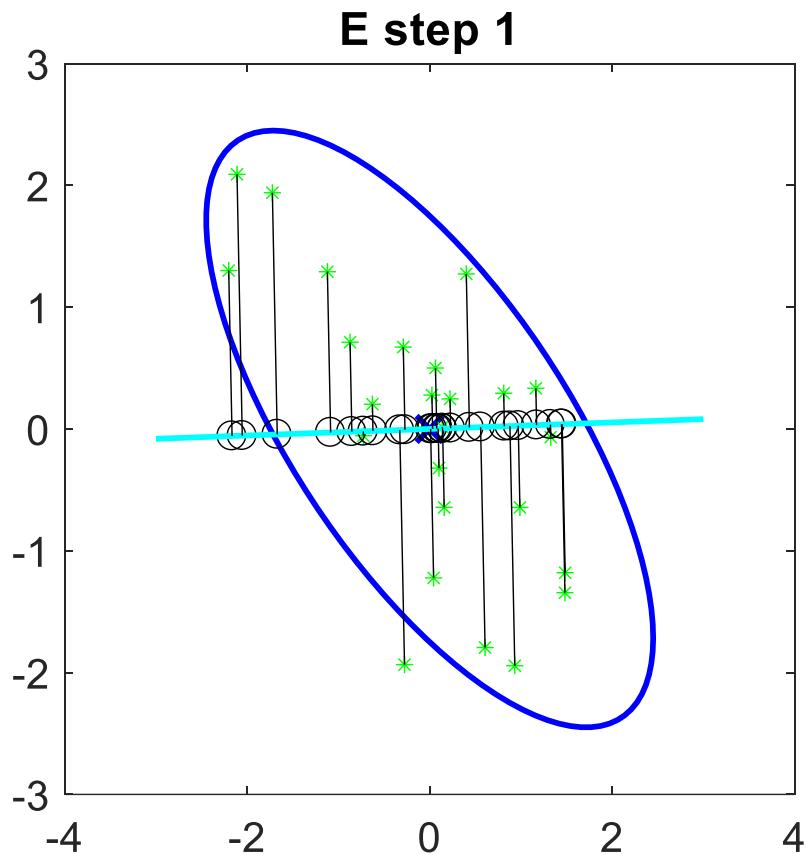
EM Algorithm for PCA: Example, D=2, M=1



[Matlab Implementation](#)



EM Algorithm for PCA: Example, D=2, M=1

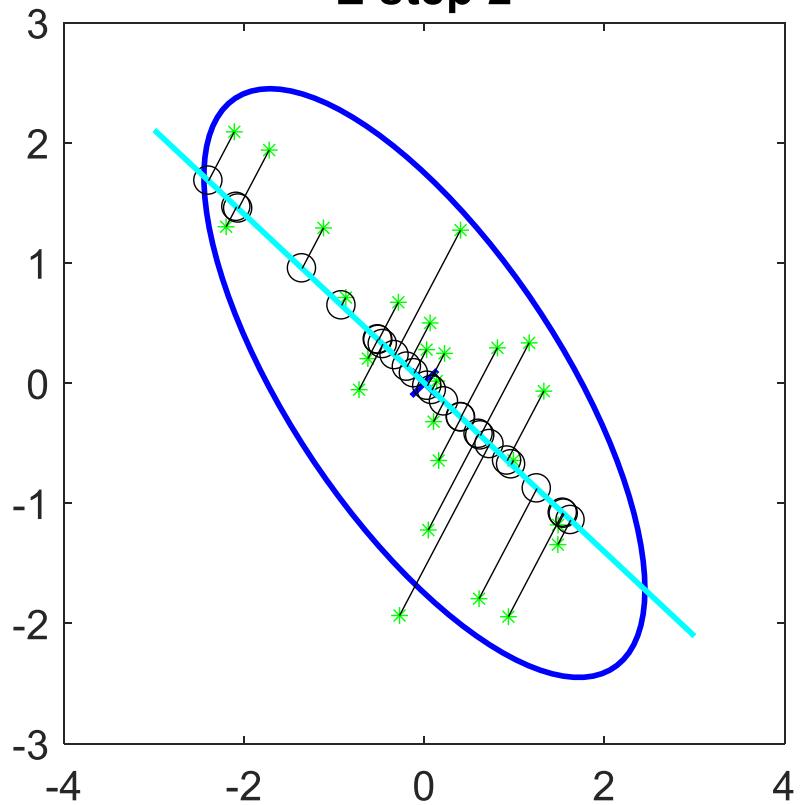


Run [pcaEmStepByStep](#)
From [PMTK3](#)

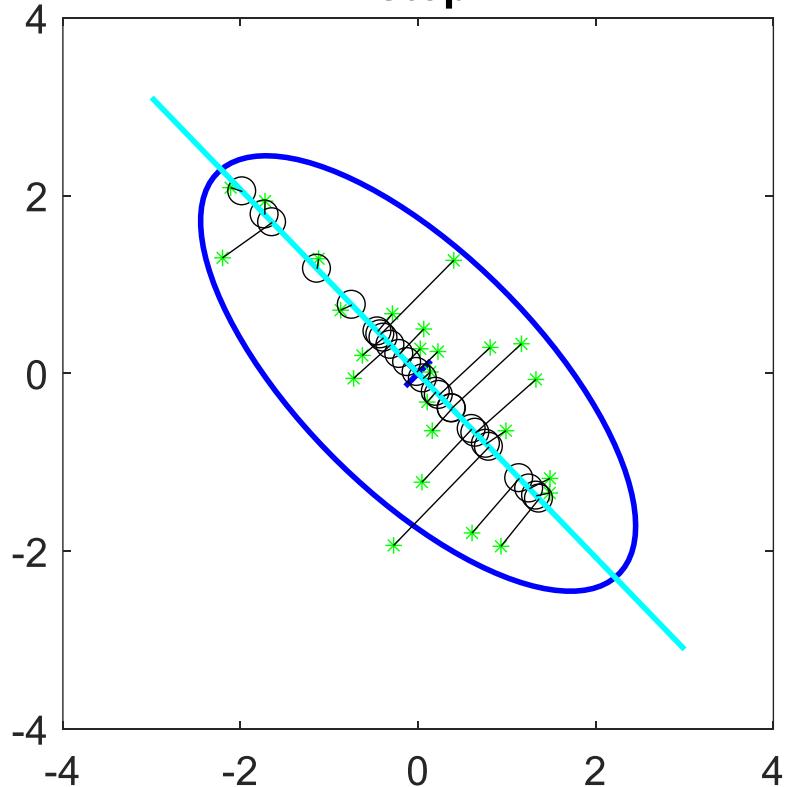


EM Algorithm for PCA: Example, D=2, M=1

E step 2



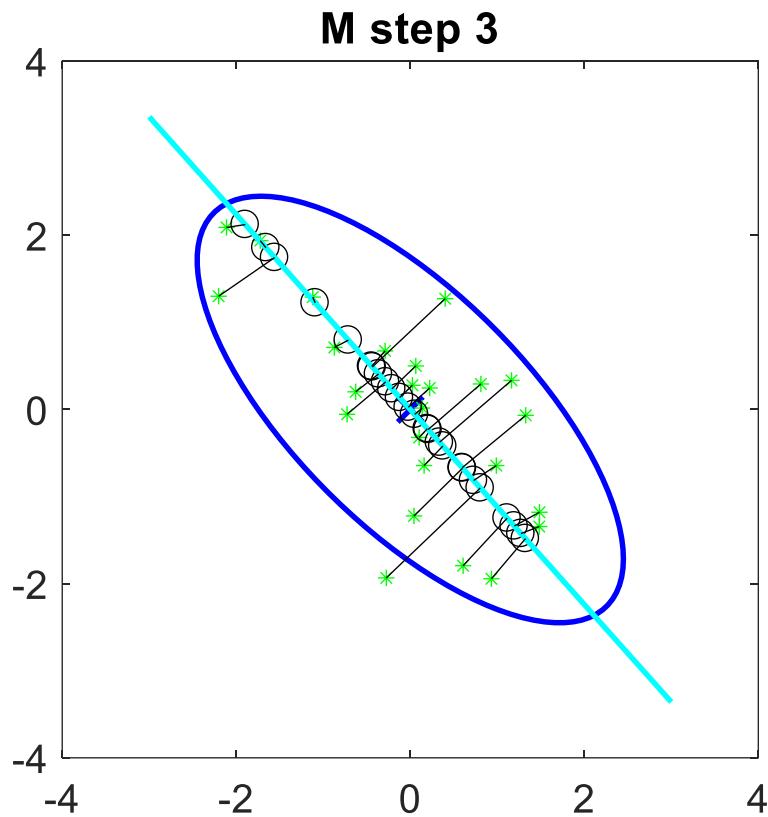
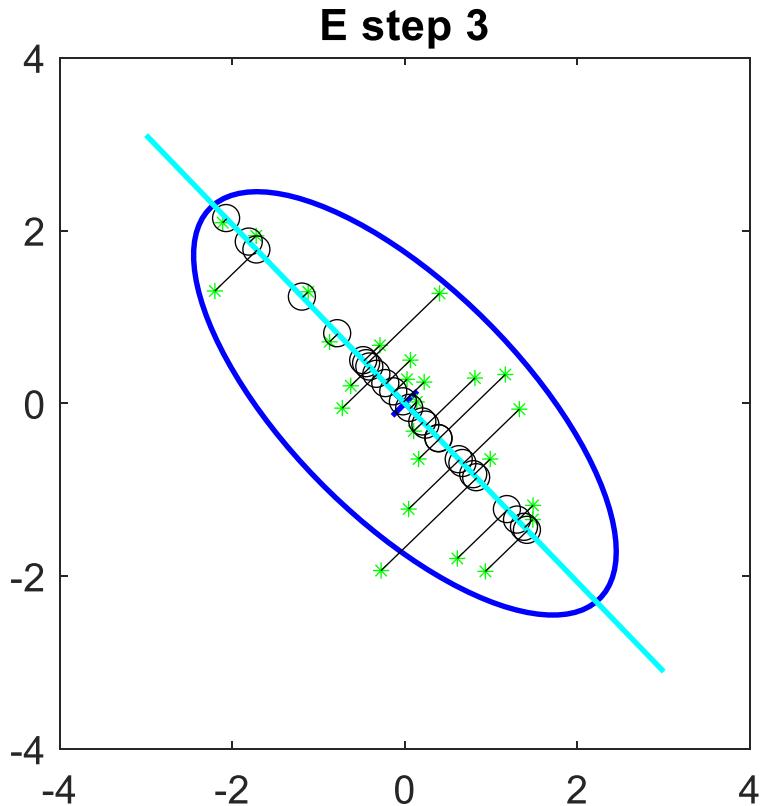
M step 2



Run [pcaEmStepByStep](#)
From [PMTK3](#)



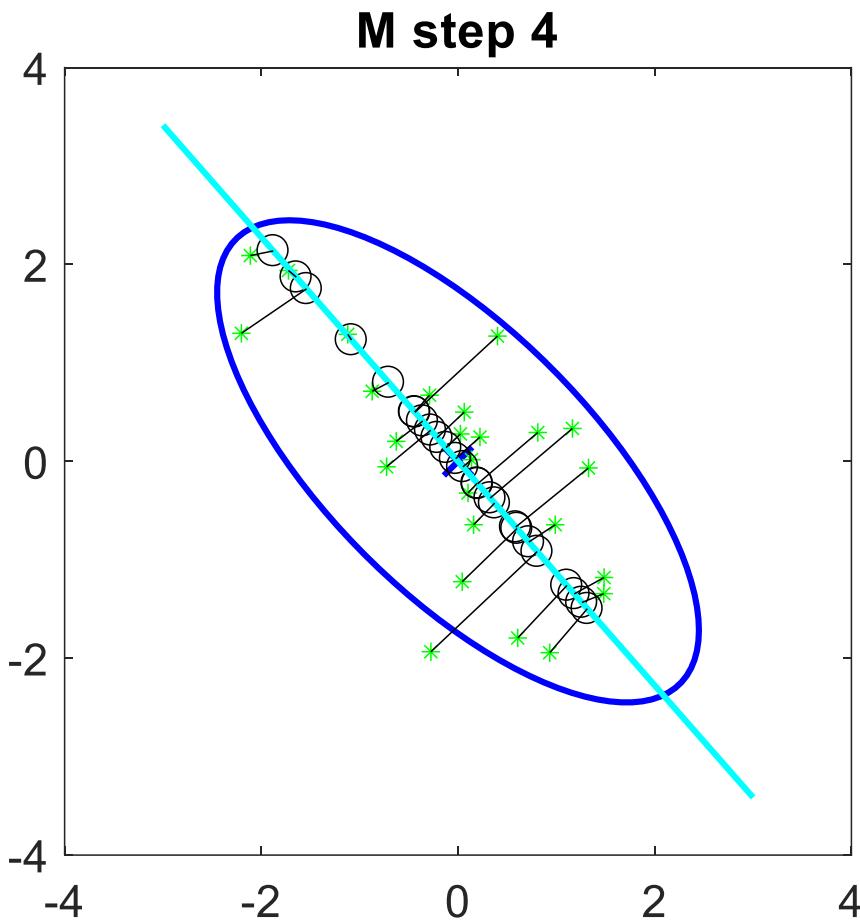
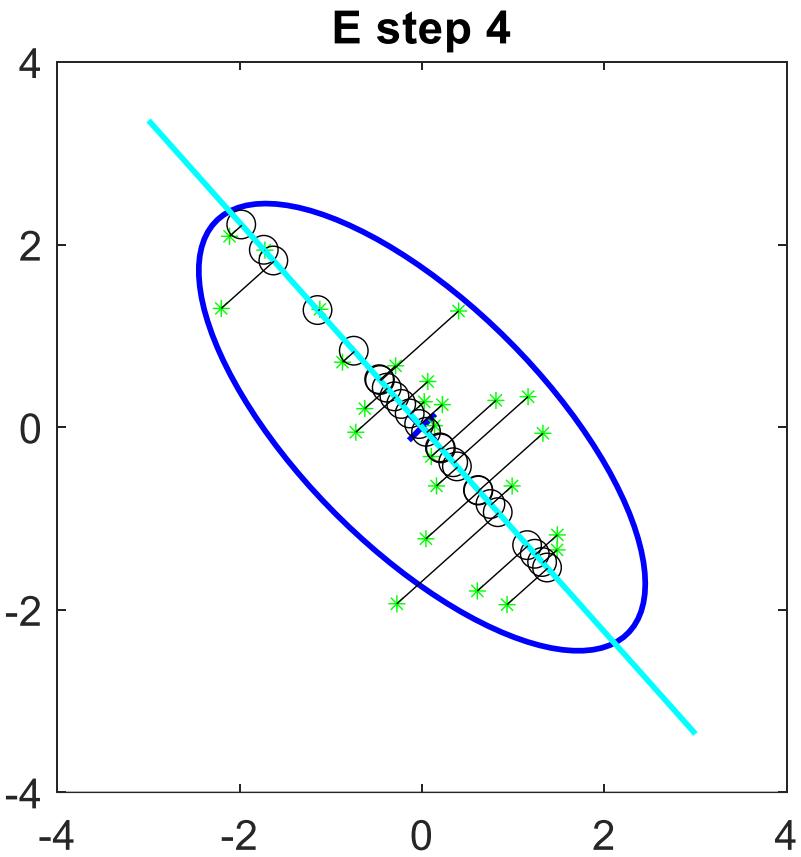
EM Algorithm for PCA: Example, D=2, M=1



Run [pcaEmStepByStep](#)
From [PMTK3](#)



EM Algorithm for PCA: Example, D=2, M=1



Run [pcaEmStepByStep](#)
From [PMTK3](#)



EM Algorithm and Standard PCA Revisited

- Let \mathbf{W} be $D \times M$ whose columns define a linear subspace of dimensionality M embedded within a data set of dimensionality D .
- Let $\boldsymbol{\mu}$ a D -dimensional vector. We approximate the data points $\mathbf{x}_n, n = 1, \dots, N$ using a linear mapping from a set of D -dimensional vectors \mathbf{z}_n as $\mathbf{Wz}_n + \boldsymbol{\mu}$.
- The sum of squares reconstruction error is:

$$J = \sum_{n=1}^N \|\mathbf{x}_n - \boldsymbol{\mu} - \mathbf{Wz}_n\|^2$$

- Minimizing wrt $\boldsymbol{\mu}$ gives:

$$0 = -\sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu} - \mathbf{Wz}_n) \Rightarrow \boldsymbol{\mu} = \bar{\mathbf{x}} - \mathbf{W}\bar{\mathbf{z}}$$

- This modifies the error as:

$$J = \sum_{n=1}^N \|(\mathbf{x}_n - \bar{\mathbf{x}}) - \mathbf{W}(\mathbf{z}_n - \bar{\mathbf{z}})\|^2$$



EM Algorithm and Standard PCA Revisited

- Let X be a $N \times D$ matrix with n^{th} row $(x_n - \bar{x})^T$.
- Similarly let Z be a $N \times M$ matrix with n^{th} row $(z_n - \bar{z})^T$.
- The cost function can be written as:

$$J = \sum_{n=1}^N \| (x_n - \bar{x}) - W(z_n - \bar{z}) \|^2 \Rightarrow J = \text{Tr} \left\{ (X - ZW^T)(X - ZW^T)^T \right\}$$

- Setting the derivative wrt Z equal to 0 gives (using $\frac{\partial}{\partial A} \text{Tr}(ABA^T) = A(B + B^T)$ and $\frac{\partial}{\partial A} \text{Tr}(AB) = B^T, \frac{\partial}{\partial A} \text{Tr}(A^T B) = B$):

$$-2XW + 2ZW^TW = 0 \Rightarrow Z = XW(W^TW)^{-1}$$

This is similar to the PCA E-Step: $\Omega = (W_{old}^T W_{old})^{-1} W_{old}^T \tilde{X}^T$

- Similarly derivative wrt W equal to 0 gives (analogous to the PCA M-Step): $W_{new} = \tilde{X}^T \Omega^T (\Omega \Omega^T)^{-1}$

$$-2X^T Z + 2WZ^T Z = 0 \Rightarrow W = X^T Z (Z^T Z)^{-1}$$

Model Selection for FA and PPCA

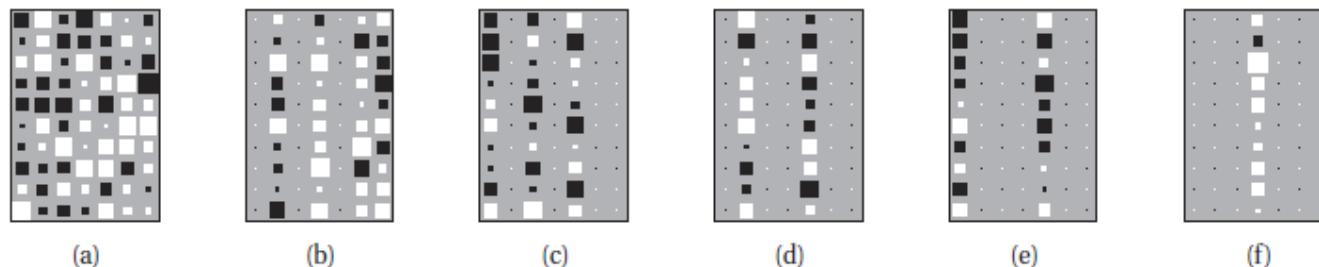
- We can in principle compute $M^* = \operatorname{argmax}_M p(M|\mathcal{D})$.
- There are two problems with this.
 - Evaluating the marginal likelihood for LVMs is quite difficult.
 - ✓ One can use **BIC** or **variational lower bounds**.
 - ✓ Alternatively, we can use the **cross-validated likelihood** as a performance measure - this can be slow, since it requires fitting each model F times, where F is the number of CV folds.
 - The second issue is the need to search over a potentially large number of models.
 - ✓ The usual approach is to perform exhaustive search over all candidate values of M .
 - ✓ However, sometimes we can set the model to its maximal size, and then use **automatic relevancy determination** combined with EM to prune out irrelevant weights.
- Minka, T. (2000a). [Automatrical choice of dimensionality for PCA](#). Technical report, MIT.
- Bishop, C. (1999). [Bayesian PCA](#). In *NIPS*.
- Ghahramani, Z. and M. Beal (2000). [Variational inference for Bayesian mixtures of factor analysers](#). In *NIPS-12*.
- Lopes, H. and M. West (2004). [Bayesian model assessment in factor analysis](#). *Statistica Sinica* 14, 41–67.
- Paisley, J. and L. Carin (2009). [Nonparametric factor analysis with beta process priors](#). In *Intl. Conf. on Machine Learning*.



Model Selection for FA and PPCA

- The ARD approach is shown in the figure using **Hinton diagrams** (for a mixture of FA using VBEM).
- The degree of sparsity depends on the amount of training data in accord with Occam's razor. When the sample size is small, the method prefers simpler models, but as the sample size gets sufficiently large, the method converges on the “correct” solution.

The data was generated from 6 clusters with intrinsic dimensionalities of 7, 4, 3, 2, 2, 1, which the method has successfully estimated.



number of points per cluster	intrinsic dimensionalities					
	1	7	4	3	2	2
8		2			1	
8	1			2		
16	1		4		2	
32	1	6	3	3	2	2
64	1	7	4	3	2	2
128	1	7	4	3	2	2

Estimated number of clusters, and their estimated dimensionalities, as a function of sample size. The VBEM algorithm found two different solutions when $N = 8$. More clusters with larger effective dimensionalities are discovered as the sample sizes increases.

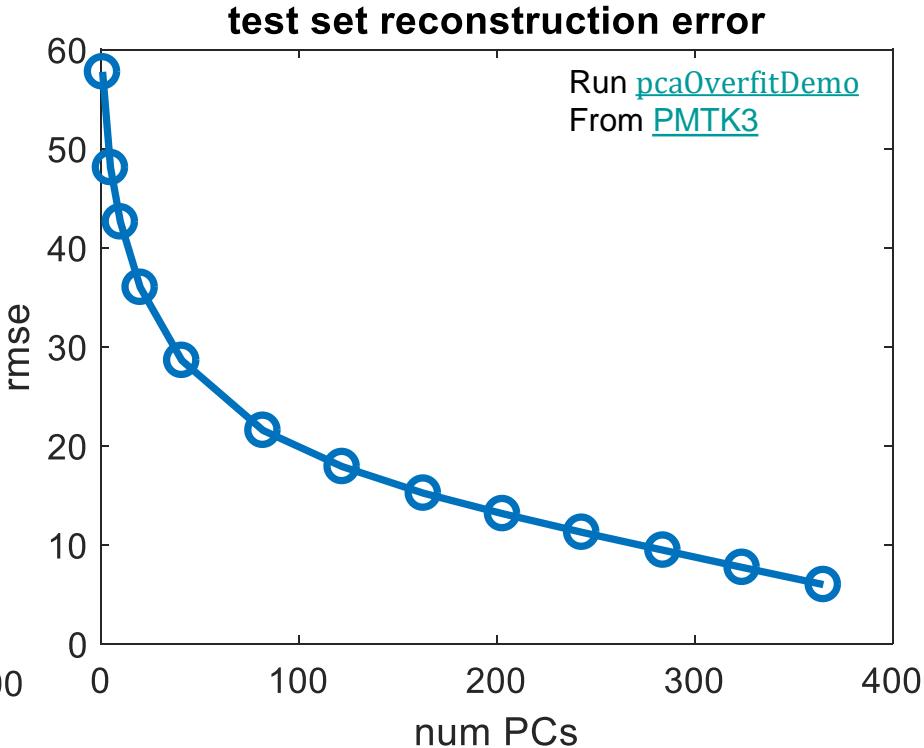
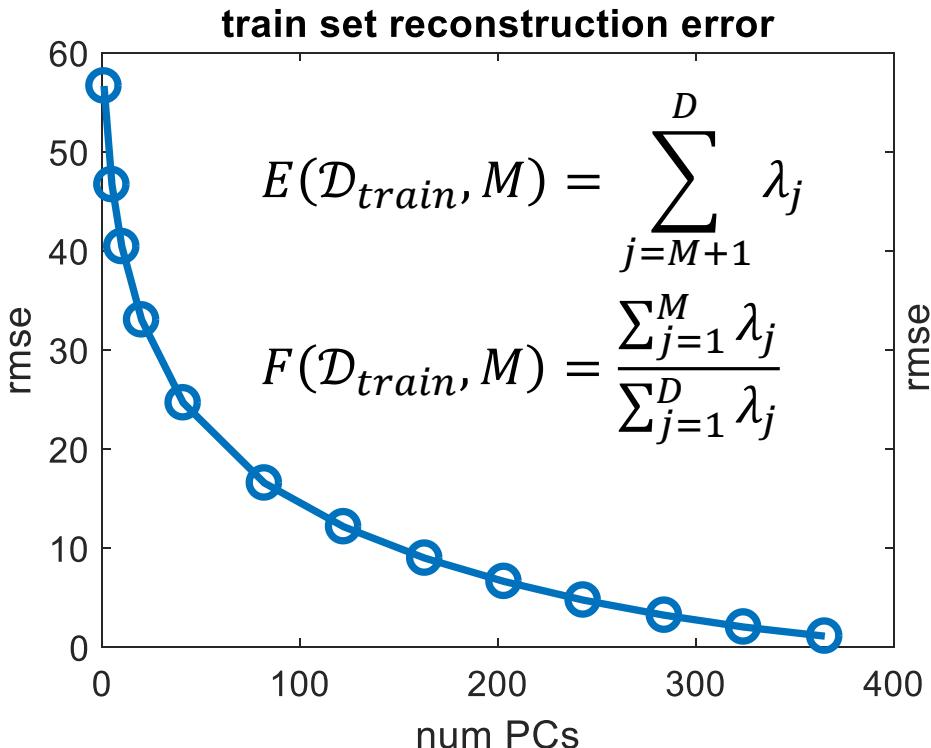
- Beal, M. (2003). [Variational Algorithms for Approximate Bayesian Inference](#). Ph.D. thesis, Gatsby Unit.

Model Selection for PCA

- Since PCA is not a probabilistic model, we cannot use any of the earlier discussed methods. An option is to **use the reconstruction error**:

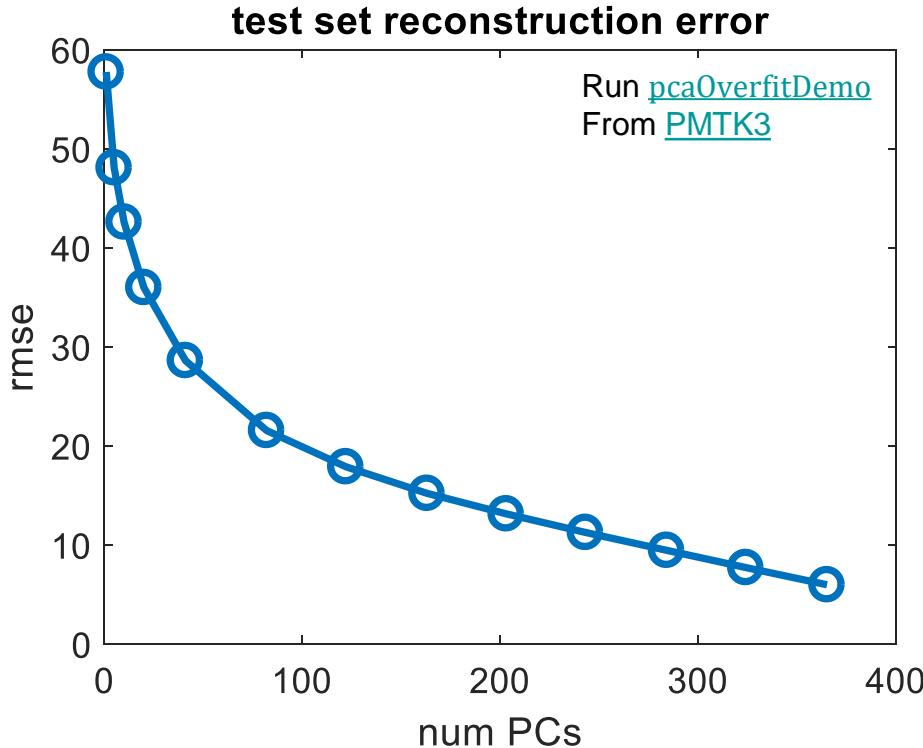
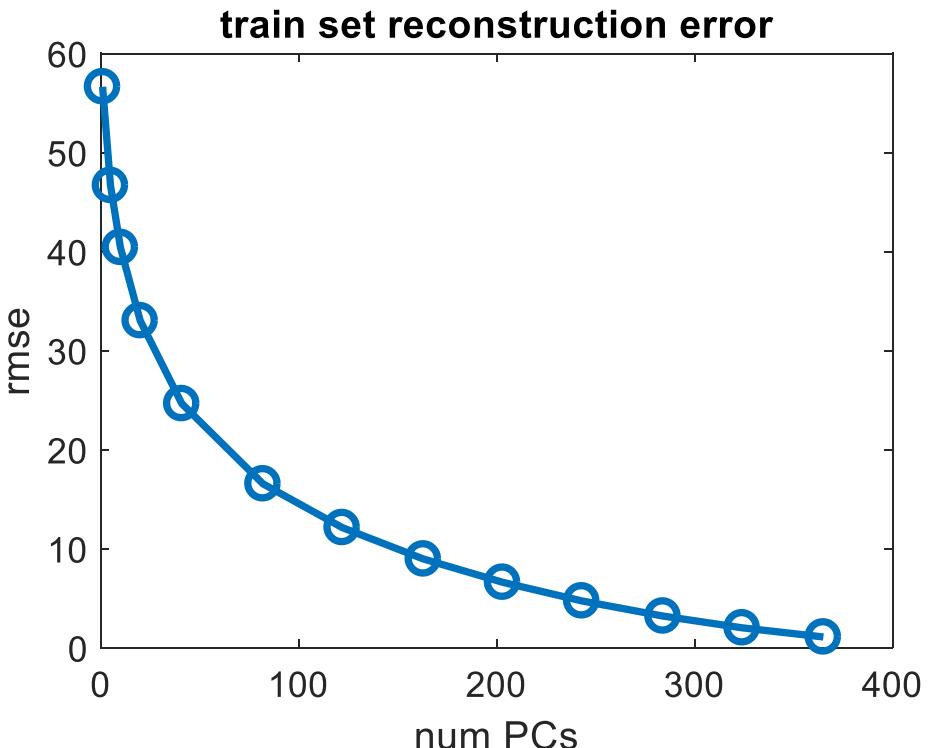
$$E(\mathcal{D}, M) = \frac{1}{|\mathcal{D}|} \sum_{i \in \mathcal{D}} \|x_i - \hat{x}_i\|^2$$

$$\hat{x}_i = \mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \mathbf{z}_i = \mathbf{W}^T(x_i - \boldsymbol{\mu})$$



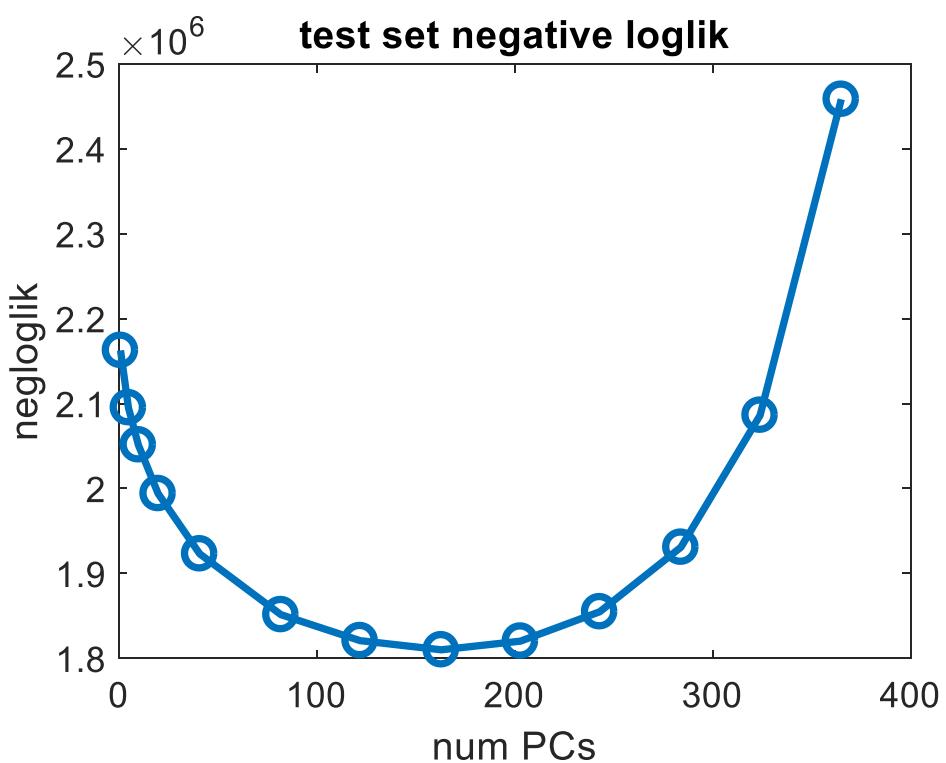
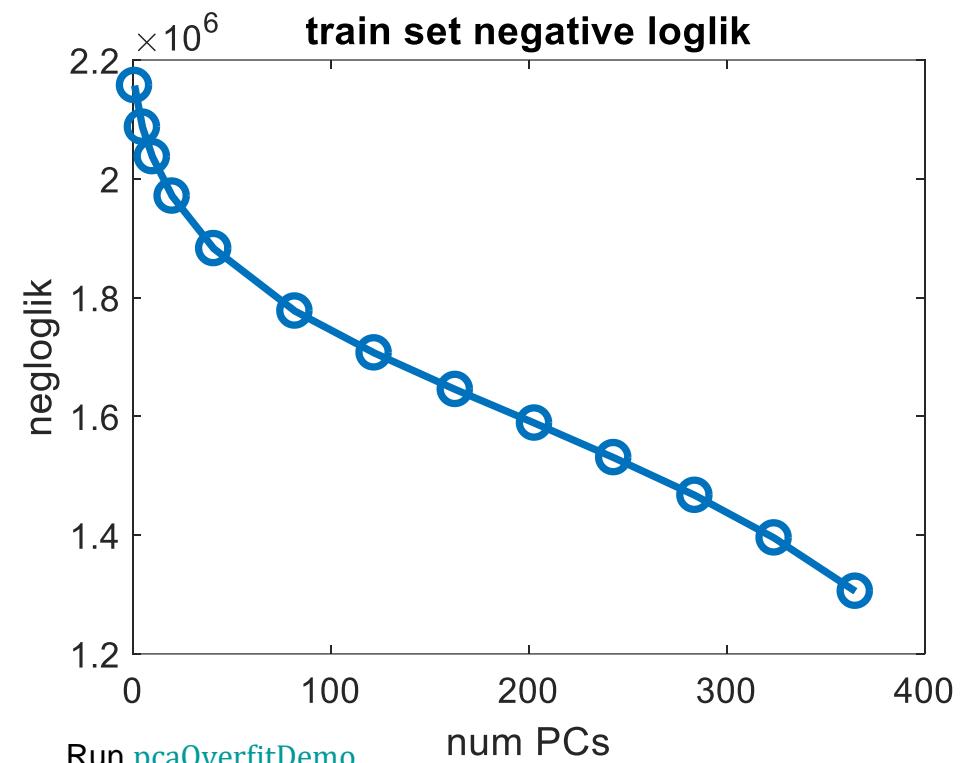
Model Selection for PCA

- If we use $M = \text{rank}(\mathbf{X})$, we get zero reconstruction error on the training set.
- We see that the **test error** continues to go down even as the model becomes more complex! Thus **we do not get the usual U-shaped curve**.
- The problem is that PCA is not a proper generative model of the data (a data compression technique). If you give it more latent dimensions it will be able to approximate the test data more accurately.



Model Selection for PCA

- When plotting the negative log likelihood on the test set we do indeed see the U shape.
- The nature of these results are similar to the ones for the K-means algorithm (also non-probabilistic method).



Run [pcaOverfitDemo](#)
From [PMTK3](#)



Profile Likelihood

- Although there is no U-shape, there is sometimes a change in the plots, from relatively large errors to relatively small.
- Let λ_k be a measure of the error incurred by a model of size k , such that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{M_{max}}$. In PCA these are the eigenvalues.
- Partitioning these values into two groups, depending on whether $k < M$ or $k > M$, where M is some threshold which we will determine. To measure the quality of M , we will use a simple change-point model, where $\lambda_k \sim \mathcal{N}(\mu_1, \sigma^2)$ if $k \leq M$, and $\lambda_k \sim \mathcal{N}(\mu_2, \sigma^2)$ if $k > M$. (It is important that σ^2 be the same in both models, to prevent overfitting in the case where one regime has less data than the other)
- Within each regime, we assume the λ_k are iid. We can fit this model for each $M = 1 : M_{max}$ by partitioning the data and computing the MLEs, using a pooled estimate of the variance.

- Zhu, M. and A. Ghodsi (2006). [Automatic dimensionality selection from the scree plot via the use of profile likelihood](#). *Computational Statistics & Data Analysis* 51, 918– 930.



Profile Likelihood

- The MLE estimates are given as:

$$\mu_1(M) = \frac{\sum_{k < M} \lambda_k}{M}, \mu_2(M) = \frac{\sum_{k > M} \lambda_k}{N - M}$$

$$\ell(M) = \sum_{k=1}^M \log \mathcal{N}(\lambda_k | \mu_1(M), \sigma^2(M)) + \sum_{k=M+1}^K \log \mathcal{N}(\lambda_k | \mu_2(M), \sigma^2(M))$$

- Finally we choose:

$$M^* = \operatorname{argmax} \ell(M)$$

- Zhu, M. and A. Ghodsi (2006). [Automatic dimensionality selection from the scree plot via the use of profile likelihood](#). *Computational Statistics & Data Analysis* 51, 918– 930.

Profile Likelihood

