# Continuous Latent Variable Models: Probabilistic and Bayesian PCA,Kernel PCA

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

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

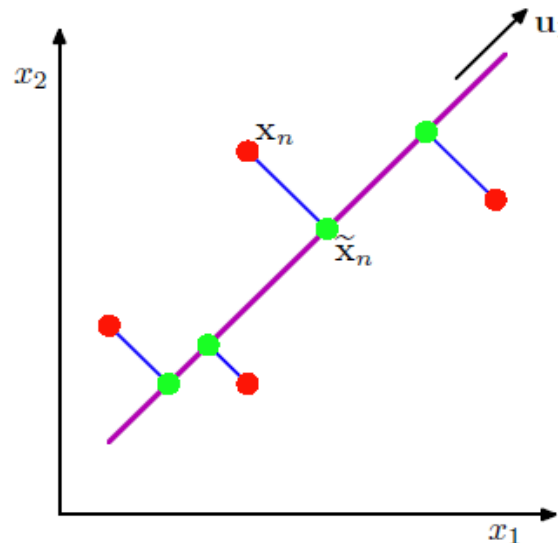November 8, 2017

# *Contents*

Following closely Chris Bishops' PRML book (Chapter 12)

# *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 with 1D Principal Subspace

❑ Consider first the projection onto 1D space ($M = 1$).

❑ We define the direction of this subspace using vector $\boldsymbol{u}_1$, which we choose to be a unit vector: $\boldsymbol{u}_1^T \boldsymbol{u}_1 = 1$

❑ Each data point $\boldsymbol{x}_n$ is then projected onto the scalar $\boldsymbol{u}_1^T \boldsymbol{x}_n$

❑ The mean of the projected data is given as: $\boldsymbol{u}_1^T \overline{\boldsymbol{x}}$

❑ Similarly, the variance of the projected data is:

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

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

# *PCA with 1D Principal Subspace*

❑ Introduce the Lagrange multiplier $\lambda_1$ and define the unconstrained maximization of

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

❑ We can see immediately that the solution satisfies:

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

❑ $\boldsymbol{u}_1$ must be an eigenvector of $\boldsymbol{S}$ with eigenvalue $\lambda_1$

❑ From the eigen-problem note that the variance of the projected data is $\boldsymbol{u}_1^T\boldsymbol{S}\boldsymbol{u}_1=\lambda_1$ so $\lambda_1$ needs to be the largest eigenvalue of $\boldsymbol{S}$.

❑ $\boldsymbol{u}_1$ is called the first principal component.
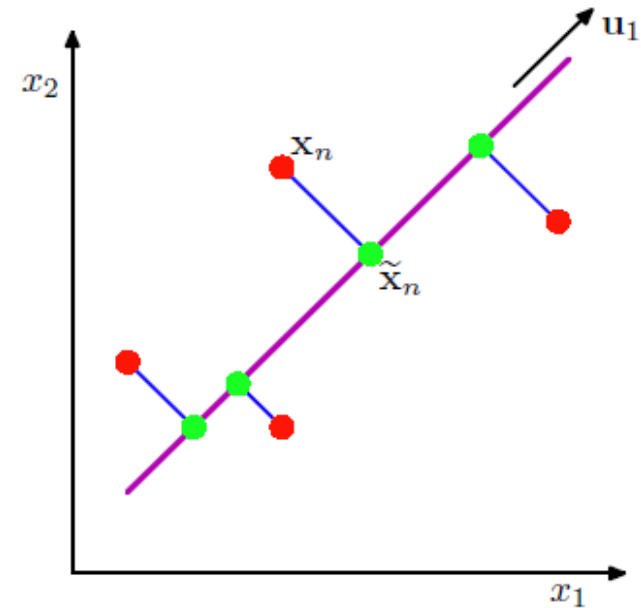
# PCA with 1D Principal Subspace

❑ 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 $\boldsymbol{u}_1, \dots, \boldsymbol{u}_M$ of the data covariance matrix $\boldsymbol{S}$ 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 $\boldsymbol{S}$) is $\mathcal{O}(MD^2)$.

# *PCA: Minimum Error Formulation*

$$x_n - \widetilde{x}_n = \sum_{i=M+1}^{D} \left\{ (x_n - \overline{x})^T u_i \right\} u_i$$

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

❑ We can also see this geometrically:

■ The projected points $\widetilde{x}_n$ 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 $\{u_i\}$ in the form

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

$$J = \frac{1}{N} \sum_{i=1}^{N} \| x_n - \tilde{x}_n \|^2$$

$$J = \frac{1}{N} \sum_{n=1}^{N} \sum_{i=M+1}^{D} \left( x_n^T u_i - \bar{x}^T u_i \right)^2 = \sum_{i=M+1}^{D} u_i^T S u_i$$

❑ The minimum of $J$ is obtained when $\{u_i\}, i = M+1, \ldots, D$ are the eigenvectors of $S$ 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 $\lambda_1$, and the min value of $J = \lambda_2$ is obtained by choosing $u_2$ corresponding to $\lambda_2$.

# *Compression of the Original Data Set*

❑ Using the earlier equations, we can derive:

$$\widetilde{x}_n = \sum_{i=1}^{M} z_{ni} u_i + \sum_{i=M+1}^{D} b_i u_i$$

$$z_{nj} = x_n^T u_j \ , \ j = 1,\dots,M$$

$$b_j = \overline{x}^T u_j \ , \ j = M+1,\dots,D$$

$$\Rightarrow \widetilde{x}_n = \sum_{i=1}^{M} (x_n^T u_i) u_i + \sum_{i=M+1}^{D} (\overline{x}^T u_i) u_i \Rightarrow$$

$$Data\ Compression: \ \widetilde{x}_n = \overline{x} + \sum_{i=1}^{M} (x_n^T u_i - \overline{x}^T u_i) u_i$$

❑ In deriving the last equation, we used the completeness of $u_i$, i.e.

$$\overline{x} = \sum_{i=1}^{D} \left( \overline{x}^T u_i \right) u_i$$

❑ This clearly shows the compression of the data set: From the $D$-dimensional $x_n$ to the $M$-dimensional vector with components $\left( x_n^T u_i - \overline{x}^T u_i \right)$.

# PCA Reconstruction

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

$$\widetilde{\boldsymbol{x}}_n = \overline{\boldsymbol{x}} + \sum_{i=1}^{M} \left( \boldsymbol{x}_n^T \boldsymbol{u}_i - \overline{\boldsymbol{x}}^T \boldsymbol{u}_i \right) \boldsymbol{u}_i$$

❑ If $M = D$, there is no dimensionality reduction but a rotation to align with the principal components e.g from $\{x_{n_1}, \ldots, x_{n_D}\}$ to $\{a_{n_1}, \ldots, a_{n_D}\}$ where:

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

# *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:

$$\boldsymbol{SU} = \boldsymbol{UL}, \quad \boldsymbol{L} = diag\left(\lambda_1, ..., \lambda_D\right), \quad \boldsymbol{U} = \left[\boldsymbol{u}_1 ... \boldsymbol{u}_D\right](orthogonal)$$

❑ For each data point $\boldsymbol{x}_n$, define a transformed value as:

$$\text{Whitening of the data}: \boldsymbol{y}_n = \boldsymbol{L}^{-1/2}\boldsymbol{U}^T\left(\boldsymbol{x}_n - \overline{\boldsymbol{x}}\right)$$

❑ The set $\boldsymbol{y}_n$ has zero mean and its covariance is the identity:

$$\frac{1}{N}\sum_{n=1}^{N}\boldsymbol{y}_n\boldsymbol{y}_n^T = \frac{1}{N}\sum_{n=1}^{N}\boldsymbol{L}^{-1/2}\boldsymbol{U}^T\left(\boldsymbol{x}_n - \overline{\boldsymbol{x}}\right)\left(\boldsymbol{x}_n - \overline{\boldsymbol{x}}\right)^T\boldsymbol{U}\boldsymbol{L}^{-1/2} = \boldsymbol{L}^{-1/2}\boldsymbol{U}^T\boldsymbol{S}\boldsymbol{U}\boldsymbol{L}^{-1/2} = \boldsymbol{I}$$
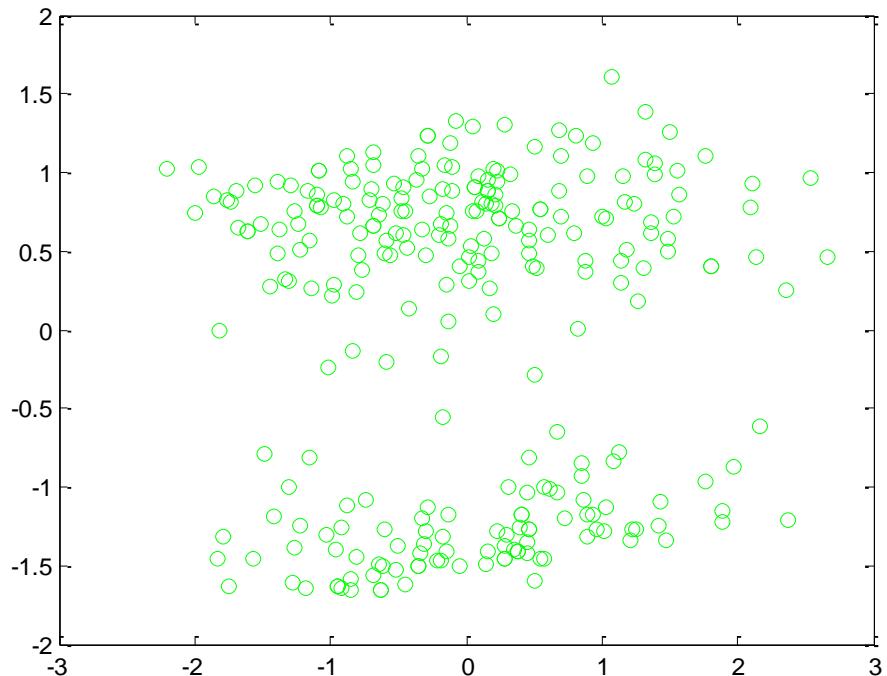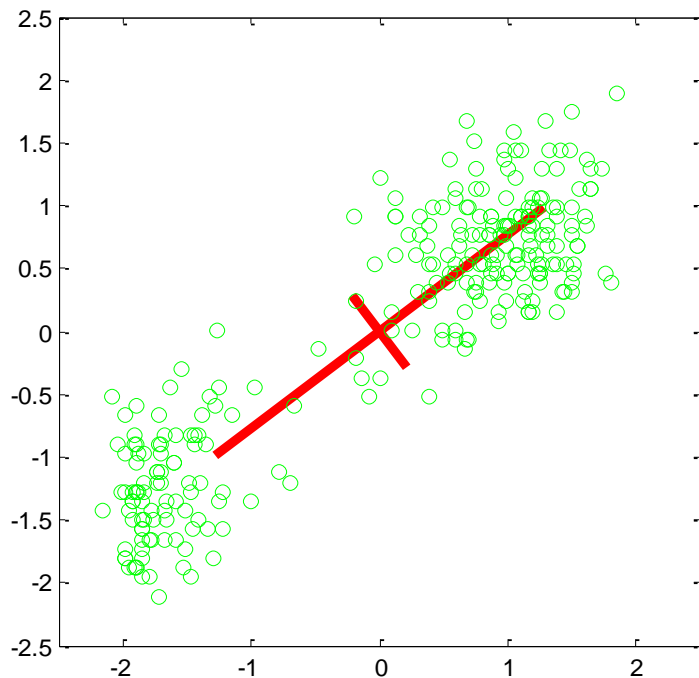
# *Old Faithful Data Set: PCA Example*

❑ 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

# *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 $N << D$, 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 $X$ to be the $N \times D$-dimensional centered data matrix, whose n<sup>th</sup> row is given by $\left( x_n - \overline{x} \right)^T$ .

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

$$S = \frac{1}{N} \sum_{n=1}^{N} \left( x_n - \overline{x} \right) \left( x_n - \overline{x} \right)^T = \frac{1}{N} X^T X$$

❑ The corresponding eigenvector equation is:

$$\frac{1}{N} X^T X u_i = \lambda_i u_i$$

❑ Pre-multipling both sides from the left by $X$ gives:

$$\frac{1}{N} X X^T \left( X u_i \right) = \lambda_i \left( X u_i \right)$$

# PCA For High Dimensional Data

$$\frac{1}{N} XX^T \left( Xu_i \right) = \lambda_i \left( Xu_i \right)$$

❑ Let us define:

$$v_i = Xu_i$$

❑ The eigenvalue problem becomes:

$$\frac{1}{N} XX^T v_i = \lambda_i v_i$$

❑ This is an eigenvector equ. for the $N \times N$ matrix $XX^T/N$

❑ 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 ($XX^T$ 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 $\boldsymbol{u}_i$, we multiply by $\boldsymbol{X}^T$:

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

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

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

❑ This approach is indeed simple:

- first evaluate $\boldsymbol{XX}^T$ 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 – 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 $z$.
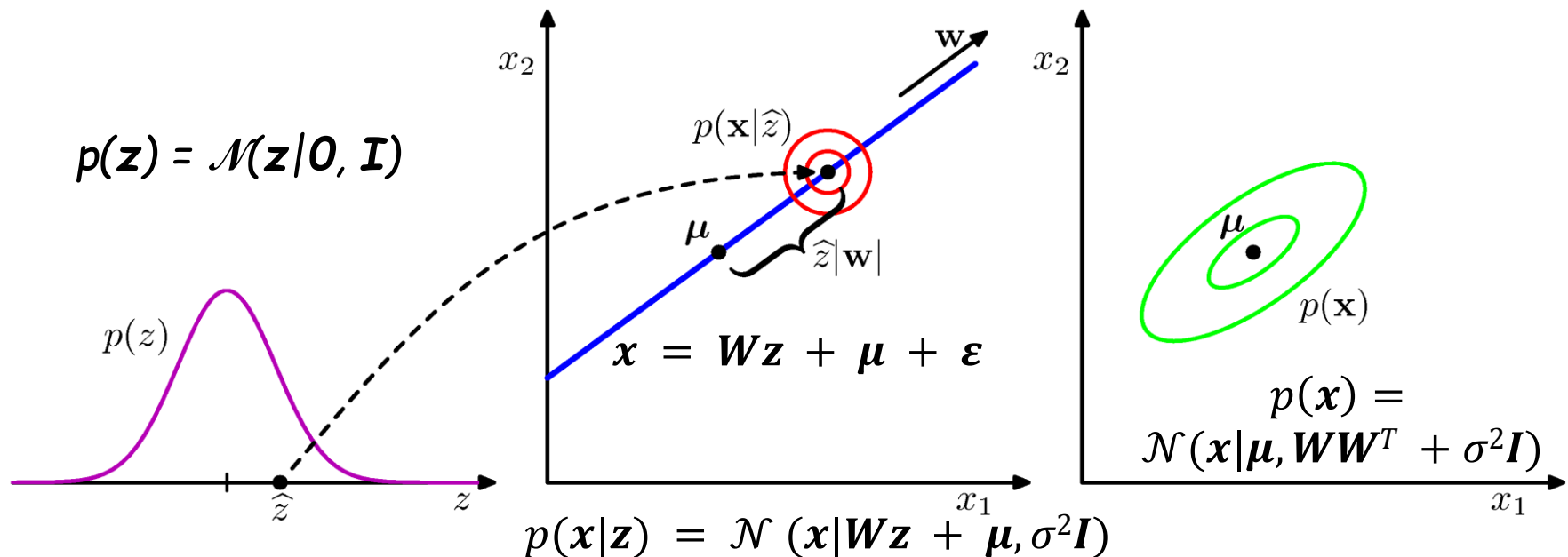
❑ The prior distribution over $z$ is given by

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

❑ 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 $\mu$.

# *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 $p(x)$.



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

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

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

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

# Probabilistic PCA – Predictive Distribution

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

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

$$\boldsymbol{C}^{-1} = \left(\boldsymbol{WW}^T + \sigma^2 \boldsymbol{I}\right)^{-1} = \sigma^{-2}\boldsymbol{I} - \sigma^{-2}\boldsymbol{WM}^{-1}\boldsymbol{W}^T$$

where

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

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

# Probabilistic PCA – Posterior Distribution

❑ The posterior distribution can be derived directly from earlier results on linear Gaussian models.

- We know $p(\boldsymbol{z}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I})$

- The conditional: $p(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \sigma^2 \boldsymbol{I})$

- From these we conclude:

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

Here we used results from an earlier lecture:

$$\left\{\begin{array}{l} p(\boldsymbol{x}) = \mathcal{N}\left(\boldsymbol{x} \mid \boldsymbol{\mu}, \Lambda^{-1}\right) \\[2mm] p(\boldsymbol{y} \mid \boldsymbol{x}) = \mathcal{N}\left(\boldsymbol{y} \mid \boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}, \boldsymbol{L}^{-1}\right) \end{array}\right. \Rightarrow$$

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

$$\Sigma = \left(\Lambda + \boldsymbol{A}^T\boldsymbol{L}\boldsymbol{A}\right)^{-1}$$

# Maximum Likelihood PCA

❑ Consider determining the model parameters using maximum likelihood. Using

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

we derive:

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

$$= -\frac{ND}{2}\ln(2\pi) - \frac{N}{2}\ln|\boldsymbol{C}| - \frac{1}{2}\sum_{n=1}^{N}(\boldsymbol{x}_n - \boldsymbol{\mu})^T \boldsymbol{C}^{-1}(\boldsymbol{x}_n - \boldsymbol{\mu})$$

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

$$\boldsymbol{\mu} = \sum_{n=1}^{N} \boldsymbol{x}_n \bigg/ N \equiv \overline{\boldsymbol{x}}$$

# *Maximum Likelihood PCA*

❑ The log-likelihood is then simplified as:

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

or

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

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

❑ Maximization wrt $W$ and $\sigma^2$ 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*

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

❑ $U_M$ is a $D \times M$ matrix whose columns are given by any subset (of size $M$) of the eigenvectors of the data covariance matrix $S$.

❑ $L_M$ is the $M \times M$ diagonal matrix with elements given by the corresponding eigenvalues $\lambda_i$ of $S$.

❑ $R$ 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 $\lambda_i$ values, the $M$ principal eigenvectors are $u_1, \ldots, u_M$. The columns of $W$ then define the principal subspace as in standard PCA.

# *Maximum Likelihood PCA*

❑ The corresponding MLE solution for $\sigma^2$ is <u>given as</u>:

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

 i.e. the average of the discarded eigenvalues.

❑ $\boldsymbol{R}$ in $\boldsymbol{W}_{ML} = \boldsymbol{U}_M \left( \boldsymbol{L}_M - \sigma^2 \boldsymbol{I} \right)^{1/2} \boldsymbol{R}$ is a rotation matrix in the $M$ dimensional latent space.

❑ Substituting this into the predictive variance $\boldsymbol{C} = \boldsymbol{W}_{ML} \boldsymbol{W}_{ML}^T + \sigma_{ML}^2 \boldsymbol{I}$ and using $\boldsymbol{R}^T \boldsymbol{R} = \boldsymbol{I}$, we see that $\boldsymbol{C}$ is independent of $\boldsymbol{R}$.

❑ The predictive density is unchanged by rotations in the latent space (statistical non-identifiability).

# *Maximum Likelihood PCA*

❑ Consider the variance of the predictive distribution $C = WW^T + \sigma^2 I$ along some direction defined by the unit vector $v$ given by $v^T C v$.

❑ Let $v$ be orthogonal to the principal subspace $U$, i.e. $v$ is given by some linear combination of the discarded eigenvectors.

❑ Then $v^T U = 0$ and hence

$$v^T C v = v^T \left( WW^T + \sigma^2 I \right) v = \sigma^2$$

❑ Thus the model predicts a noise variance orthogonal to the principal subspace which <u>it was shown </u>to be the average of the discarded eigenvalues.

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

# Maximum Likelihood PCA

❑ Consider the variance of the predictive distribution $C = WW^T + \sigma^2 I$ along some direction defined by the unit vector $v$ given by $v^T C v$.

❑ Let now $v = u_i$ where $u_i$ is one of the retained eigenvectors defining the principal subspace.

❑ Then using $W_{ML} = U_M \left( L_M - \sigma^2 I \right)^{1/2} R$, and $u_i^T u_j = 0, j = 1, \ldots, M$ we see that:

$$v^T C v = u_i^T \left( WW^T + \sigma^2 I \right) u_i = \left( \lambda_i - \sigma^2 \right) + \sigma^2 = \lambda_i$$

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 $\sigma^2$. Variance is `lost' in the projections.

# Maximum Likelihood PCA

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

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

❑ We construct the MLE model by solving the underlying eigenvalue problem for the data covariance matrix and evaluate $W$ and $\sigma^2$ from the Eqs. above. We choose $R = I$.

❑ Note: if the MLE model is found using optimization methods or via the EM algorithm (see following slides), $R$ is arbitrary and the columns of $W$ 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 vs Standard PCA*

❑ In the limit $\sigma^2 \to 0$, the posterior mean becomes:

$$\mathbb{E}[z \mid x] = M^{-1} W_{ML}^T (x - \overline{x}), \quad M = W_{ML}^T W_{ML} + \sigma^2 I \to W_{ML}^T W_{ML} \Rightarrow$$

$$\mathbb{E}[z \mid x] = \left( W_{ML}^T W_{ML} \right)^{-1} W_{ML}^T (x - \overline{x})$$

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

$$\mathbb{E}[z \mid x] = \left( W_{ML}^T W_{ML} \right)^{-1} W_{ML}^T (x - \overline{x}) = L_M^{-1/2} U_M^T (x - \overline{x})$$

❑ This is an orthogonal projection of the data point onto the latent space, i.e. for the limit $\sigma^2 \to 0$, we recover the standard PCA model. The posterior covariance
$\mathrm{var}[z \mid x] = \sigma^2 M^{-1} \to 0$ 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*

❑ 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 $\boldsymbol{C}$ depends on $\boldsymbol{W}$ ($D \times M$), and σ²: $\boldsymbol{C} = \boldsymbol{W}\boldsymbol{W}^T + \sigma^2 \boldsymbol{I}$ total parameters $DM + 1$.

  ▪ We need to subtract the redundancy associated with rotations of the coordinate system in the latent space.

  ▪ $\boldsymbol{R}$ 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. $\boldsymbol{R}$ has a total of $M(M-1)/2$ independent parameters.

❑ The number of degrees of freedom in $\boldsymbol{C}$ grows linearly with $D$

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

# 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.

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

❑ The n[th] row of $\boldsymbol{Z}$ is given by $\boldsymbol{z}_n$.

❑ Recall that $p(\boldsymbol{z}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{0},\boldsymbol{I})$, $p(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{Wz} + \boldsymbol{\mu}, \sigma^2\boldsymbol{I})$

❑ We can now write the expectation with respect to the posterior distribution over the latent variables.

# EM Algorithms for PCA

$$\mathbb{E}\left[\ln p\left(\boldsymbol{X}, \boldsymbol{Z} \mid \boldsymbol{\mu}, \boldsymbol{W}, \sigma^2\right)\right] = -\sum_{n=1}^{N}\left\{\begin{array}{l}\dfrac{D}{2}\ln\left(2\pi\sigma^2\right) + \dfrac{1}{2}Tr\left(\mathbb{E}\left[z_n z_n^T\right]\right) \\[2mm] +\dfrac{1}{2\sigma^2}\left\|x_n - \boldsymbol{\mu}\right\|^2 - \dfrac{1}{\sigma^2}\mathbb{E}\left[z_n\right]^T \boldsymbol{W}^T\left(x_n - \boldsymbol{\mu}\right) \\[2mm] +\dfrac{1}{2\sigma^2}Tr\left(\mathbb{E}\left[z_n z_n^T\right]\boldsymbol{W}^T \boldsymbol{W}\right) + \dfrac{M}{2\ln\left(2\pi\right)}\end{array}\right\}$$

❑ E-Step: We use the old parameters to evaluate:

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

❑ This follows directly from

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

together with the standard result

$$\mathbb{E}\left[z_n z_n^T\right] = \mathrm{cov}\left[z_n\right] + \mathbb{E}\left[z_n\right]\mathbb{E}\left[z_n\right]^T$$

# EM Algorithms for PCA

$$\mathbb{E}\left[\ln p\left(\boldsymbol{X}, \boldsymbol{Z} \mid \boldsymbol{\mu}, \boldsymbol{W}, \sigma^2\right)\right] = -\sum_{n=1}^{N} \left\{ \begin{array}{l} \dfrac{D}{2}\ln\left(2\pi\sigma^2\right) + \dfrac{1}{2}Tr\left(\mathbb{E}\left[\boldsymbol{z}_n \boldsymbol{z}_n^T\right]\right) \\ + \dfrac{1}{2\sigma^2}\|\boldsymbol{x}_n - \boldsymbol{\mu}\|^2 - \dfrac{1}{\sigma^2}\mathbb{E}\left[\boldsymbol{z}_n\right]^T \boldsymbol{W}^T\left(\boldsymbol{x}_n - \boldsymbol{\mu}\right) \\ + \dfrac{1}{2\sigma^2}Tr\left(\mathbb{E}\left[\boldsymbol{z}_n \boldsymbol{z}_n^T\right]\boldsymbol{W}^T\boldsymbol{W}\right) + \dfrac{M}{2\ln\left(2\pi\right)} \end{array} \right\}$$

❑ **M-Step:** We maximize with respect to $\boldsymbol{W}$ and $\sigma^2$ keeping the posterior statistics fixed. We obtain:

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

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

# *Proof of the M-Step Equations*

$$\mathbb{E}\left[\ln p\left(\boldsymbol{X},\boldsymbol{Z}\mid\boldsymbol{\mu},\boldsymbol{W},\sigma^2\right)\right]=-\sum_{n=1}^{N}\left\{\begin{array}{l}\dfrac{D}{2}\ln\left(2\pi\sigma^2\right)+\dfrac{1}{2}Tr\left(\mathbb{E}\left[\boldsymbol{z}_n\boldsymbol{z}_n^T\right]\right)\\[4mm]+\dfrac{1}{2\sigma^2}\|\boldsymbol{x}_n-\boldsymbol{\mu}\|^2-\dfrac{1}{\sigma^2}\mathbb{E}\left[\boldsymbol{z}_n\right]^T\boldsymbol{W}^T\left(\boldsymbol{x}_n-\boldsymbol{\mu}\right)\\[4mm]+\dfrac{1}{2\sigma^2}Tr\left(\mathbb{E}\left[\boldsymbol{z}_n\boldsymbol{z}_n^T\right]\boldsymbol{W}^T\boldsymbol{W}\right)+\dfrac{M}{2\ln\left(2\pi\right)}\end{array}\right\}$$

❑ The two M-equations are derived by setting the derivatives wrt $\boldsymbol{W}$ and $\sigma^2$ equal to zero:

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

$$\frac{\partial\mathbb{E}\left[\ln p\left(\boldsymbol{X},\boldsymbol{Z}\mid\boldsymbol{\mu},\boldsymbol{W},\sigma^2\right)\right]}{\partial\sigma^2}=\sum_{n=1}^{N}\left\{\begin{array}{l}-\dfrac{D}{2\sigma^2}-\dfrac{1}{\sigma^4}\mathbb{E}\left[\boldsymbol{z}_n\right]^T\boldsymbol{W}^T\left(\boldsymbol{x}_n-\boldsymbol{\mu}\right)+\dfrac{1}{2\sigma^4}\|\boldsymbol{x}_n-\boldsymbol{\mu}\|^2\\[4mm]+\dfrac{1}{2\sigma^4}Tr\left(\mathbb{E}\left[\boldsymbol{z}_n\boldsymbol{z}_n^T\right]\boldsymbol{W}^T\boldsymbol{W}\right)\end{array}\right\}=0$$

❑ Here we used $\dfrac{\partial}{\partial\boldsymbol{A}}Tr(\boldsymbol{A}\boldsymbol{B}\boldsymbol{A}^T)=\boldsymbol{A}(\boldsymbol{B}+\boldsymbol{B}^T),\ \dfrac{\partial}{\partial\boldsymbol{A}}Tr(\boldsymbol{A}\boldsymbol{B})=\boldsymbol{B}^T$

# EM Algorithms for PCA

❑ Initialize the parameters

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

$$\mathbb{E}\left[z_n\right] = M^{-1}W^T\left(x_n - \overline{x}\right) \qquad \mathbb{E}\left[z_n z_n^T\right] = \sigma^2 M^{-1} + \mathbb{E}\left[z_n\right]\mathbb{E}\left[z_n^T\right]$$

❑ Revise the parameter values in the M-Step.

$$W_{new} = \left[\sum_{n=1}^{N}\left(x_n - \overline{x}\right)\mathbb{E}\left[z_n\right]^T\right]\left[\sum_{n=1}^{N}\mathbb{E}\left[z_n z_n^T\right]\right]^{-1}$$

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

# 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

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

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: 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 $\sigma^2 \to 0$*

❑ When $\sigma^2 \to 0$, EM corresponds to standard PCA ([Roweis, 1998](#))

❑ Defining $\widetilde{X}$ a matrix of size $N \times D$ whose $n^{th}$ row is given by $x_n - \overline{x}$.

❑ Defining a matrix $\Omega$ of size $M \times N$ whose $n^{th}$ column is given by the vector $\mathbb{E}[z_n]$.

❑ The E-Step for $\sigma^2 \to 0$ becomes

$$\mathbb{E}\left[ z_n \right] = M^{-1} W^T \left( x_n - \overline{x} \right) \Rightarrow \quad \Omega = \left( W_{old}^T W_{old} \right)^{-1} W_{old}^T \widetilde{X}^T$$

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 $\sigma^2 \to 0$

❑ When $\sigma^2 \to 0$, EM corresponds to standard PCA

❑ $\widetilde{X}$ an $N \times D$ matrix whose $n^{th}$ row is given by $x_n - \overline{x}$.

❑ $\Omega$ a matrix $M \times N$ whose $n^{th}$ column is given by $\mathbb{E}[z_n]$.

❑ Noting that

$$\sigma^2 \to 0 \Rightarrow \mathbb{E}\left[z_n z_n^T\right] = \sigma^2 M^{-1} + \mathbb{E}\left[z_n\right]\mathbb{E}\left[z_n\right]^T \to \mathbb{E}\left[z_n\right]\mathbb{E}\left[z_n\right]^T$$

the M-Step takes the form:

$$W_{new} = \left[\sum_{n=1}^{N}\left(x_n - \overline{x}\right)\mathbb{E}\left[z_n\right]^T\right]\left[\sum_{n=1}^{N}\mathbb{E}\left[z_n z_n^T\right]\right]^{-1} \Rightarrow W_{new} = \widetilde{X}^T \Omega^T (\Omega \Omega^T)^{-1}$$

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.



$$\boldsymbol{\Omega} = \left(\boldsymbol{W}_{old}^T \boldsymbol{W}_{old}\right)^{-1} \boldsymbol{W}_{old}^T \widetilde{\boldsymbol{X}}^T$$

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

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

# 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 for PCA: Example, D=2, M=1



Run pcaEmStepByStep
From PMTK3

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



E step 4

M step 4

Run pcaEmStepByStep
From PMTK3

# *Bayesian PCA*

❑ It involves a specific choice of prior over $W$ that allows surplus dimensions in the principal subspace to be pruned out of the model (automatic relevance determination, or ARD)

❑ Specifically, we define an independent Gaussian prior over each column of $W$, which represent the vectors defining the principal subspace.

❑ Each such Gaussian has an independent variance governed by a precision hyperparameter $a_i$ where $w_i$ is the i<sup>th</sup> column of $W$.

$$p(W \mid a) = \prod_{i=1}^{M} \left( \frac{\alpha_i}{2\pi} \right)^{D/2} \exp\left\{ -\frac{1}{2} \alpha_i w_i^T w_i \right\}$$

# Bayesian PCA: Effective Dimensionality

❑ Find $a_i$ iteratively by maximizing the marginal likelihood function in which $W$ has been integrated out.

$$p(X \mid \boldsymbol{\alpha}, \mu, \sigma^2) = \int p(X \mid W, \mu, \sigma^2) \, p(W \mid \boldsymbol{\alpha}) \, dW$$

❑ In the optimization, some of the $a_i \to \infty$, with the corresponding $w_i \to 0$ (the posterior distribution becomes a delta function at the origin) giving a sparse solution.

❑ The effective dimensionality of the principal subspace is then the number of finite $a_i$'s. The corresponding $w_i$ can be thought of as relevant for modeling the data distribution.

❑ The Bayesian approach is making a trade-off between improving the fit to the data by using a larger number of $w_i$'s with their corresponding eigenvalues $\lambda_i$ each tuned to the data, and reducing the complexity of the model by suppressing some of the $w_i$'s.

# *Maximizing the Marginal Likelihood*

❑ The values of $a_i$ are re-estimated during training by maximizing the marginal likelihood

$$p(X \mid \boldsymbol{\alpha}, \mu, \sigma^2) = \int p(X \mid W, \mu, \sigma^2) p(W \mid \boldsymbol{\alpha}) dW$$

where

$$\ln p\left(X \mid \boldsymbol{\mu}, W, \sigma^2\right) = -\frac{ND}{2}\ln(2\pi) - \frac{N}{2}\ln|C| - \frac{1}{2}\sum_{n=1}^{N}(x_n - \boldsymbol{\mu})^T C^{-1}(x_n - \boldsymbol{\mu})$$

with

$$C = WW^T + \sigma^2 I$$

For simplicity $\mu$ and $\sigma^2$ are treated as parameters to be estimated, rather than defining priors over these parameters.

# Bayesian PCA: Re-estimation Equations

$$p(\boldsymbol{X} \mid \boldsymbol{\alpha}, \mu, \sigma^2) = \int p(\boldsymbol{X} \mid \boldsymbol{W}, \mu, \sigma^2) \, p(\boldsymbol{W} \mid \boldsymbol{\alpha}) d\boldsymbol{W}$$

❑ Because this integration is intractable, we make use of the Laplace approximation.

❑ Assume that the posterior distribution is sharply peaked, as will occur for sufficiently large data sets, then the re-estimation equations obtained by maximizing the marginal likelihood (see earlier lecture notes on model evidence) with respect to $\alpha_i$ take the simple form

$$\alpha_i^{new} = \frac{D}{w_i^T w_i}$$

noting that the dimensionality of $\boldsymbol{w}_i$ is $D$.

# Bayesian PCA: EM Implementation

❑ These re estimations are interleaved with the EM algorithm updates for determining $W$ and $\sigma^2$.

❑ The E-Step equations are again given as

$$\mathbb{E}\left[z_n\right] = M^{-1}W^T\left(x_n - \overline{x}\right) \qquad \mathbb{E}\left[z_n z_n^T\right] = \sigma^2 M^{-1} + \mathbb{E}\left[z_n\right]\mathbb{E}\left[z_n\right]^T$$

❑ Similarly, the M-Step equation for $\sigma^2$ is again given by

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

❑ The only change is to the M-Step equation for $W$:

$$W_{new} = \left[\sum_{n=1}^{N}\left(x_n - \overline{x}\right)\mathbb{E}\left[z_n\right]^T\right]\left[\sum_{n=1}^{N}\mathbb{E}\left[z_n z_n^T\right] + \sigma^2 A\right]^{-1}, \quad A = diag\left(\alpha_i\right)$$

❑ The value of $\mu$ is given by the sample mean, as before.

# *Bayesian PCA*

❑ For $M = D - 1$, if all $a_i$'s are finite, the model represents a full-covariance Gaussian.

❑ For $M = D - 1$, if all $a_i$ → ∞ the model is equivalent to an isotropic Gaussian, and so the model can encompass all permissible values for the effective dimensionality of the principal subspace.

❑ Using smaller values of M saves on computational cost but limits the maximum dimensionality of the subspace.

❑ A comparison of the Bayesian PCA algorithm with standard probabilistic PCA is shown next.

# Bayesian PCA: Hinton Diagrams of W

❑ 'Hinton' diagrams of the matrix $W$. Each element of the matrix is depicted as a square (white for positive & black for negative values) whose area is proportional to the magnitude of that element.



$W$:10×9                    $W$:10×3                    Matlab Implementation

❑ The synthetic data set comprises 300 data points in $D =$ 10 dimensions sampled from a Gaussian distribution having std 1.0 in 3 directions and std 0.5 in the remaining 7 directions.

# Bayesian PCA: Hinton Diagrams for W

❑ Left-hand plot: result from maximum likelihood probabilistic PCA.

❑ Right-hand plot: corresponding result from Bayesian PCA.

❑ The Bayesian model discovers the appropriate dimensionality by suppressing the 6 surplus degrees of freedom.

Matlab Implementation

# Kernel PCA

❑ **Kernel trick:** take an algorithm expressed in terms of scalar products of the form $x^T x$ and generalize it by replacing the scalar products with a nonlinear kernel.

❑ Apply this technique of kernel substitution to PCA, obtaining a nonlinear algorithm called kernel PCA (Scholkopf et al. 1998)

❑ Consider a data set $\{x_n\}, n = 1, \ldots, N$, of dimension $D$. Assume that we have already subtracted the sample mean from each $x_n$, so that $\sum_n x_n = 0$ .

❑ Start by expressing PCA in such a form that the data $\{x_n\}$ appear only in the form of the scalar products $x_n^T x_m$.

▪ Scholkopf, B., A. Smola, and K.-R. Muller (1998). Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation* **10**(5), 1299–1319.

# *Kernel PCA*

❑ Recall that the principal components are defined by the normalized eigenvectors $\boldsymbol{u}_i$ of the $N \times N$ covariance matrix

$$\boldsymbol{S}\boldsymbol{u}_i = \lambda_i \boldsymbol{u}_i, \ \text{where}: S = \frac{1}{N}\sum_{n=1}^{N} \boldsymbol{x}_n \boldsymbol{x}_n^T$$

❑ Consider a nonlinear transformation $\boldsymbol{\phi}(\boldsymbol{x})$ into an $M$-dimensional feature space, so that each data point $\boldsymbol{x}_n$ is thereby projected onto a point $\boldsymbol{\phi}(\boldsymbol{x}_n)$.

❑ Now perform standard PCA in the feature space.

  ▪ This implicitly defines a nonlinear principal component model in the original data space.

# Kernel PCA

❑ A data set in the original data space (left) is projected by $\boldsymbol{\phi}(\boldsymbol{x})$ into a feature space (right).

❑ By performing PCA in the feature space, we obtain the principal components ($\boldsymbol{v}_1$ being the 1st one).

# *Kernel PCA*

❑ The green lines in feature space indicate the linear projections onto $\boldsymbol{v}_1$, which correspond to nonlinear projections in the original data space.

❑ In general, it is not possible to represent the nonlinear principal component by a vector in $\boldsymbol{x}$ space.

# Kernel PCA

❑ For the moment, let us assume that the projected data set also has zero mean, so that $\sum_{n=1}^{N} \phi(x_n) = 0$. We shall return to this point shortly.

❑ The $M \times M$ sample covariance matrix in feature space is given by

$$C = \frac{1}{N} \sum_{n=1}^{N} \phi(x_n) \phi(x_n)^T$$

❑ The eigenvector expansion is defined as:

$$Cv_i = \lambda_i v_i, \, i = 1, ..., M$$

❑ Can we solve this eigenvalue problem without having to work explicitly in the feature space?

# Kernel PCA

$$C = \frac{1}{N} \sum_{n=1}^{N} \boldsymbol{\phi}(\boldsymbol{x}_n)\boldsymbol{\phi}(\boldsymbol{x}_n)^T \qquad\qquad C\boldsymbol{v}_i = \lambda_i \boldsymbol{v}_i, \, i = 1,...,M$$

❑ From the definition of $\boldsymbol{C}$, the eigenvector equations tells us that $\boldsymbol{v}_i$ satisfies

$$\frac{1}{N} \sum_{n=1}^{N} \boldsymbol{\phi}(\boldsymbol{x}_n)\left\{\boldsymbol{\phi}(\boldsymbol{x}_n)^T \boldsymbol{v}_i\right\} = \lambda_i \boldsymbol{v}_i, \, i = 1,...,M$$

❑ We see that (provided $\lambda_i > 0$), the vector $v_i$ is given by a linear combination of the $\boldsymbol{\phi}(\boldsymbol{x}_n)$ and so can be written in the form

$$\boldsymbol{v}_i = \sum_{n=1}^{N} a_{in}\boldsymbol{\phi}(\boldsymbol{x}_n), \, i = 1,...,M$$

# Kernel PCA

$$\frac{1}{N}\sum_{n=1}^{N}\boldsymbol{\phi}(\boldsymbol{x}_n)\left\{\boldsymbol{\phi}(\boldsymbol{x}_n)^T\boldsymbol{v}_i\right\}=\lambda_i\boldsymbol{v}_i,\, i=1,...,M$$

❑ Substituting

$$\boldsymbol{v}_i=\sum_{n=1}^{N}a_{in}\boldsymbol{\phi}(\boldsymbol{x}_n),\, i=1,...,M$$

back into the eigenvector equation, we obtain

$$\left[\frac{1}{N}\sum_{n=1}^{N}\boldsymbol{\phi}(\boldsymbol{x}_n)\boldsymbol{\phi}(\boldsymbol{x}_n)^T\right]\sum_{m=1}^{N}a_{im}\boldsymbol{\phi}(\boldsymbol{x}_m)=\lambda_i\sum_{n=1}^{N}a_{in}\boldsymbol{\phi}(\boldsymbol{x}_n),\, i=1,...,M$$

❑ The key step is now to express this in terms of the kernel function

$$K_{nm}\equiv k\left(\boldsymbol{x}_n,\boldsymbol{x}_m\right)=\boldsymbol{\phi}(\boldsymbol{x}_n)^T\boldsymbol{\phi}(\boldsymbol{x}_m)$$

which we do by multiplying both sides by $\boldsymbol{\phi}(\boldsymbol{x}_l)^T$.

# Kernel PCA

$$\frac{1}{N}\sum_{n=1}^{N}\boldsymbol{\phi}(\boldsymbol{x}_n)\boldsymbol{\phi}(\boldsymbol{x}_n)^T\sum_{m=1}^{N}a_{im}\boldsymbol{\phi}(\boldsymbol{x}_m)=\lambda_i\sum_{n=1}^{N}a_{in}\boldsymbol{\phi}(\boldsymbol{x}_n),\, i=1,...,M$$

❑ Multiplying both sides by $\boldsymbol{\phi}(\boldsymbol{x}_l)^T$ to give

$$\frac{1}{N}\sum_{n=1}^{N}\boldsymbol{\phi}(\boldsymbol{x}_l)^T\boldsymbol{\phi}(\boldsymbol{x}_n)\boldsymbol{\phi}(\boldsymbol{x}_n)^T\sum_{m=1}^{N}a_{im}\boldsymbol{\phi}(\boldsymbol{x}_m)=\lambda_i\sum_{n=1}^{N}a_{in}\boldsymbol{\phi}(\boldsymbol{x}_l)^T\boldsymbol{\phi}(\boldsymbol{x}_n),\, i=1,...,M \Rightarrow$$

$$\frac{1}{N}\sum_{n=1}^{N}k\left(\boldsymbol{x}_l,\boldsymbol{x}_n\right)\sum_{m=1}^{N}a_{im}k\left(\boldsymbol{x}_n,\boldsymbol{x}_m\right)=\lambda_i\sum_{n=1}^{N}a_{in}k\left(\boldsymbol{x}_l,\boldsymbol{x}_n\right),\, i=1,...,M \Rightarrow$$

❑ This can be written in a matrix form as:

$$\boldsymbol{K}^2\boldsymbol{a}_i=\lambda_i N\boldsymbol{K}\boldsymbol{a}_i$$

$\boldsymbol{a}_i$ $N$-dimensional column vector with elements $a_{in}, n = 1,..,N$.

❑ Can simplify $\boldsymbol{K}\boldsymbol{a}_i = \lambda_i N \boldsymbol{a}_i$ . These Eqs. differ by eigen-vectors of $\boldsymbol{K}$ with $0$ eigenvalues that don't affect the principal components projection.

# Kernel PCA

$$\mathbf{K}\mathbf{a}_i = \lambda_i N\mathbf{a}_i$$

❑ We can find solutions for $\mathbf{a}_i$ by solving the above eigenvalue problem.

❑ The normalization condition for the coefficients $\mathbf{a}_i$ is obtained by requiring that the eigenvectors in feature space be normalized.

❑ Using

and

$$\mathbf{v}_i = \sum_{n=1}^{N} a_{in}\boldsymbol{\phi}(\mathbf{x}_n), \; i = 1,...,M$$

we have:

$$\mathbf{K}\mathbf{a}_i = \lambda_i N\mathbf{a}_i$$

$$1 = \mathbf{v}_i^T\mathbf{v}_i = \frac{1}{N}\sum_{n=1}^{N}\sum_{m=1}^{N} a_{in}a_{im}\boldsymbol{\phi}(\mathbf{x}_n)^T\boldsymbol{\phi}(\mathbf{x}_m) = \mathbf{a}_i^T\mathbf{K}\mathbf{a}_i = \lambda_i N\mathbf{a}_i^T\mathbf{a}_i \implies \boxed{1 = \lambda_i N\mathbf{a}_i^T\mathbf{a}_i}$$

# Kernel PCA

❑ The principal component projections can now be cast in terms of the kernel function.

❑ Using

$$v_i = \sum_{n=1}^{N} a_{in} \phi(x_n), \, i = 1, ..., M$$

the projection of a point $x$ onto eigenvector $i$ is given by

$$y_i(x) = \phi(x)^T v_i = \sum_{n=1}^{N} a_{in} \phi(x)^T \phi(x_n) = \sum_{n=1}^{N} a_{in} k(x, x_n)$$

# *Kernel PCA*

❑ In the original $D$-dimensional $\boldsymbol{x}$ space: $D$ orthogonal eigenvectors and thus at most $D$ linear principal components.

❑ The dimensionality $M$ of the feature space, however, can be much larger than $D$ (even infinite), and thus we can find a number of nonlinear principal components that can exceed $D$.

❑ Note, however, that the number of nonzero eigenvalues cannot exceed the number $N$ of data points:

  ▪ The covariance matrix in feature space has rank at most equal to $N$.

  ▪ Recall that kernel PCA involves the eigenvector expansion of the $N \times N$ matrix $\boldsymbol{K}$.

# Kernel PCA

❑ So far we have assumed that the projected data set given by $\phi(x_n)$ has zero mean, which in general will not be the case.

❑ We cannot simply compute and then subtract off the mean, since we want to formulate the algorithm purely in terms of the kernel function.

❑ The projected data points after centralizing are given by

$$\phi(x_n) = \phi(x_n) - \frac{1}{N} \sum_{l=1}^{N} \phi(x_l)$$

❑ The corresponding elements of the Gram matrix $\widetilde{K}_{nm} = \widetilde{\phi}(x_n)^T \widetilde{\phi}(x_m)$ can be computed as shown next.

# Kernel PCA

$$\widetilde{K}_{nm} = \widetilde{\boldsymbol{\phi}}(x_n)^T \widetilde{\boldsymbol{\phi}}(x_m) = \left( \boldsymbol{\phi}(\boldsymbol{x}_n) - \frac{1}{N} \sum_{l=1}^{N} \boldsymbol{\phi}(\boldsymbol{x}_l) \right)^T \left( \boldsymbol{\phi}(\boldsymbol{x}_m) - \frac{1}{N} \sum_{l=1}^{N} \boldsymbol{\phi}(\boldsymbol{x}_l) \right)$$

$$= \boldsymbol{\phi}(\boldsymbol{x}_n)^T \boldsymbol{\phi}(\boldsymbol{x}_m) - \frac{1}{N} \sum_{l=1}^{N} \boldsymbol{\phi}(\boldsymbol{x}_n)^T \boldsymbol{\phi}(\boldsymbol{x}_l)$$

$$- \frac{1}{N} \sum_{l=1}^{N} \boldsymbol{\phi}(\boldsymbol{x}_l)^T \boldsymbol{\phi}(\boldsymbol{x}_m) + \frac{1}{N^2} \sum_{j=1}^{N} \sum_{l=1}^{N} \boldsymbol{\phi}(\boldsymbol{x}_j)^T \boldsymbol{\phi}(\boldsymbol{x}_l)$$

$$= k(\boldsymbol{x}_n, \boldsymbol{x}_m) - \frac{1}{N} \sum_{l=1}^{N} k(\boldsymbol{x}_n, \boldsymbol{x}_l) - \frac{1}{N} \sum_{l=1}^{N} k(\boldsymbol{x}_l, \boldsymbol{x}_m) + \frac{1}{N^2} \sum_{j=1}^{N} \sum_{l=1}^{N} k(\boldsymbol{x}_j, \boldsymbol{x}_l)$$

❑ We can express this in matrix notation as:

$$\widetilde{K} = K - \mathbf{1}_N K - K \mathbf{1}_N + \mathbf{1}_N K \mathbf{1}_N$$

where $\mathbf{1}_N$ denotes the $N \times N$ matrix in which every element takes the value $1/N$.

# Kernel PCA

$$\widetilde{K} = K - 1_N K - K 1_N + 1_N K 1_N$$

❑ Thus we can evaluate $\widetilde{K}$ using only the kernel function and then use $\widetilde{K}$ to determine the eigenvalues and eigenvectors.

❑ Note that the standard PCA algorithm is recovered as a special case if we use a linear kernel $k(x, x') = x^T x'$.

❑ We often use Gaussian kernels of the form:

$$k(x, x') = \exp\left( -\frac{\|x - x'\|}{2\sigma^2} \right)$$

# *Kernel PCA*

❑ KPCA, with a Gaussian kernel ($2\sigma^2 = 0.1$) applied to a synthetic data in 2D, shows the first 8 eigenfunctions and eigenvalues.

❑ The lines along which the projection onto the corresponding principal component is constant are shown:

$$\phi(x)^T v_i = \sum_{n=1}^{N} a_{in} k(x, x_n)$$

Eigenvalue=21.72    Eigenvalue=21.65    Eigenvalue=4.11    Eigenvalue=3.93



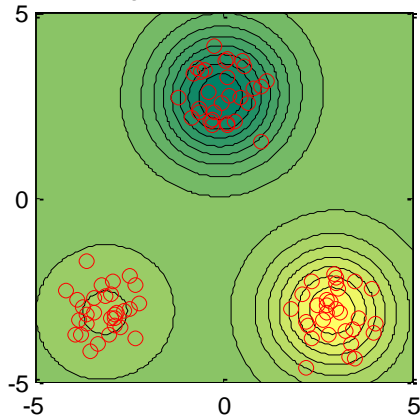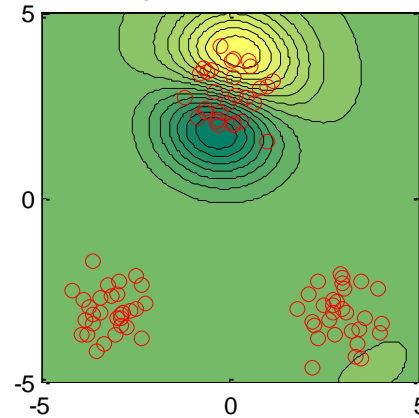Eigenvalue=3.66    Eigenvalue=3.09    Eigenvalue=2.60    Eigenvalue=2.53

# *Kernel PCA*

❑ Note the first 2 eigenvectors separate the 3 clusters, the next three eigenvectors split each of the cluster into halves, and the following three eigenvectors again split the clusters into halves along directions orthogonal to the previous splits.



Eigenvalue=21.72

Eigenvalue=21.65

Eigenvalue=4.11

Eigenvalue=3.93

Eigenvalue=3.66

Eigenvalue=3.09

Eigenvalue=2.60

Eigenvalue=2.53

Matlab Implementation

# Kernel PCA



Matlab Implementation

# Kernel PCA: Disadvatages

❑ KPCA involves finding the eigenvectors of the $N \times N$ matrix $\widetilde{K}$ rather than the $D \times D$ matrix $S$ of PCA.

▪ In practice, for large $N$ approximations are used.

❑ In PCA, we retain some reduced number $L < D$ of eigenvectors and then approximate a data vector $x_n$ by its projection $\hat{x}_n$ onto the $L$-dimensional principal subspace, defined by

$$\hat{x}_n = \sum_{i=1}^{L} (x_n^T u_i) \, u_i$$

❑ In KPCA, this in general is not possible.

# *Kernel PCA: Preimage problem*

$$\widehat{x}_n = \sum_{i=1}^{L} (x_n^T u_i)\, u_i$$

❑ This is not possible in KPCA. Indeed note that the mapping $\phi(x)$ maps the $D$-dimensional $x$ space into a $D$-dimensional manifold in the $M$-dimensional feature space $\phi$.

❑ The vector $x$ is known as the pre-image of $\phi(x)$. The projection of points in feature space onto the linear PCA subspace in that space will typically not lie on the nonlinear $D$-dimensional manifold and so will not have a corresponding pre-image in data space.

❑ Techniques have been proposed for finding approximate pre-images (Bakir et al., 2004)

▪ Bakir, G. H., J. Weston, and B. Scholkopf (2004). Learning to find pre-images. In S. Thrun, L. K. Saul, and B. Scholkopf (Eds.), *Advances in Neural Information Processing Systems*, Volume 16, pp. 449–456. MIT Press.