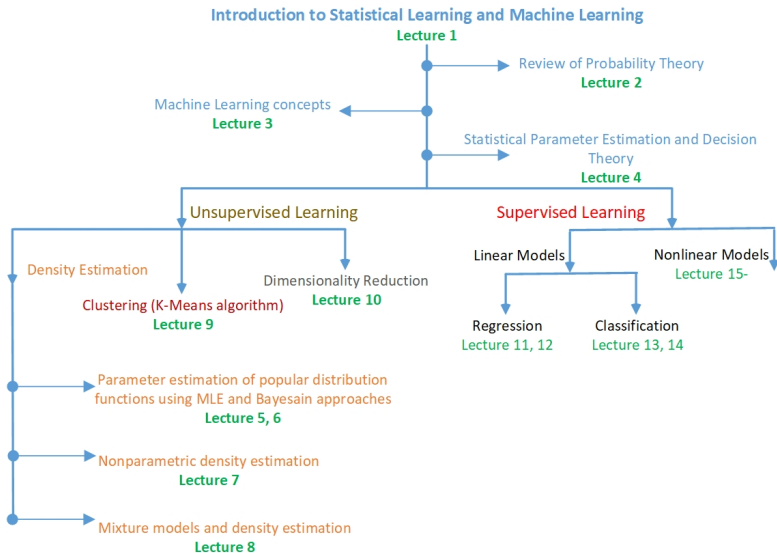


# Statistical Learning and Machine Learning

## Lecture 10 - Principal Component Analysis

September 18, 2021

# Course overview and where do we stand



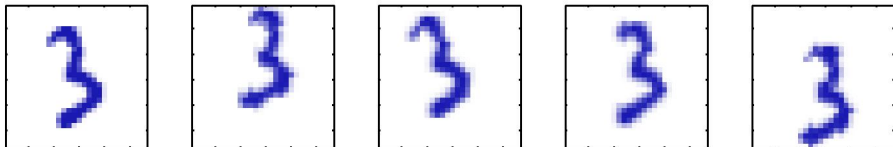
# Main topics of the lecture

- Principal Component Analysis (PCA) and applications
- Minimum Variance Formulation of PCA
- Probabilistic PCA

# Intrinsic dimensionality of data

Data forming data sets usually lie close to a manifold of much lower dimensionality than that of the original data space. In this example:

- translation
- rotation
- zero-padding



# Principal Component Analysis

PCA is a technique for determining a linear data transformation (data projection) of the form:

$$y = U^T x. \quad (18)$$

It is widely used for:

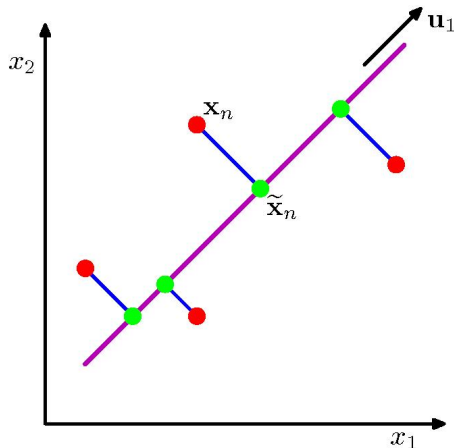
- dimensionality reduction
- lossy data compression
- feature extraction
- data visualization.

There are two ways to define it as:

- the data projection preserving the maximum variance of the data set
- the data projection minimizing the average transformation cost.

# Principal Component Analysis

The space on which the data are projected is called principal subspace. In this example the 2D data points are projected to the magenta line.



# PCA: Maximum variance formulation

Let  $\mathcal{D} = \{x_1, \dots, x_N\}$  be a set of  $N$  data points  $x_n \in \mathbb{R}^D$ . Our goal is to:

- project the data onto a space with dimensionality  $M < D$  (we assume that  $M$  is given)
- maximize the variance of the projected data.

# PCA: Maximum variance formulation ( $M = 1$ )

We use a vector  $u_1 \in \mathbb{R}^D$  to project the vector  $x_n \in \mathbb{R}^D$  to a value  $y_n$ :

- we are interested only in the direction of the line on which  $y_n, n = 1, \dots, N$  lie
- we choose a unit vector:  $u_1^T u_1 = 1$ .

The mean value and the variance of the projected data are:

$$\bar{y} = \frac{1}{N} \sum_{n=1}^N u_1^T x_n = u_1^T \bar{x} \quad (19)$$

$$\text{var}[y] = \frac{1}{N} \sum_{n=1}^N (u_1^T x_n - u_1^T \bar{x})^2 = u_1^T S u_1 \quad (20)$$

where  $S$  is the covariance matrix of the original data:

$$S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T \quad (21)$$



# PCA: Maximum variance formulation ( $M = 1$ )

We want to maximize  $\text{var}[y]$  under the constraint that  $\mathbf{u}_1^T \mathbf{u}_1 = 1$ . To do so, we introduce a Lagrange multiplier and maximize for:

$$\mathcal{J}_{PCA} = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1). \quad (22)$$

Setting  $\frac{\partial \mathcal{J}_{PCA}}{\partial \mathbf{u}_1} = 0$ :

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1 \quad (23)$$

Thus,  $\mathbf{u}_1$  is an eigenvector of  $\mathbf{S}$ .

By multiplying both sides of (23) with  $\mathbf{u}_1^T$  we get:  $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1$

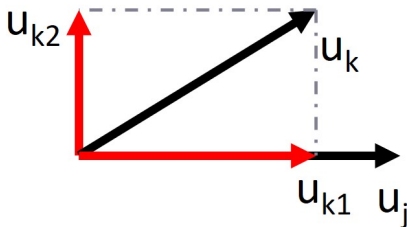
To maximize the variance, we set  $\mathbf{u}_1$  equal to the eigenvector of  $\mathbf{S}$  corresponding to the largest eigenvalue  $\lambda_1$ .

$\mathbf{u}_1$  is called the *first principal component*.

# PCA: Maximum variance formulation ( $M > 1$ )

We define additional principal components in an incremental fashion:

- $u_1$  is the first principal component of  $S$
- $u_k$ ,  $k > 1$ , such that
  - $u_k^T u_k = 1$
  - $u_k^T u_j = 0$ ,  $\forall j < k$ .



The data projection matrix is formed by the eigenvectors corresponding to the  $M$  largest eigenvalues of  $S$  ( $U = [u_1, \dots, u_M]$ ).

# PCA: Minimum-error formulation

We define a new coordinate system  $\{\mathbf{u}_j\}$ ,  $j = 1, \dots, D$  that satisfy  $\mathbf{u}_j^T \mathbf{u}_j = 1$ , and  $\mathbf{u}_i^T \mathbf{u}_j = 0$  for  $i \neq j$ .

The data points  $\mathbf{x}_n$ ,  $n = 1, \dots, N$  can be expressed in the new coordinate system by:

$$\mathbf{x}_n = \sum_{i=1}^D \alpha_{ni} \mathbf{u}_i \quad (24)$$

The above corresponds to:

- a coordinate change from  $\mathbf{x}_n \in \mathbb{R}^D$  to  $\boldsymbol{\alpha}_n \in \mathbb{R}^D$
- $\alpha_{nj} = \mathbf{u}_j^T \mathbf{x}_n$  being the dimension of the  $n$ -th data point in the new coordinate system

Thus the original data points can be reconstructed by:

$$\mathbf{x}_n = \sum_{i=1}^D (\mathbf{u}_i^T \mathbf{x}_n) \mathbf{u}_i \quad (25)$$

However, we want to approximate the data points using only  $M$  data point-dependent dimensions:

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i \quad (26)$$

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

We use the following cost function:

$$\begin{aligned}\mathcal{J} &= \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 \\ &= \frac{1}{N} \sum_{n=1}^N \left\| \mathbf{x}_n - \left( \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i \right) \right\|^2\end{aligned}\quad (27)$$

which is a function of  $\{\mathbf{u}_i\}$ ,  $\{z_{ni}\}$  and  $\{b_i\}$ ,  $n = 1, \dots, N$ ,  $i = 1, \dots, D$ .

# PCA: Minimum-error formulation

Minimization w.r.t.  $\{z_{nj}\}$ :

$$\frac{\partial \mathcal{J}}{\partial z_{nj}} = 0 \Rightarrow z_{nj} = \mathbf{x}_n^T \mathbf{u}_j, \quad j = 1, \dots, M \quad (28)$$

Minimization w.r.t.  $\{b_j\}$ :

$$\frac{\partial \mathcal{J}}{\partial b_j} = 0 \Rightarrow b_j = \bar{\mathbf{x}}^T \mathbf{u}_j, \quad j = M + 1, \dots, D \quad (29)$$

Substituting the above to  $\mathcal{J}$  and using  $\mathbf{x}_n = \sum_{i=1}^D (\mathbf{u}_i^T \mathbf{x}_n) \mathbf{u}_i$  we get:

$$\mathbf{x}_n - \tilde{\mathbf{x}}_n = \sum_{i=M+1}^D \{(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{u}_i\} \mathbf{u}_i \quad (30)$$

# PCA: Minimum-error formulation

Thus:

$$\mathcal{J} = \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D (\mathbf{x}_n^T \mathbf{u}_i - \bar{\mathbf{x}}^T \mathbf{u}_i)^2 = \sum_{i=M+1}^D \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i \quad (31)$$

For  $\mathbf{u}_j$ ,  $j = M+1, \dots, D$ , we want to minimize  $\mathcal{J}$  under the constraint that  $\mathbf{u}_j^T \mathbf{u}_j = 1$ . To do so, we introduce a Lagrange multiplier and minimize for:

$$\tilde{\mathcal{J}} = \mathbf{u}_j^T \mathbf{S} \mathbf{u}_j + \lambda_j (1 - \mathbf{u}_j^T \mathbf{u}_j). \quad (32)$$

To minimize  $\tilde{\mathcal{J}}$ , we set  $\mathbf{u}_i$ ,  $i = M+1, \dots, D$  equal to the eigenvectors of  $\mathbf{S}$  corresponding to the smallest eigenvalues  $\lambda_i$ .

The distortion is equal to:  $\mathcal{J} = \sum_{i=M+1}^D \lambda_i$ .

# PCA for high-dimensional data

The application of eigen-analysis to the matrix  $S \in \mathbb{R}^{D \times D}$  has a time complexity of  $O(D^3)$ . When  $D$  is large, application of PCA is very slow.

Let  $X \in \mathbb{R}^{N \times D}$  be the (centered) data matrix. Then we have:

- $S = \frac{1}{N}X^T X$
- The original eigen-analysis problem is:  $\frac{1}{N}X^T X u_i = \lambda_i u_i$
- Multiply with  $X$  both sides:  $\frac{1}{N}X X^T (X u_i) = \lambda_i (X u_i)$
- Set  $(X u_i) = v_i$  and we get:  $\frac{1}{N}X X^T v_i = \lambda_i v_i$
- Solve the above problem which has time complexity  $O(N^3)$
- Calculate the original eigenvectors  $u_i$  by:  $u_i = \frac{1}{(N\lambda_i)^{1/2}} X^T v_i$

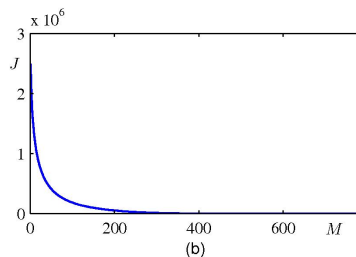
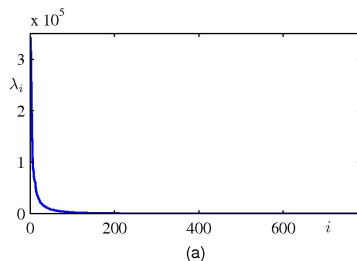
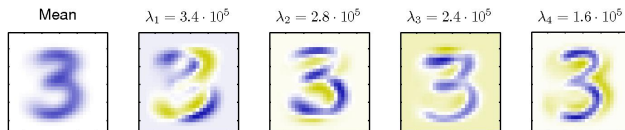
Thus, when  $N < D$  we can apply PCA faster (with time complexity  $O(N^3)$ ).



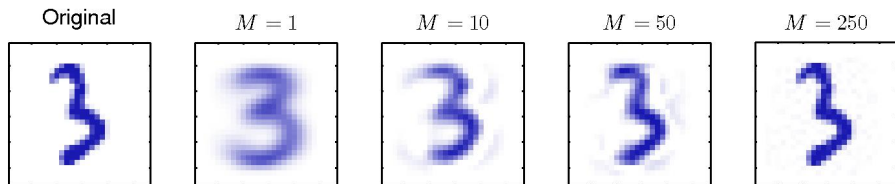
# PCA: Applications

Data compression:

$$\tilde{\mathbf{x}}_n = \hat{\mathbf{x}} + \sum_{i=1}^M (\mathbf{x}_n^T \mathbf{u}_i - \hat{\mathbf{x}}^T \mathbf{u}_i) \mathbf{u}_i \quad (33)$$



Data reconstruction for a varying number of principal components



Data pre-processing:

- when the data dimensions correspond to measurements of different types (leading to different scaling)
- normalization helps in appropriately scaling the various dimensions

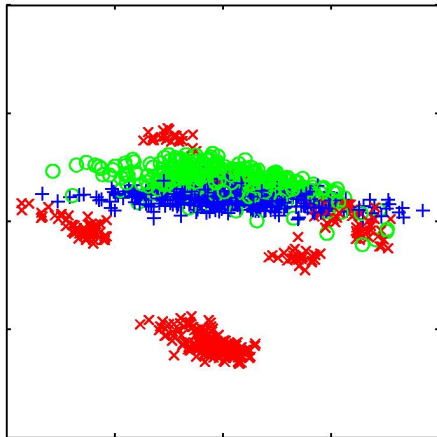
PCA-based normalization (leading to de-correlated dimensions)

- 1 Calculate the data covariance matrix  $S$
- 2 Perform eigen-analysis:  $SU = UL$
- 3 Transform the data:  $y_n = L^{-1/2}U^T(x_n - \bar{x})$
- 4 Now:  $\frac{1}{N} \sum_{n=1}^N y_n y_n^T = I$

The above data transformation is called *data whitening*.

# PCA: Applications

Data visualization: Project the data to 2D or 3D.



# Probabilistic PCA

We introduce variable  $z$  corresponding to the principal component subspace and the distributions:

$$p(z) = \mathcal{N}(z|0, I) \quad (34)$$

$$p(x|z) = \mathcal{N}(x|Wz + \mu, \sigma^2 I) \quad (35)$$

where  $W \in \mathbb{R}^{D \times M}$  and  $\mu \in \mathbb{R}^D$

The marginal distribution of  $x$  is  $p(x) = \mathcal{N}(x|\mu, C)$ , where:

$$\mu = \mathbb{E}[Wz + \mu + \epsilon] = \mathbb{E}[x] \quad (36)$$

$$\begin{aligned} C &= \mathbb{E}[(Wz + \epsilon)(Wz + \epsilon^T)] \\ &= \mathbb{E}[Wzz^T W^T] + \mathbb{E}[\epsilon\epsilon^T] = WW^T + \sigma^2 I \end{aligned} \quad (37)$$

# Probabilistic PCA: Maximum Likelihood

The parameters of the Probabilistic PCA model are:

$$\boldsymbol{\mu}, \quad \mathbf{W} \quad \text{and} \quad \sigma^2 \quad (38)$$

Given a set of data points  $\mathbf{X} = \{\mathbf{x}_n\}$  we want to estimate the above parameters.

The log-likelihood function is:

$$\begin{aligned} \ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \sum_{n=1}^N \ln p(\mathbf{x}_n|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{C}| \\ &\quad - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \end{aligned} \quad (39)$$

Setting  $\frac{\partial \ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2)}{\partial \boldsymbol{\mu}} = 0$ , we get  $\boldsymbol{\mu} = \bar{\mathbf{x}}$ .

# Probabilistic PCA: Maximum Likelihood

Maximization w.r.t.  $W$  and  $\sigma$  is more complex, but has an exact solution:

$$W_{ML} = U_M(L_M - \sigma^2 I)^{1/2} R \quad (40)$$

$$\sigma_{ML}^2 = \frac{1}{D - M} \sum_{i=M+1}^D \lambda_i \quad (41)$$

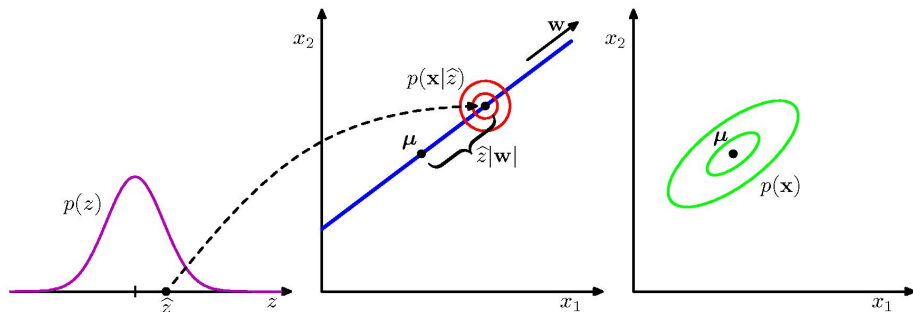
where:

- $U_M \in \mathbb{R}^{D \times M}$  is formed by the  $M$  eigenvectors corresponding to the largest eigenvalues  $\lambda_i$ ,  $i = 1, \dots, M$
- $L_M = \text{diag}(\lambda_i)$  is formed by the  $M$  largest eigenvalues
- $R \in \mathbb{R}^{M \times M}$  is an arbitrary orthogonal matrix.

# Probabilistic PCA: Generative model

Generate a data point  $\mathbf{x}$  by:

- 1 drawing a random value  $\hat{z}$  using  $p(z)$
- 2 using  $\hat{z}$ , draw a vector  $\mathbf{x}$  from an isotropic Gaussian  $\mathcal{N}(W\hat{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I})$





In cases of working on high-dimensional spaces ( $D$  is large) it may be easier to use an iterative method to estimate the parameters of PPCA.

The log-likelihood function for  $X = [x_1, \dots, x_N]$  and  $Z = [z_1, \dots, z_N]$  is:

$$\ln p(X, Z | \mu, W, \sigma^2) = \sum_{n=1}^N \{ \ln p(x_n | z_n) + \ln p(z_n) \} \quad (42)$$

When  $\sigma^2 \rightarrow 0$ , we obtain the standard PCA model.

**E step:**

$$\mathbb{E}[z_n] = M^{-1}W^T(x_n - \bar{x}) \quad (43)$$

$$\mathbb{E}[z_n z_n^T] = \sigma^2 M^{-1} + \mathbb{E}[z_n] \mathbb{E}[z_n]^T \quad (44)$$

where  $M = W^T W + \sigma^2 I$ .

**M step:**

$$W_{new} = \left[ \sum_{n=1}^N (x_n - \bar{x}) \mathbb{E}[z_n]^T \right] \left[ \sum_{n=1}^N \mathbb{E}[z_n z_n^T] \right]^{-1} \quad (45)$$

$$\begin{aligned} \sigma_{new}^2 &= \frac{1}{ND} \sum_{n=1}^N \{ \|x_n - \bar{x}_n\|^2 - 2 \mathbb{E}[z_n]^T W_{new} (x_n - \bar{x}) \\ &\quad + \text{Tr}(\mathbb{E}[z_n z_n^T] W_{new}^T W_{new}) \} \end{aligned} \quad (46)$$