

# Linear Dimensionality Reduction: Principal Component Analysis

Piyush Rai

Machine Learning (CS771A)

Sept 2, 2016

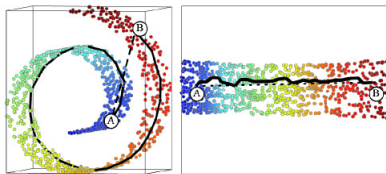
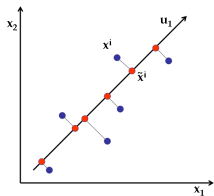
# Dimensionality Reduction

# Dimensionality Reduction

- Usually considered an **unsupervised learning** method

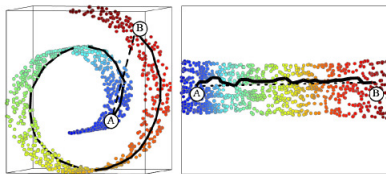
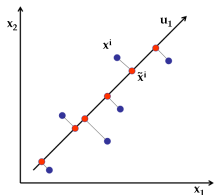
# Dimensionality Reduction

- Usually considered an **unsupervised learning** method
- Used for learning the **low-dimensional structures** in the data



# Dimensionality Reduction

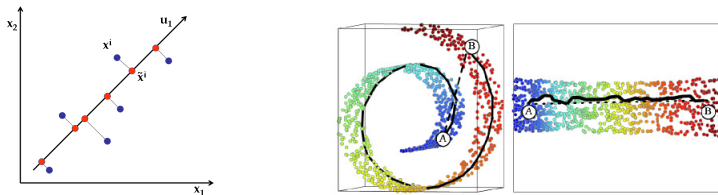
- Usually considered an **unsupervised learning** method
- Used for learning the **low-dimensional structures** in the data



- Also useful for **“feature learning”** or **“representation learning”** (learning a better, often smaller-dimensional, representation of the data), e.g.,

# Dimensionality Reduction

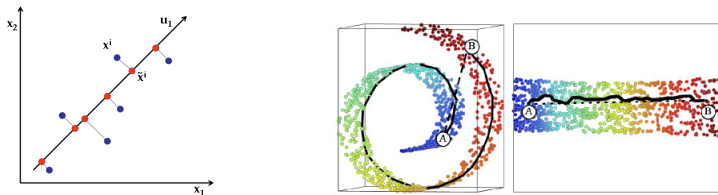
- Usually considered an **unsupervised learning** method
- Used for learning the **low-dimensional structures** in the data



- Also useful for “**feature learning**” or “**representation learning**” (learning a better, often smaller-dimensional, representation of the data), e.g.,
  - Documents using using topic vectors instead of bag-of-words vectors

# Dimensionality Reduction

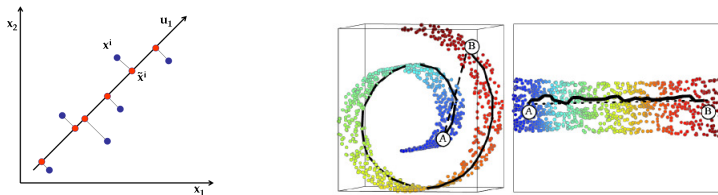
- Usually considered an **unsupervised learning** method
- Used for learning the **low-dimensional structures** in the data



- Also useful for “**feature learning**” or “**representation learning**” (learning a better, often smaller-dimensional, representation of the data), e.g.,
  - Documents using using topic vectors instead of bag-of-words vectors
  - Images using their constituent parts (faces - eigenfaces)

# Dimensionality Reduction

- Usually considered an **unsupervised learning** method
- Used for learning the **low-dimensional structures** in the data

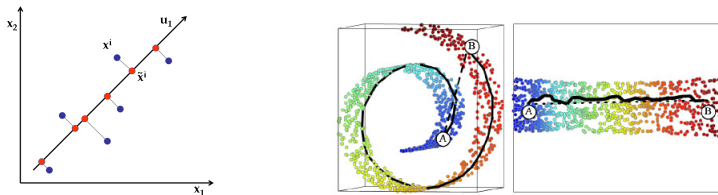


- Also useful for “**feature learning**” or “**representation learning**” (learning a better, often smaller-dimensional, representation of the data), e.g.,
  - Documents using using topic vectors instead of bag-of-words vectors
  - Images using their constituent parts (faces - eigenfaces)
- Can be used for **speeding up** learning algorithms



# Dimensionality Reduction

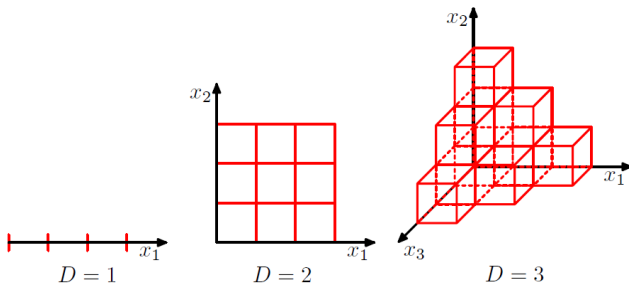
- Usually considered an **unsupervised learning** method
- Used for learning the **low-dimensional structures** in the data



- Also useful for “**feature learning**” or “**representation learning**” (learning a better, often smaller-dimensional, representation of the data), e.g.,
  - Documents using using topic vectors instead of bag-of-words vectors
  - Images using their constituent parts (faces - eigenfaces)
- Can be used for **speeding up** learning algorithms
- Can be used for **data compression**

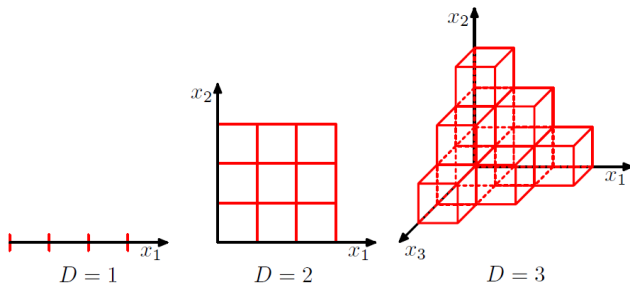
# Curse of Dimensionality

- Exponentially large # of examples required to “fill up” high-dim spaces



# Curse of Dimensionality

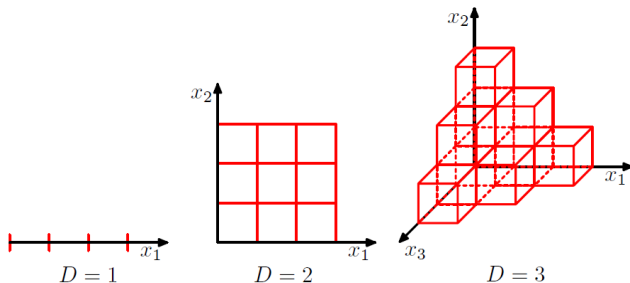
- Exponentially large # of examples required to “fill up” high-dim spaces



- Fewer dimensions  $\Rightarrow$  Less chances of overfitting  $\Rightarrow$  Better generalization

# Curse of Dimensionality

- Exponentially large # of examples required to “fill up” high-dim spaces



- Fewer dimensions  $\Rightarrow$  Less chances of overfitting  $\Rightarrow$  Better generalization
- Dimensionality reduction is a way to beat the curse of dimensionality

# Linear Dimensionality Reduction

- A **projection matrix**  $\mathbf{U} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_K]$  of size  $D \times K$  defines  $K$  **linear projection directions**, each  $\mathbf{u}_k \in \mathbb{R}^D$ , for the  $D$  dim. data (assume  $K < D$ )

# Linear Dimensionality Reduction

- A **projection matrix**  $\mathbf{U} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_K]$  of size  $D \times K$  defines  $K$  **linear projection directions**, each  $\mathbf{u}_k \in \mathbb{R}^D$ , for the  $D$  dim. data (assume  $K < D$ )
- Can use  $\mathbf{U}$  to transform  $\mathbf{x}_n \in \mathbb{R}^D$  into  $\mathbf{z}_n \in \mathbb{R}^K$  as shown below

The diagram illustrates the transformation of a data vector  $\mathbf{x}_n$  into a projected vector  $\mathbf{z}_n$  using the projection matrix  $\mathbf{U}^T$ . The dimensions of the vectors and matrix are indicated above them:

- $\mathbf{z}_n$  is a  $K \times 1$  vector, represented by a red vertical bar.
- $\mathbf{U}^T$  is a  $K \times D$  matrix, represented by a gray rectangle. It contains  $K$  rows, each representing the transpose of a projection direction  $\mathbf{u}_k$ , labeled  $\mathbf{u}_1^T$  through  $\mathbf{u}_K^T$ .
- $\mathbf{x}_n$  is a  $D \times 1$  vector, represented by a blue vertical bar.

The transformation is shown as a matrix multiplication:  $\mathbf{z}_n = \mathbf{U}^T \mathbf{x}_n$ .

# Linear Dimensionality Reduction

- A **projection matrix**  $\mathbf{U} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_K]$  of size  $D \times K$  defines  $K$  **linear projection directions**, each  $\mathbf{u}_k \in \mathbb{R}^D$ , for the  $D$  dim. data (assume  $K < D$ )
- Can use  $\mathbf{U}$  to transform  $\mathbf{x}_n \in \mathbb{R}^D$  into  $\mathbf{z}_n \in \mathbb{R}^K$  as shown below

Diagram illustrating the linear transformation  $\mathbf{z}_n = \mathbf{U}^T \mathbf{x}_n$ . The diagram shows a red vertical bar representing  $\mathbf{z}_n$  (size  $K \times 1$ ) equal to a gray matrix representing  $\mathbf{U}^T$  (size  $K \times D$ ) multiplied by a blue vertical bar representing  $\mathbf{x}_n$  (size  $D \times 1$ ). The matrix  $\mathbf{U}^T$  is shown with rows  $\mathbf{u}_1^T$  and  $\mathbf{u}_K^T$ .

- Note that  $\mathbf{z}_n = \mathbf{U}^T \mathbf{x}_n = [\mathbf{u}_1^T \mathbf{x}_n \ \mathbf{u}_2^T \mathbf{x}_n \ \dots \ \mathbf{u}_K^T \mathbf{x}_n]$  is a  $K$ -dim projection of  $\mathbf{x}_n$

# Linear Dimensionality Reduction

- A **projection matrix**  $\mathbf{U} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_K]$  of size  $D \times K$  defines  $K$  **linear projection directions**, each  $\mathbf{u}_k \in \mathbb{R}^D$ , for the  $D$  dim. data (assume  $K < D$ )
- Can use  $\mathbf{U}$  to transform  $\mathbf{x}_n \in \mathbb{R}^D$  into  $\mathbf{z}_n \in \mathbb{R}^K$  as shown below

The diagram illustrates the transformation of a data point  $\mathbf{x}_n$  into its projection  $\mathbf{z}_n$ . It shows a red vertical bar representing  $\mathbf{z}_n$  of size  $K \times 1$ , followed by an equals sign, a gray rectangular box representing the projection matrix  $\mathbf{U}^T$  of size  $K \times D$ , followed by an asterisk, and finally a blue vertical bar representing the data point  $\mathbf{x}_n$  of size  $D \times 1$ . The matrix  $\mathbf{U}^T$  is shown with rows labeled  $\mathbf{u}_1^T$  and  $\mathbf{u}_K^T$ .

- Note that  $\mathbf{z}_n = \mathbf{U}^T \mathbf{x}_n = [\mathbf{u}_1^T \mathbf{x}_n \ \mathbf{u}_2^T \mathbf{x}_n \ \dots \ \mathbf{u}_K^T \mathbf{x}_n]$  is a  $K$ -dim projection of  $\mathbf{x}_n$ 
  - $\mathbf{z}_n \in \mathbb{R}^K$  is also called low-dimensional **“embedding”** of  $\mathbf{x}_n \in \mathbb{R}^D$



# Linear Dimensionality Reduction

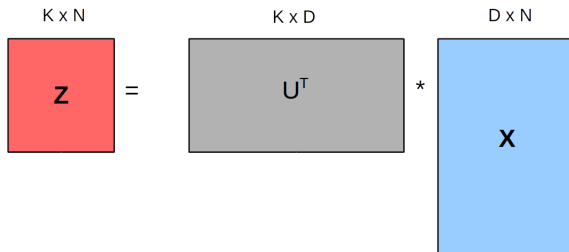
- $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$  is  $D \times N$  matrix denoting all the  $N$  data points

# Linear Dimensionality Reduction

- $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$  is  $D \times N$  matrix denoting all the  $N$  data points
- $\mathbf{Z} = [\mathbf{z}_1 \ \mathbf{z}_2 \ \dots \ \mathbf{z}_N]$  is  $K \times N$  matrix denoting **embeddings** of data points

# Linear Dimensionality Reduction

- $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$  is  $D \times N$  matrix denoting all the  $N$  data points
- $\mathbf{Z} = [\mathbf{z}_1 \ \mathbf{z}_2 \ \dots \ \mathbf{z}_N]$  is  $K \times N$  matrix denoting **embeddings** of data points
- With this notation, the figure on previous slide can be re-drawn as below



A diagram illustrating the matrix equation  $\mathbf{Z} = \mathbf{U}^T \mathbf{X}$ . It consists of three colored rectangles: a red rectangle on the left labeled  $\mathbf{Z}$  with dimensions  $K \times N$  above it; a gray rectangle in the middle labeled  $\mathbf{U}^T$  with dimensions  $K \times D$  above it; and a blue rectangle on the right labeled  $\mathbf{X}$  with dimensions  $D \times N$  above it. An equals sign is placed between the red and gray rectangles, and an asterisk is placed between the gray and blue rectangles.

# Linear Dimensionality Reduction

- $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$  is  $D \times N$  matrix denoting all the  $N$  data points
- $\mathbf{Z} = [\mathbf{z}_1 \ \mathbf{z}_2 \ \dots \ \mathbf{z}_N]$  is  $K \times N$  matrix denoting **embeddings** of data points
- With this notation, the figure on previous slide can be re-drawn as below

$$\begin{matrix} K \times N \\ \text{Z} \end{matrix} = \begin{matrix} K \times D \\ \mathbf{U}^T \end{matrix} * \begin{matrix} D \times N \\ \mathbf{X} \end{matrix}$$

- How do we learn the “best” projection matrix  $\mathbf{U}$ ?

# Linear Dimensionality Reduction

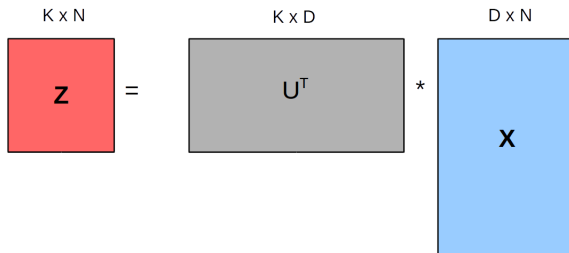
- $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$  is  $D \times N$  matrix denoting all the  $N$  data points
- $\mathbf{Z} = [\mathbf{z}_1 \ \mathbf{z}_2 \ \dots \ \mathbf{z}_N]$  is  $K \times N$  matrix denoting **embeddings** of data points
- With this notation, the figure on previous slide can be re-drawn as below

$$\begin{matrix} K \times N \\ \text{Z} \end{matrix} = \begin{matrix} K \times D \\ \text{U}^T \end{matrix} * \begin{matrix} D \times N \\ \text{X} \end{matrix}$$

- How do we learn the “best” projection matrix  $\mathbf{U}$ ?
- What criteria should we optimize for when learning  $\mathbf{U}$ ?

# Linear Dimensionality Reduction

- $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$  is  $D \times N$  matrix denoting all the  $N$  data points
- $\mathbf{Z} = [\mathbf{z}_1 \ \mathbf{z}_2 \ \dots \ \mathbf{z}_N]$  is  $K \times N$  matrix denoting **embeddings** of data points
- With this notation, the figure on previous slide can be re-drawn as below



A diagram illustrating the matrix equation  $\mathbf{Z} = \mathbf{U}^T \mathbf{X}$ . It consists of three colored rectangles: a red rectangle on the left labeled  $\mathbf{Z}$  with dimensions  $K \times N$  above it; a gray rectangle in the middle labeled  $\mathbf{U}^T$  with dimensions  $K \times D$  above it; and a blue rectangle on the right labeled  $\mathbf{X}$  with dimensions  $D \times N$  above it. An equals sign is placed between the red and gray rectangles, and an asterisk is placed between the gray and blue rectangles.

- How do we learn the “best” projection matrix  $\mathbf{U}$ ?
- What criteria should we optimize for when learning  $\mathbf{U}$ ?
- Principal Component Analysis (PCA) is an algorithm for doing this

# Principal Component Analysis

# Principal Component Analysis (PCA)

- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)



# Principal Component Analysis (PCA)

- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as

# Principal Component Analysis (PCA)

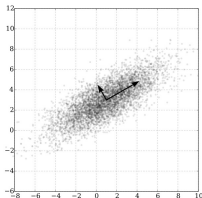
- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as
  - Learning projection directions that capture **maximum variance** in data

# Principal Component Analysis (PCA)

- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as
  - Learning projection directions that capture **maximum variance** in data
  - Learning projection directions that result in **smallest reconstruction error**

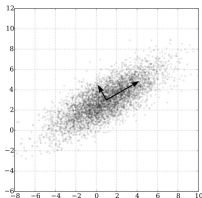
# Principal Component Analysis (PCA)

- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as
  - Learning projection directions that capture **maximum variance** in data
  - Learning projection directions that result in **smallest reconstruction error**
- Can also be seen as **changing the basis** in which the data is represented (and transforming the features such that **new features become decorrelated**)



# Principal Component Analysis (PCA)

- A classic linear dim. reduction method (Pearson, 1901; Hotelling, 1930)
- Can be seen as
  - Learning projection directions that capture **maximum variance** in data
  - Learning projection directions that result in **smallest reconstruction error**
- Can also be seen as **changing the basis** in which the data is represented (and transforming the features such that **new features become decorrelated**)

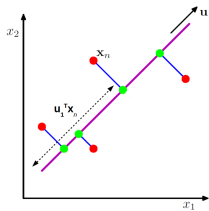


- Also related to other classic methods, e.g., **Factor Analysis** (Spearman, 1904)

# PCA as Maximizing Variance

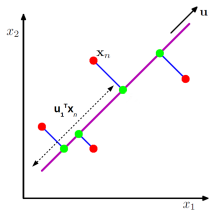
# Variance Captured by Projections

- Consider projecting  $\mathbf{x}_n \in \mathbb{R}^D$  on a **one-dim subspace** defined by  $\mathbf{u}_1 \in \mathbb{R}^D$
- Projection/embedding of  $\mathbf{x}_n$  along a one-dim subspace  $\mathbf{u}_1 = \mathbf{u}_1^\top \mathbf{x}_n$  (location of the green point along the purple line representing  $\mathbf{u}_1$ )



# Variance Captured by Projections

- Consider projecting  $\mathbf{x}_n \in \mathbb{R}^D$  on a **one-dim subspace** defined by  $\mathbf{u}_1 \in \mathbb{R}^D$
- Projection/embedding of  $\mathbf{x}_n$  along a one-dim subspace  $\mathbf{u}_1 = \mathbf{u}_1^\top \mathbf{x}_n$  (location of the green point along the purple line representing  $\mathbf{u}_1$ )

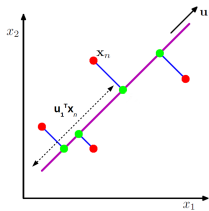


- Mean of projections of all the data:  $\frac{1}{N} \sum_{n=1}^N \mathbf{u}_1^\top \mathbf{x}_n = \mathbf{u}_1^\top (\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n) = \mathbf{u}_1^\top \boldsymbol{\mu}$



# Variance Captured by Projections

- Consider projecting  $\mathbf{x}_n \in \mathbb{R}^D$  on a **one-dim subspace** defined by  $\mathbf{u}_1 \in \mathbb{R}^D$
- Projection/embedding of  $\mathbf{x}_n$  along a one-dim subspace  $\mathbf{u}_1 = \mathbf{u}_1^\top \mathbf{x}_n$  (location of the green point along the purple line representing  $\mathbf{u}_1$ )

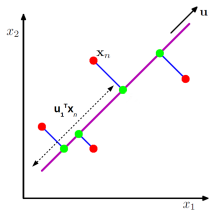


- Mean of projections of all the data:  $\frac{1}{N} \sum_{n=1}^N \mathbf{u}_1^\top \mathbf{x}_n = \mathbf{u}_1^\top (\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n) = \mathbf{u}_1^\top \boldsymbol{\mu}$
- Variance** of the projected data (“spread” of the green points)

$$\frac{1}{N} \sum_{n=1}^N (\mathbf{u}_1^\top \mathbf{x}_n - \mathbf{u}_1^\top \boldsymbol{\mu})^2 = \frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^\top (\mathbf{x}_n - \boldsymbol{\mu})\}^2 = \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$$

# Variance Captured by Projections

- Consider projecting  $\mathbf{x}_n \in \mathbb{R}^D$  on a **one-dim subspace** defined by  $\mathbf{u}_1 \in \mathbb{R}^D$
- Projection/embedding of  $\mathbf{x}_n$  along a one-dim subspace  $\mathbf{u}_1 = \mathbf{u}_1^\top \mathbf{x}_n$  (location of the green point along the purple line representing  $\mathbf{u}_1$ )

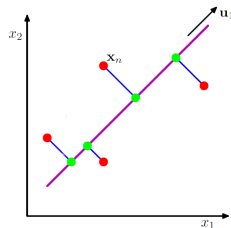


- Mean of projections of all the data:  $\frac{1}{N} \sum_{n=1}^N \mathbf{u}_1^\top \mathbf{x}_n = \mathbf{u}_1^\top (\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n) = \mathbf{u}_1^\top \boldsymbol{\mu}$
- Variance** of the projected data (“spread” of the green points)

$$\frac{1}{N} \sum_{n=1}^N (\mathbf{u}_1^\top \mathbf{x}_n - \mathbf{u}_1^\top \boldsymbol{\mu})^2 = \frac{1}{N} \sum_{n=1}^N \{\mathbf{u}_1^\top (\mathbf{x}_n - \boldsymbol{\mu})\}^2 = \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$$

- S** is the  $D \times D$  **data covariance matrix**:  $\mathbf{s} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top$ . If data already centered ( $\boldsymbol{\mu} = 0$ ) then  $\mathbf{s} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top = \frac{1}{N} \mathbf{X}^\top \mathbf{X}$

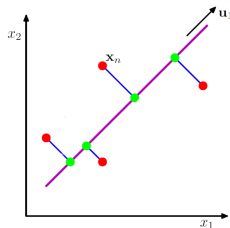
# Direction of Maximum Variance



- We want  $\mathbf{u}_1$  s.t. the variance of the projected data is maximized

$$\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$$

# Direction of Maximum Variance

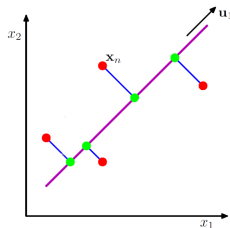


- We want  $\mathbf{u}_1$  s.t. the variance of the projected data is maximized

$$\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$$

- To prevent trivial solution (max var. = infinite), assume  $\|\mathbf{u}_1\| = 1 = \mathbf{u}_1^\top \mathbf{u}_1$

# Direction of Maximum Variance



- We want  $\mathbf{u}_1$  s.t. the variance of the projected data is maximized

$$\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$$

- To prevent trivial solution (max var. = infinite), assume  $\|\mathbf{u}_1\| = 1 = \mathbf{u}_1^\top \mathbf{u}_1$
- We will find  $\mathbf{u}_1$  by solving the following constrained opt. problem

$$\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^\top \mathbf{u}_1)$$

where  $\lambda_1$  is a Lagrange multiplier

# Direction of Maximum Variance

- The objective function:  $\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^\top \mathbf{u}_1)$
- Taking the derivative w.r.t.  $\mathbf{u}_1$  and setting to zero gives

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

# Direction of Maximum Variance

- The objective function:  $\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^\top \mathbf{u}_1)$
- Taking the derivative w.r.t.  $\mathbf{u}_1$  and setting to zero gives

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

- Thus  $\mathbf{u}_1$  is an eigenvector of  $\mathbf{S}$  (with corresponding eigenvalue  $\lambda_1$ )

# Direction of Maximum Variance

- The objective function:  $\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^\top \mathbf{u}_1)$
- Taking the derivative w.r.t.  $\mathbf{u}_1$  and setting to zero gives

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

- Thus  $\mathbf{u}_1$  is an eigenvector of  $\mathbf{S}$  (with corresponding eigenvalue  $\lambda_1$ )
- But which of  $\mathbf{S}$ 's ( $D$  possible) eigenvectors it is?



# Direction of Maximum Variance

- The objective function:  $\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^\top \mathbf{u}_1)$
- Taking the derivative w.r.t.  $\mathbf{u}_1$  and setting to zero gives

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

- Thus  $\mathbf{u}_1$  is an eigenvector of  $\mathbf{S}$  (with corresponding eigenvalue  $\lambda_1$ )
- But which of  $\mathbf{S}$ 's ( $D$  possible) eigenvectors it is?
- Note that since  $\mathbf{u}_1^\top \mathbf{u}_1 = 1$ , the variance of projected data is

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

# Direction of Maximum Variance

- The objective function:  $\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^\top \mathbf{u}_1)$
- Taking the derivative w.r.t.  $\mathbf{u}_1$  and setting to zero gives

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

- Thus  $\mathbf{u}_1$  is an eigenvector of  $\mathbf{S}$  (with corresponding eigenvalue  $\lambda_1$ )
- But which of  $\mathbf{S}$ 's ( $D$  possible) eigenvectors it is?
- Note that since  $\mathbf{u}_1^\top \mathbf{u}_1 = 1$ , the variance of projected data is

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

- Var. is maximized when  $\mathbf{u}_1$  is the (top) eigenvector with largest eigenvalue

# Direction of Maximum Variance

- The objective function:  $\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^\top \mathbf{u}_1)$

- Taking the derivative w.r.t.  $\mathbf{u}_1$  and setting to zero gives

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

- Thus  $\mathbf{u}_1$  is an eigenvector of  $\mathbf{S}$  (with corresponding eigenvalue  $\lambda_1$ )

- But which of  $\mathbf{S}$ 's ( $D$  possible) eigenvectors it is?

- Note that since  $\mathbf{u}_1^\top \mathbf{u}_1 = 1$ , the variance of projected data is

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

- Var. is maximized when  $\mathbf{u}_1$  is the (top) eigenvector with largest eigenvalue
- The top eigenvector  $\mathbf{u}_1$  is also known as the first Principal Component (PC)

# Direction of Maximum Variance

- The objective function:  $\arg \max_{\mathbf{u}_1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1(1 - \mathbf{u}_1^\top \mathbf{u}_1)$
- Taking the derivative w.r.t.  $\mathbf{u}_1$  and setting to zero gives

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

- Thus  $\mathbf{u}_1$  is an eigenvector of  $\mathbf{S}$  (with corresponding eigenvalue  $\lambda_1$ )
- But which of  $\mathbf{S}$ 's ( $D$  possible) eigenvectors it is?
- Note that since  $\mathbf{u}_1^\top \mathbf{u}_1 = 1$ , the variance of projected data is

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

- Var. is maximized when  $\mathbf{u}_1$  is the (top) eigenvector with largest eigenvalue
- The top eigenvector  $\mathbf{u}_1$  is also known as the first Principal Component (PC)
- Other directions can also be found likewise (with each being orthogonal to all previous ones) using the eigendecomposition of  $\mathbf{S}$  (this is PCA)

# Principal Component Analysis

- Steps in Principal Component Analysis

# Principal Component Analysis

- **Steps in Principal Component Analysis**

- Center the data (subtract the mean  $\mu = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$  from each data point)

# Principal Component Analysis

- **Steps in Principal Component Analysis**

- Center the data (subtract the mean  $\mu = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$  from each data point)
- Compute the covariance matrix  $\mathbf{S}$  using the centered data as

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^T \quad (\text{note: } \mathbf{X} \text{ assumed } D \times N \text{ here})$$

# Principal Component Analysis

- **Steps in Principal Component Analysis**

- Center the data (subtract the mean  $\mu = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$  from each data point)
- Compute the covariance matrix  $\mathbf{S}$  using the centered data as

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^T \quad (\text{note: } \mathbf{X} \text{ assumed } D \times N \text{ here})$$

- Do an eigendecomposition of the covariance matrix  $\mathbf{S}$



# Principal Component Analysis

- **Steps in Principal Component Analysis**

- Center the data (subtract the mean  $\mu = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$  from each data point)
- Compute the covariance matrix  $\mathbf{S}$  using the centered data as

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^T \quad (\text{note: } \mathbf{X} \text{ assumed } D \times N \text{ here})$$

- Do an eigendecomposition of the covariance matrix  $\mathbf{S}$
- Take first  $K$  leading eigenvectors  $\{\mathbf{u}_k\}_{k=1}^K$  with eigenvalues  $\{\lambda_k\}_{k=1}^K$

# Principal Component Analysis

- **Steps in Principal Component Analysis**

- Center the data (subtract the mean  $\mu = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$  from each data point)
- Compute the covariance matrix  $\mathbf{S}$  using the centered data as

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top \quad (\text{note: } \mathbf{X} \text{ assumed } D \times N \text{ here})$$

- Do an eigendecomposition of the covariance matrix  $\mathbf{S}$
- Take first  $K$  leading eigenvectors  $\{\mathbf{u}_k\}_{k=1}^K$  with eigenvalues  $\{\lambda_k\}_{k=1}^K$
- The final  $K$  dim. projection/embedding of data is given by

$$\mathbf{Z} = \mathbf{U}^\top \mathbf{X}$$

where  $\mathbf{U} = [\mathbf{u}_1 \ \dots \ \mathbf{u}_K]$  is  $D \times K$  and embedding matrix  $\mathbf{Z}$  is  $K \times N$

# Principal Component Analysis

- **Steps in Principal Component Analysis**

- **Center the data** (subtract the mean  $\mu = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$  from each data point)
- Compute the covariance matrix  $\mathbf{S}$  using the centered data as

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top \quad (\text{note: } \mathbf{X} \text{ assumed } D \times N \text{ here})$$

- Do an eigendecomposition of the covariance matrix  $\mathbf{S}$
- **Take first  $K$  leading eigenvectors  $\{\mathbf{u}_k\}_{k=1}^K$  with eigenvalues  $\{\lambda_k\}_{k=1}^K$**
- The final  $K$  dim. projection/embedding of data is given by

$$\mathbf{Z} = \mathbf{U}^\top \mathbf{X}$$

where  $\mathbf{U} = [\mathbf{u}_1 \ \dots \ \mathbf{u}_K]$  is  $D \times K$  and embedding matrix  $\mathbf{Z}$  is  $K \times N$

- **A word about notation:** If  $\mathbf{X}$  is  $N \times D$ , then  $\mathbf{S} = \frac{1}{N} \mathbf{X}^\top \mathbf{X}$  (needs to be  $D \times D$ ) and the embedding will be computed as  $\mathbf{Z} = \mathbf{X} \mathbf{U}$  where  $\mathbf{Z}$  is  $N \times K$

# PCA as Minimizing the Reconstruction Error

# Data as Combination of Basis Vectors

- Assume *complete* orthonormal basis vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_D$ , each  $\mathbf{u}_d \in \mathbb{R}^D$

# Data as Combination of Basis Vectors

- Assume *complete* orthonormal basis vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_D$ , each  $\mathbf{u}_d \in \mathbb{R}^D$
- We can represent each data point  $\mathbf{x}_n \in \mathbb{R}^D$  **exactly** using this new basis

$$\mathbf{x}_n = \sum_{k=1}^D z_{nk} \mathbf{u}_k$$

# Data as Combination of Basis Vectors

- Assume *complete* orthonormal basis vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_D$ , each  $\mathbf{u}_d \in \mathbb{R}^D$
- We can represent each data point  $\mathbf{x}_n \in \mathbb{R}^D$  **exactly** using this new basis

$$\mathbf{x}_n = \sum_{k=1}^D z_{nk} \mathbf{u}_k$$
$$\begin{bmatrix} x_{n1} \\ x_{n2} \\ \vdots \\ x_{nD} \end{bmatrix} = \begin{bmatrix} | & | \\ \mathbf{u}_1 & \mathbf{u}_2 \\ | & | \end{bmatrix} \begin{bmatrix} | \\ \mathbf{u}_D \\ | \end{bmatrix} * \begin{bmatrix} z_{n1} \\ z_{n2} \\ \vdots \\ z_{nD} \end{bmatrix}$$

# Data as Combination of Basis Vectors

- Assume *complete* orthonormal basis vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_D$ , each  $\mathbf{u}_d \in \mathbb{R}^D$
- We can represent each data point  $\mathbf{x}_n \in \mathbb{R}^D$  **exactly** using this new basis

$$\mathbf{x}_n = \sum_{k=1}^D z_{nk} \mathbf{u}_k$$
$$\begin{bmatrix} x_{n1} \\ x_{n2} \\ \vdots \\ x_{nD} \end{bmatrix} = \begin{bmatrix} | & | \\ \mathbf{u}_1 & \mathbf{u}_2 \\ | & | \end{bmatrix} \begin{bmatrix} | \\ \mathbf{u}_D \\ | \end{bmatrix} * \begin{bmatrix} z_{n1} \\ z_{n2} \\ \vdots \\ z_{nD} \end{bmatrix}$$

- Denoting  $\mathbf{z}_n = [z_{n1} \ z_{n2} \ \dots \ z_{nD}]^\top$ ,  $\mathbf{U} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_D]$ , and using  $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_D$

$$\mathbf{x}_n = \mathbf{U} \mathbf{z}_n \quad \text{and} \quad \mathbf{z}_n = \mathbf{U}^\top \mathbf{x}_n$$



# Data as Combination of Basis Vectors

- Assume *complete* orthonormal basis vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_D$ , each  $\mathbf{u}_d \in \mathbb{R}^D$
- We can represent each data point  $\mathbf{x}_n \in \mathbb{R}^D$  **exactly** using this new basis

$$\mathbf{x}_n = \sum_{k=1}^D z_{nk} \mathbf{u}_k$$
$$\begin{bmatrix} x_{n1} \\ x_{n2} \\ \vdots \\ x_{nD} \end{bmatrix} = \begin{bmatrix} | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & & \mathbf{u}_D \\ | & | & & | \end{bmatrix} * \begin{bmatrix} z_{n1} \\ z_{n2} \\ \vdots \\ z_{nD} \end{bmatrix}$$

- Denoting  $\mathbf{z}_n = [z_{n1} \ z_{n2} \ \dots \ z_{nD}]^\top$ ,  $\mathbf{U} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_D]$ , and using  $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_D$

$$\mathbf{x}_n = \mathbf{U} \mathbf{z}_n \quad \text{and} \quad \mathbf{z}_n = \mathbf{U}^\top \mathbf{x}_n$$

- Also note that each component of vector  $\mathbf{z}_n$  is  $z_{nk} = \mathbf{u}_k^\top \mathbf{x}_n$

# Reconstruction of Data from Projections

- Reconstruction of  $\mathbf{x}_n$  from  $\mathbf{z}_n$  will be **exact** if we use all the  $D$  basis vectors

# Reconstruction of Data from Projections

- Reconstruction of  $\mathbf{x}_n$  from  $\mathbf{z}_n$  will be **exact** if we use all the  $D$  basis vectors
- Will be **approximate** if we only use  $K < D$  basis vectors:  $\mathbf{x}_n \approx \sum_{k=1}^K z_{nk} \mathbf{u}_k$

# Reconstruction of Data from Projections

- Reconstruction of  $\mathbf{x}_n$  from  $\mathbf{z}_n$  will be **exact** if we use all the  $D$  basis vectors
- Will be **approximate** if we only use  $K < D$  basis vectors:  $\mathbf{x}_n \approx \sum_{k=1}^K z_{nk} \mathbf{u}_k$
- Let's use  $K = 1$  basis vector. Then the **one-dim embedding** of  $\mathbf{x}_n$  is

$$\mathbf{z}_n = \mathbf{u}_1^\top \mathbf{x}_n \quad (\text{note: this will just be a scalar})$$

# Reconstruction of Data from Projections

- Reconstruction of  $\mathbf{x}_n$  from  $\mathbf{z}_n$  will be **exact** if we use all the  $D$  basis vectors
- Will be **approximate** if we only use  $K < D$  basis vectors:  $\mathbf{x}_n \approx \sum_{k=1}^K z_{nk} \mathbf{u}_k$
- Let's use  $K = 1$  basis vector. Then the **one-dim embedding** of  $\mathbf{x}_n$  is

$$\mathbf{z}_n = \mathbf{u}_1^\top \mathbf{x}_n \quad (\text{note: this will just be a scalar})$$

- We can now try “reconstructing”  $\mathbf{x}_n$  from its embedding  $\mathbf{z}_n$  as follows

$$\tilde{\mathbf{x}}_n = \mathbf{u}_1 \mathbf{z}_n = \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n$$

# Reconstruction of Data from Projections

- Reconstruction of  $\mathbf{x}_n$  from  $\mathbf{z}_n$  will be **exact** if we use all the  $D$  basis vectors
- Will be **approximate** if we only use  $K < D$  basis vectors:  $\mathbf{x}_n \approx \sum_{k=1}^K z_{nk} \mathbf{u}_k$
- Let's use  $K = 1$  basis vector. Then the **one-dim embedding** of  $\mathbf{x}_n$  is

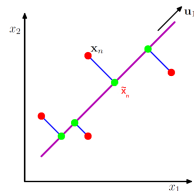
$$\mathbf{z}_n = \mathbf{u}_1^\top \mathbf{x}_n \quad (\text{note: this will just be a scalar})$$

- We can now try “reconstructing”  $\mathbf{x}_n$  from its embedding  $\mathbf{z}_n$  as follows

$$\tilde{\mathbf{x}}_n = \mathbf{u}_1 \mathbf{z}_n = \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n$$

- Total error or “loss” in reconstructing all the data points

$$L(\mathbf{u}_1) = \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2 = \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n\|^2$$



# Direction with Best Reconstruction

- We want to find  $\mathbf{u}_1$  that minimizes the reconstruction error

$$L(\mathbf{u}_1) = \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{u}_1 \mathbf{u}_1^T \mathbf{x}_n\|^2$$

# Direction with Best Reconstruction

- We want to find  $\mathbf{u}_1$  that minimizes the reconstruction error

$$\begin{aligned} L(\mathbf{u}_1) &= \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n\|^2 \\ &= \sum_{n=1}^N \{ \mathbf{x}_n^\top \mathbf{x}_n + (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n)^\top (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n) - 2 \mathbf{x}_n^\top \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n \} \end{aligned}$$



# Direction with Best Reconstruction

- We want to find  $\mathbf{u}_1$  that minimizes the reconstruction error

$$\begin{aligned}L(\mathbf{u}_1) &= \sum_{n=1}^N ||\mathbf{x}_n - \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n||^2 \\&= \sum_{n=1}^N \{ \mathbf{x}_n^\top \mathbf{x}_n + (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n)^\top (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n) - 2 \mathbf{x}_n^\top \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n \} \\&= \sum_{n=1}^N -\mathbf{u}_1^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{u}_1 \quad (\text{using } \mathbf{u}_1^\top \mathbf{u}_1 = 1 \text{ and ignoring constants w.r.t. } \mathbf{u}_1)\end{aligned}$$

# Direction with Best Reconstruction

- We want to find  $\mathbf{u}_1$  that minimizes the reconstruction error

$$\begin{aligned}L(\mathbf{u}_1) &= \sum_{n=1}^N ||\mathbf{x}_n - \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n||^2 \\&= \sum_{n=1}^N \{\mathbf{x}_n^\top \mathbf{x}_n + (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n)^\top (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n) - 2\mathbf{x}_n^\top \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n\} \\&= \sum_{n=1}^N -\mathbf{u}_1^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{u}_1 \quad (\text{using } \mathbf{u}_1^\top \mathbf{u}_1 = 1 \text{ and ignoring constants w.r.t. } \mathbf{u}_1)\end{aligned}$$

- Thus the problem is equivalent to the following **maximization**

# Direction with Best Reconstruction

- We want to find  $\mathbf{u}_1$  that minimizes the reconstruction error

$$\begin{aligned}L(\mathbf{u}_1) &= \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n\|^2 \\&= \sum_{n=1}^N \{\mathbf{x}_n^\top \mathbf{x}_n + (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n)^\top (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n) - 2\mathbf{x}_n^\top \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n\} \\&= \sum_{n=1}^N -\mathbf{u}_1^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{u}_1 \quad (\text{using } \mathbf{u}_1^\top \mathbf{u}_1 = 1 \text{ and ignoring constants w.r.t. } \mathbf{u}_1)\end{aligned}$$

- Thus the problem is equivalent to the following maximization

$$\arg \max_{\mathbf{u}_1: \|\mathbf{u}_1\|^2=1} \mathbf{u}_1^\top \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{u}_1 = \arg \max_{\mathbf{u}_1: \|\mathbf{u}_1\|^2=1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$$

where  $\mathbf{S}$  is the covariance matrix of the data (data assumed centered)

# Direction with Best Reconstruction

- We want to find  $\mathbf{u}_1$  that minimizes the reconstruction error

$$\begin{aligned}L(\mathbf{u}_1) &= \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n\|^2 \\&= \sum_{n=1}^N \{\mathbf{x}_n^\top \mathbf{x}_n + (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n)^\top (\mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n) - 2\mathbf{x}_n^\top \mathbf{u}_1 \mathbf{u}_1^\top \mathbf{x}_n\} \\&= \sum_{n=1}^N -\mathbf{u}_1^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{u}_1 \quad (\text{using } \mathbf{u}_1^\top \mathbf{u}_1 = 1 \text{ and ignoring constants w.r.t. } \mathbf{u}_1)\end{aligned}$$

- Thus the problem is equivalent to the following maximization

$$\arg \max_{\mathbf{u}_1: \|\mathbf{u}_1\|^2=1} \mathbf{u}_1^\top \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{u}_1 = \arg \max_{\mathbf{u}_1: \|\mathbf{u}_1\|^2=1} \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$$

where  $\mathbf{S}$  is the covariance matrix of the data (data assumed centered)

- It's the same objective that we had when we maximized the variance

# How many Principal Components to Use?

- Eigenvalue  $\lambda_k$  measures the variance captured by the corresponding PC  $\mathbf{u}_k$

# How many Principal Components to Use?

- Eigenvalue  $\lambda_k$  measures the variance captured by the corresponding PC  $\mathbf{u}_k$
- The “left-over” variance will therefore be

$$\sum_{k=K+1}^D \lambda_k$$

# How many Principal Components to Use?

- Eigenvalue  $\lambda_k$  measures the variance captured by the corresponding PC  $\mathbf{u}_k$
- The “left-over” variance will therefore be

$$\sum_{k=K+1}^D \lambda_k$$

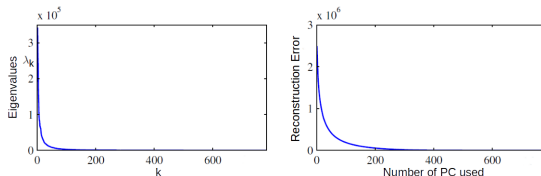
- Can choose  $K$  by looking at what fraction of variance is captured by the first  $K$  PCs

# How many Principal Components to Use?

- Eigenvalue  $\lambda_k$  measures the variance captured by the corresponding PC  $\mathbf{u}_k$
- The “left-over” variance will therefore be

$$\sum_{k=K+1}^D \lambda_k$$

- Can choose  $K$  by looking at what fraction of variance is captured by the first  $K$  PCs
- Another direct way is to look at the spectrum of the eigenvalues plot, or the plot of reconstruction error vs number of PC



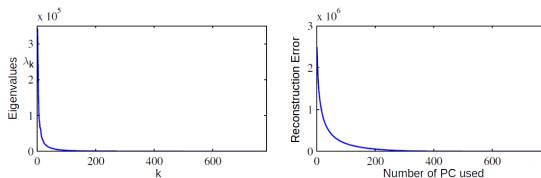


# How many Principal Components to Use?

- Eigenvalue  $\lambda_k$  measures the variance captured by the corresponding PC  $\mathbf{u}_k$
- The “left-over” variance will therefore be

$$\sum_{k=K+1}^D \lambda_k$$

- Can choose  $K$  by looking at what fraction of variance is captured by the first  $K$  PCs
- Another direct way is to look at the spectrum of the eigenvalues plot, or the plot of reconstruction error vs number of PC



- Can also use other criteria such as AIC/BIC (or more advanced probabilistic approaches to PCA using nonparametric Bayesian methods)

# PCA as Matrix Factorization

- Note that PCA represents each  $\mathbf{x}_n$  as  $\mathbf{x}_n = \mathbf{U}\mathbf{z}_n$

# PCA as Matrix Factorization

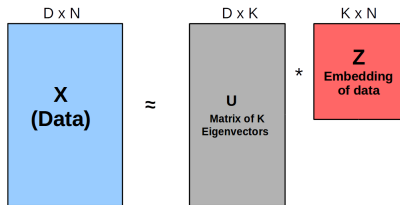
- Note that PCA represents each  $\mathbf{x}_n$  as  $\mathbf{x}_n = \mathbf{U}\mathbf{z}_n$
- When using only  $K < D$  components,  $\mathbf{x}_n \approx \mathbf{U}\mathbf{z}_n$

# PCA as Matrix Factorization

- Note that PCA represents each  $\mathbf{x}_n$  as  $\mathbf{x}_n = \mathbf{U}\mathbf{z}_n$
- When using only  $K < D$  components,  $\mathbf{x}_n \approx \mathbf{U}\mathbf{z}_n$
- For all the  $N$  data points, we can write the same as

$$\mathbf{X} \approx \mathbf{U}\mathbf{Z}$$

where  $\mathbf{X}$  is  $D \times N$ ,  $\mathbf{U}$  is  $D \times K$  and  $\mathbf{Z}$  is  $K \times N$

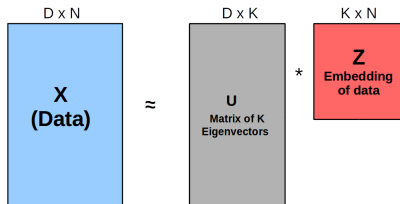


# PCA as Matrix Factorization

- Note that PCA represents each  $\mathbf{x}_n$  as  $\mathbf{x}_n = \mathbf{U}\mathbf{z}_n$
- When using only  $K < D$  components,  $\mathbf{x}_n \approx \mathbf{U}\mathbf{z}_n$
- For all the  $N$  data points, we can write the same as

$$\mathbf{X} \approx \mathbf{U}\mathbf{Z}$$

where  $\mathbf{X}$  is  $D \times N$ ,  $\mathbf{U}$  is  $D \times K$  and  $\mathbf{Z}$  is  $K \times N$



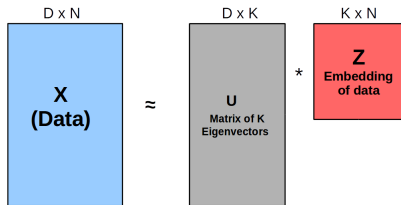
- The above approx. is equivalent to a **low-rank matrix factorization** of  $\mathbf{X}$

# PCA as Matrix Factorization

- Note that PCA represents each  $\mathbf{x}_n$  as  $\mathbf{x}_n = \mathbf{U}\mathbf{z}_n$
- When using only  $K < D$  components,  $\mathbf{x}_n \approx \mathbf{U}\mathbf{z}_n$
- For all the  $N$  data points, we can write the same as

$$\mathbf{X} \approx \mathbf{U}\mathbf{Z}$$

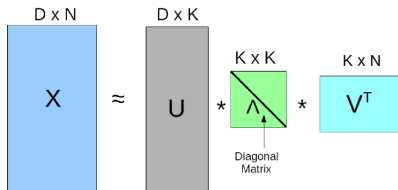
where  $\mathbf{X}$  is  $D \times N$ ,  $\mathbf{U}$  is  $D \times K$  and  $\mathbf{Z}$  is  $K \times N$



- The above approx. is equivalent to a **low-rank matrix factorization** of  $\mathbf{X}$ 
  - Also closely related to Singular Value Decomposition (SVD); see next slide

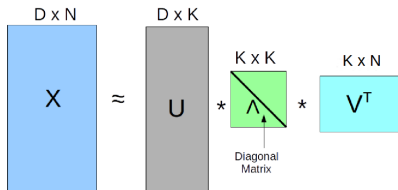
# PCA and SVD

- A rank- $K$  SVD approximates a data matrix  $\mathbf{X}$  as follows:  $\mathbf{X} \approx \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$



# PCA and SVD

- A rank- $K$  SVD approximates a data matrix  $\mathbf{X}$  as follows:  $\mathbf{X} \approx \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$

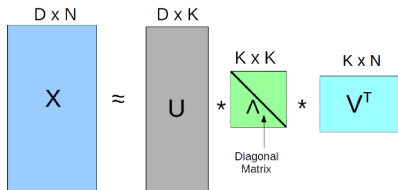


- $\mathbf{U}$  is  $D \times K$  matrix with top  $K$  left singular vectors of  $\mathbf{X}$



# PCA and SVD

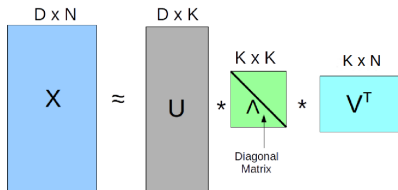
- A rank- $K$  SVD approximates a data matrix  $\mathbf{X}$  as follows:  $\mathbf{X} \approx \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$



- $\mathbf{U}$  is  $D \times K$  matrix with top  $K$  left singular vectors of  $\mathbf{X}$
- $\mathbf{\Lambda}$  is a  $K \times K$  diagonal matrix (with top  $K$  singular values)

# PCA and SVD

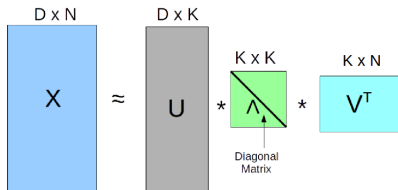
- A rank- $K$  SVD approximates a data matrix  $\mathbf{X}$  as follows:  $\mathbf{X} \approx \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$



- $\mathbf{U}$  is  $D \times K$  matrix with top  $K$  left singular vectors of  $\mathbf{X}$
- $\mathbf{\Lambda}$  is a  $K \times K$  diagonal matrix (with top  $K$  singular values)
- $\mathbf{V}$  is  $N \times K$  matrix with top  $K$  right singular vectors of  $\mathbf{X}$

# PCA and SVD

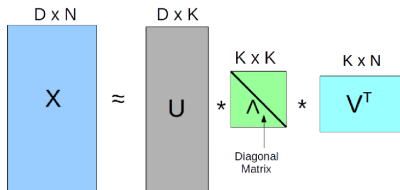
- A rank- $K$  SVD approximates a data matrix  $\mathbf{X}$  as follows:  $\mathbf{X} \approx \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$



- $\mathbf{U}$  is  $D \times K$  matrix with top  $K$  left singular vectors of  $\mathbf{X}$
- $\mathbf{\Lambda}$  is a  $K \times K$  diagonal matrix (with top  $K$  singular values)
- $\mathbf{V}$  is  $N \times K$  matrix with top  $K$  right singular vectors of  $\mathbf{X}$
- Rank- $K$  SVD is based on minimizing the reconstruction error

# PCA and SVD

- A rank- $K$  SVD approximates a data matrix  $\mathbf{X}$  as follows:  $\mathbf{X} \approx \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$

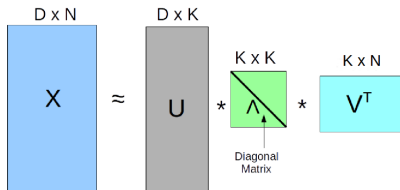


- $\mathbf{U}$  is  $D \times K$  matrix with top  $K$  left singular vectors of  $\mathbf{X}$
- $\mathbf{\Lambda}$  is a  $K \times K$  diagonal matrix (with top  $K$  singular values)
- $\mathbf{V}$  is  $N \times K$  matrix with top  $K$  right singular vectors of  $\mathbf{X}$
- Rank- $K$  SVD is based on minimizing the reconstruction error

$$\|\mathbf{X} - \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T\|$$

# PCA and SVD

- A rank- $K$  SVD approximates a data matrix  $\mathbf{X}$  as follows:  $\mathbf{X} \approx \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$



- $\mathbf{U}$  is  $D \times K$  matrix with top  $K$  left singular vectors of  $\mathbf{X}$
- $\mathbf{\Lambda}$  is a  $K \times K$  diagonal matrix (with top  $K$  singular values)
- $\mathbf{V}$  is  $N \times K$  matrix with top  $K$  right singular vectors of  $\mathbf{X}$
- Rank- $K$  SVD is based on minimizing the reconstruction error

$$\|\mathbf{X} - \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T\|$$

- PCA is equivalent to the **best rank- $K$  SVD** after centering the data

# PCA: Some Comments

- The idea of approximating each data point as a combination of basis vectors

$$\mathbf{x}_n \approx \sum_{k=1}^K z_{nk} \mathbf{u}_k \quad \text{or} \quad \mathbf{X} \approx \mathbf{UZ}$$

is also popularly known as “**Dictionary Learning**” in signal/image processing; the learned basis vectors represent the “Dictionary”

# PCA: Some Comments

- The idea of approximating each data point as a combination of basis vectors

$$\mathbf{x}_n \approx \sum_{k=1}^K z_{nk} \mathbf{u}_k \quad \text{or} \quad \mathbf{X} \approx \mathbf{UZ}$$

is also popularly known as “**Dictionary Learning**” in signal/image processing; the learned basis vectors represent the “Dictionary”

- Some examples:
  - Each face in a collection as a combination of a small no of “eigenfaces”

# PCA: Some Comments

- The idea of approximating each data point as a combination of basis vectors

$$\mathbf{x}_n \approx \sum_{k=1}^K z_{nk} \mathbf{u}_k \quad \text{or} \quad \mathbf{X} \approx \mathbf{UZ}$$

is also popularly known as “**Dictionary Learning**” in signal/image processing; the learned basis vectors represent the “Dictionary”

- Some examples:
  - Each face in a collection as a combination of a small no of “eigenfaces”

$$\begin{array}{ccc} \mathbf{X} \text{ (D} \times \text{N)} & & \mathbf{U} \text{ (D} \times \text{K)} \\ \left( \begin{array}{c} \text{[face 1]} \quad \dots \quad \text{[face N]} \end{array} \right) & \approx & \left( \begin{array}{c} \text{[eigenface 1]} \quad \text{[eigenface 2]} \quad \dots \quad \text{[eigenface K]} \end{array} \right) \end{array} \begin{array}{c} \mathbf{Z} \text{ (K} \times \text{N)} \\ \left( \begin{array}{c} \text{[z}_{11} \text{]} \quad \dots \quad \text{[z}_{1N} \text{]} \\ \vdots \\ \text{[z}_{K1} \text{]} \quad \dots \quad \text{[z}_{KN} \text{]} \end{array} \right) \end{array}$$

- Each document in a collection as a comb. of a small no of “topics”



# PCA: Some Comments

- The idea of approximating each data point as a combination of basis vectors

$$\mathbf{x}_n \approx \sum_{k=1}^K z_{nk} \mathbf{u}_k \quad \text{or} \quad \mathbf{X} \approx \mathbf{UZ}$$

is also popularly known as “**Dictionary Learning**” in signal/image processing; the learned basis vectors represent the “Dictionary”

- Some examples:

- Each face in a collection as a combination of a small no of “eigenfaces”

$$\begin{array}{ccc} \mathbf{X} \text{ (D} \times \text{N)} & & \mathbf{U} \text{ (D} \times \text{K)} \end{array} \quad \begin{array}{c} \mathbf{Z} \text{ (K} \times \text{N)} \\ \left( \begin{array}{c} \text{z}_1 \dots \text{z}_N \end{array} \right) \end{array}$$

$\left( \begin{array}{c} \text{img}_1 \dots \text{img}_N \end{array} \right) \approx \left( \begin{array}{c} \text{eigenface}_1 \dots \text{eigenface}_K \end{array} \right) \left( \begin{array}{c} \text{z}_1 \dots \text{z}_N \end{array} \right)$

- Each document in a collection as a comb. of a small no of “topics”
- Each gene-expression sample as a comb. of a small no of “genetic pathways”

# PCA: Some Comments

- The idea of approximating each data point as a combination of basis vectors

$$\mathbf{x}_n \approx \sum_{k=1}^K z_{nk} \mathbf{u}_k \quad \text{or} \quad \mathbf{X} \approx \mathbf{UZ}$$

is also popularly known as “**Dictionary Learning**” in signal/image processing; the learned basis vectors represent the “Dictionary”

- Some examples:

- Each face in a collection as a combination of a small no of “eigenfaces”

$$\mathbf{X} \text{ (D} \times \text{N)} \quad \mathbf{U} \text{ (D} \times \text{K)} \quad \mathbf{Z} \text{ (K} \times \text{N)}$$
$$\left( \begin{array}{c|c|c} \text{img 1} & \dots & \text{img N} \end{array} \right) \approx \left( \begin{array}{c|c|c|c|c} \text{eigenface 1} & \text{eigenface 2} & \text{eigenface 3} & \text{eigenface 4} & \text{eigenface 5} \end{array} \right) \left( \begin{array}{c|c|c} \mathbf{z}_1 & \dots & \mathbf{z}_N \end{array} \right)$$

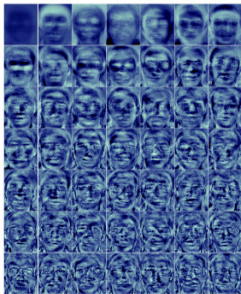
- Each document in a collection as a comb. of a small no of “topics”
- Each gene-expression sample as a comb. of a small no of “genetic pathways”
- The “eigenfaces”, “topics”, “genetic pathways”, etc. are the “basis vectors”, which can be learned from data using PCA/SVD or other similar methods

# PCA: Example

Original Collection of Images



K=49 Eigenvectors  
("eigenfaces") learned  
by PCA on this data



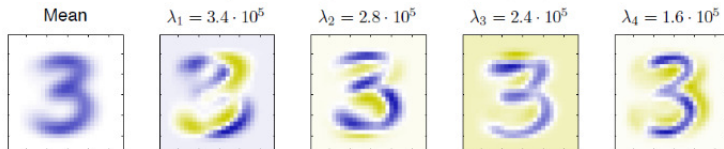
Each image's reconstructed version



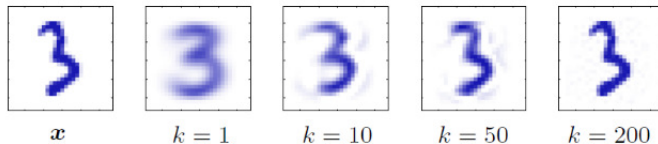
# PCA: Example

$16 \times 16$  pixel images of handwritten 3s (as vectors in  $\mathbb{R}^{256}$ )

Mean  $\mu$  and eigenvectors  $v_1, v_2, v_3, v_4$



Reconstructions:



Each input image now represented by **just k numbers**  
(combination weights of each of the k eigenvectors)

# PCA: Limitations and Extensions

- A linear projection method

# PCA: Limitations and Extensions

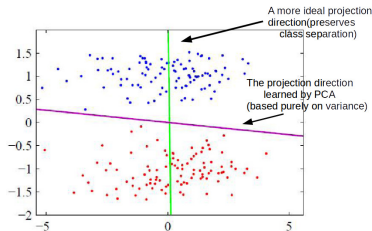
- A linear projection method
  - Won't work well if data can't be approximated by a linear subspace

# PCA: Limitations and Extensions

- A linear projection method
  - Won't work well if data can't be approximated by a linear subspace
  - But PCA can be kernelized easily (Kernel PCA)

# PCA: Limitations and Extensions

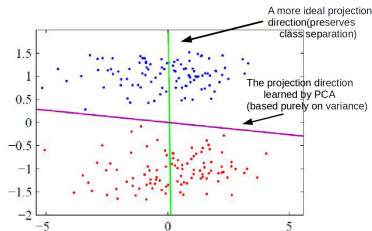
- A linear projection method
  - Won't work well if data can't be approximated by a linear subspace
  - But PCA can be kernelized easily (Kernel PCA)
- Variance based projection directions can sometimes be suboptimal (e.g., if we want to preserve class separation, e.g., when doing classification)





# PCA: Limitations and Extensions

- A linear projection method
  - Won't work well if data can't be approximated by a linear subspace
  - But PCA can be kernelized easily (Kernel PCA)
- Variance based projection directions can sometimes be suboptimal (e.g., if we want to preserve class separation, e.g., when doing classification)



- PCA relies on eigendecomposition of an  $D \times D$  covariance matrix
  - Can be slow if done naively. Takes  $O(D^3)$  time
  - Many faster methods exist (e.g., Power Method)