Linear Dimensionality Reduction: Principal Component Analysis

Piyush Rai

Machine Learning (CS771A)

Sept 2, 2016

イロト イポト イヨト イヨト

-

イロト 不得 とくほと くほど

• Usually considered an unsupervised learning method

э.

イロト 不得 とくほと くほど

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



イロト イヨト イヨト イヨト

- 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.,

< ロ > < 同 > < 回 > < 回 >

- 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

< ロ > < 同 > < 回 > < 回 >

- 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)

< ロ > < 同 > < 回 > < 回 >

- 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

(日) (周) (ヨ) (ヨ)

- 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

(日) (周) (ヨ) (ヨ)

#### **Curse of Dimensionality**

• Exponentially large # of examples required to "fill up" high-dim spaces



э.

イロン 不同と 不同と 不同と

#### **Curse of Dimensionality**

• Exponentially large # of examples required to "fill up" high-dim spaces



 $\bullet~\mbox{Fewer dimensions}$   $\Rightarrow~\mbox{Less chances of overfitting}$   $\Rightarrow~\mbox{Better generalization}$ 

-

イロン 不同と イヨン イヨン

#### **Curse of Dimensionality**

• Exponentially large # of examples required to "fill up" high-dim spaces



 $\bullet~\mbox{Fewer dimensions}$   $\Rightarrow~\mbox{Less chances of overfitting}$   $\Rightarrow~\mbox{Better generalization}$ 

• Dimensionality reduction is a way to beat the curse of dimensionality

イロト イポト イヨト イヨト

• 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)

-

ヘロト 人間 とくほ とくほ とう

- 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 **U** to transform  $\boldsymbol{x}_n \in \mathbb{R}^D$  into  $\boldsymbol{z}_n \in \mathbb{R}^K$  as shown below



イロト イポト イヨト イヨト

- 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 **U** to transform  $\boldsymbol{x}_n \in \mathbb{R}^D$  into  $\boldsymbol{z}_n \in \mathbb{R}^K$  as shown below



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

・ロン ・雪と ・ヨン・

- 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 **U** to transform  $\boldsymbol{x}_n \in \mathbb{R}^D$  into  $\boldsymbol{z}_n \in \mathbb{R}^K$  as shown below



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

Machine Learning (CS771A)

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

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

•  $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$  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

◆□▶ ◆□▶ ◆□▶ ◆□▶ ●□

•  $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$  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



э.

ヘロト 人間 とくほ とくほ とう

•  $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$  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



• How do we learn the "best" projection matrix U?

・ロト ・同ト ・ヨト ・ヨト 一日

•  $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$  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



- How do we learn the "best" projection matrix U?
- What criteria should we optimize for when learning U?

・ロン ・雪と ・ヨン・

•  $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]$  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



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

Machine Learning (CS771A)

・ロン ・雪と ・ヨン・

イロン 不得と 不足と 不足と 一足

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

-

・ロト ・ 同ト ・ ヨト ・ ヨト

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

-

ヘロト 人間 とくほ とくほ とう

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

-

ヘロト 人間 とくほ とくほ とう

- 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

-

- 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)



イロト イポト イヨト イヨト

- 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)

Machine Learning (CS771A)

# PCA as Maximizing Variance

- Consider projecting  $\boldsymbol{x}_n \in \mathbb{R}^D$  on a one-dim subspace defined by  $\boldsymbol{u}_1 \in \mathbb{R}^D$
- Projection/embedding of x<sub>n</sub> along a one-dim subspace u<sub>1</sub> = u<sub>1</sub><sup>⊤</sup>x<sub>n</sub> (location of the green point along the purple line representing u<sub>1</sub>)



イロン 不同と 不同と 不同と

- Consider projecting  $\boldsymbol{x}_n \in \mathbb{R}^D$  on a one-dim subspace defined by  $\boldsymbol{u}_1 \in \mathbb{R}^D$
- Projection/embedding of x<sub>n</sub> along a one-dim subspace u<sub>1</sub> = u<sub>1</sub><sup>⊤</sup>x<sub>n</sub> (location of the green point along the purple line representing u<sub>1</sub>)



• Mean of projections of all the data:  $\frac{1}{N}\sum_{n=1}^{N} u_1^\top x_n = u_1^\top (\frac{1}{N}\sum_{n=1}^{N} x_n) = u_1^\top \mu$ 

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

- Consider projecting  $\boldsymbol{x}_n \in \mathbb{R}^D$  on a one-dim subspace defined by  $\boldsymbol{u}_1 \in \mathbb{R}^D$
- Projection/embedding of x<sub>n</sub> along a one-dim subspace u<sub>1</sub> = u<sub>1</sub><sup>⊤</sup>x<sub>n</sub> (location of the green point along the purple line representing u<sub>1</sub>)



- Mean of projections of all the data:  $\frac{1}{N}\sum_{n=1}^{N} u_1^\top x_n = u_1^\top (\frac{1}{N}\sum_{n=1}^{N} x_n) = u_1^\top \mu$
- Variance of the projected data ("spread" of the green points)

$$\frac{1}{N}\sum_{n=1}^{N}\left(\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}-\boldsymbol{u}_{1}^{\top}\boldsymbol{\mu}\right)^{2}=\frac{1}{N}\sum_{n=1}^{N}\left\{\boldsymbol{u}_{1}^{\top}(\boldsymbol{x}_{n}-\boldsymbol{\mu})\right\}^{2}=\boldsymbol{u}_{1}^{\top}\boldsymbol{\mathsf{S}}\boldsymbol{u}_{1}$$

イロト イポト イヨト イヨト

- Consider projecting  $\boldsymbol{x}_n \in \mathbb{R}^D$  on a one-dim subspace defined by  $\boldsymbol{u}_1 \in \mathbb{R}^D$
- Projection/embedding of x<sub>n</sub> along a one-dim subspace u<sub>1</sub> = u<sub>1</sub><sup>⊤</sup>x<sub>n</sub> (location of the green point along the purple line representing u<sub>1</sub>)



- Mean of projections of all the data:  $\frac{1}{N}\sum_{n=1}^{N} u_1^\top x_n = u_1^\top (\frac{1}{N}\sum_{n=1}^{N} x_n) = u_1^\top \mu$
- Variance of the projected data ("spread" of the green points)

$$\frac{1}{N}\sum_{n=1}^{N}\left(\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}-\boldsymbol{u}_{1}^{\top}\boldsymbol{\mu}\right)^{2}=\frac{1}{N}\sum_{n=1}^{N}\left\{\boldsymbol{u}_{1}^{\top}(\boldsymbol{x}_{n}-\boldsymbol{\mu})\right\}^{2}=\boldsymbol{u}_{1}^{\top}\boldsymbol{\mathsf{S}}\boldsymbol{u}_{1}$$

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

Machine Learning (CS771A)

#### **Direction of Maximum Variance**



• We want  $u_1$  s.t. the variance of the projected data is maximized  $\arg\max_{u_1} u_1^\top \mathbf{S} u_1$ 

3

イロン 不同と イヨン イヨン

#### **Direction of Maximum Variance**



- We want  $u_1$  s.t. the variance of the projected data is maximized  $\arg \max_{u_1} u_1^\top \mathbf{S} u_1$
- To prevent trivial solution (max var. = infinite), assume  $||\boldsymbol{u}_1|| = 1 = \boldsymbol{u}_1^\top \boldsymbol{u}_1$

イロン 不同と イヨン イヨン


- We want  $u_1$  s.t. the variance of the projected data is maximized  $\arg \max_{u_1} u_1^\top \mathbf{S} u_1$
- To prevent trivial solution (max var. = infinite), assume  $||\boldsymbol{u}_1|| = 1 = \boldsymbol{u}_1^\top \boldsymbol{u}_1$
- We will find  $u_1$  by solving the following constrained opt. problem

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

where  $\lambda_1$  is a Lagrange multiplier

Machine Learning (CS771A)

メポト イヨト イヨト

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

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

-

ヘロト 人間 ト 人造 ト 人造 トー

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

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

• Thus  $u_1$  is an eigenvector of **S** (with corresponding eigenvalue  $\lambda_1$ )

イロン 不得 とくほど 不良 とうほう

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

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

- Thus  $u_1$  is an eigenvector of **S** (with corresponding eigenvalue  $\lambda_1$ )
- But which of **S**'s (*D* possible) eigenvectors it is?

イロン 不通 とうほう イロン しゅう

- The objective function:  $\arg \max_{u_1} u_1^\top S u_1 + \lambda_1 (1 u_1^\top u_1)$
- Taking the derivative w.r.t.  $\boldsymbol{u}_1$  and setting to zero gives

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

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

$$oldsymbol{u}_1^{ op} \mathbf{S} oldsymbol{u}_1 = \lambda_1$$

◆□▶ ◆□▶ ◆□▶ ◆□▶ ●□

- The objective function:  $\arg \max_{u_1} u_1^\top S u_1 + \lambda_1 (1 u_1^\top u_1)$
- Taking the derivative w.r.t.  $\boldsymbol{u}_1$  and setting to zero gives

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

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

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

• Var. is maximized when  $u_1$  is the (top) eigenvector with largest eigenvalue

◆□▶ ◆□▶ ◆□▶ ◆□▶ ●□

- The objective function:  $\arg \max_{u_1} u_1^\top S u_1 + \lambda_1 (1 u_1^\top u_1)$
- Taking the derivative w.r.t.  $\boldsymbol{u}_1$  and setting to zero gives

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

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

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

- Var. is maximized when  $u_1$  is the (top) eigenvector with largest eigenvalue
- The top eigenvector  $u_1$  is also known as the first Principal Component (PC)

- The objective function:  $\arg \max_{u_1} u_1^\top S u_1 + \lambda_1 (1 u_1^\top u_1)$
- Taking the derivative w.r.t.  $\boldsymbol{u}_1$  and setting to zero gives

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

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

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

- Var. is maximized when  $u_1$  is the (top) eigenvector with largest eigenvalue
- The top eigenvector  $u_1$  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 **S** (this is PCA)

・ロト ・ 同 ト ・ 三 ト ・ 三 ・ つへの

• Steps in Principal Component Analysis

ъ.

イロト イポト イヨト イヨト

- Steps in Principal Component Analysis
  - Center the data (subtract the mean  $\mu = \frac{1}{N} \sum_{n=1}^{N} x_n$  from each data point)

・ロン ・御 と く ヨ と く ヨ と 二 ヨー

- Steps in Principal Component Analysis
  - Center the data (subtract the mean  $\mu = \frac{1}{N} \sum_{n=1}^{N} x_n$  from each data point)
  - $\bullet\,$  Compute the covariance matrix  ${\boldsymbol{S}}$  using the centered data as

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

・ロト ・ 同ト ・ ヨト ・ ヨト - ヨー

- Steps in Principal Component Analysis
  - Center the data (subtract the mean  $\mu = \frac{1}{N} \sum_{n=1}^{N} x_n$  from each data point)
  - $\bullet\,$  Compute the covariance matrix  ${\boldsymbol{S}}$  using the centered data as

$$\mathbf{S} = rac{1}{N} \mathbf{X} \mathbf{X}^ op$$
 (note:  $\mathbf{X}$  assumed  $D imes N$  here)

 $\bullet\,$  Do an eigendecomposition of the covariance matrix  ${\bf S}$ 

◆□▶ ◆□▶ ◆□▶ ◆□▶ ●□

- Steps in Principal Component Analysis
  - Center the data (subtract the mean  $\mu = \frac{1}{N} \sum_{n=1}^{N} x_n$  from each data point)
  - $\bullet\,$  Compute the covariance matrix  ${\boldsymbol{S}}$  using the centered data as

$$\mathbf{S} = rac{1}{N} \mathbf{X} \mathbf{X}^ op$$
 (note:  $\mathbf{X}$  assumed  $D imes N$  here)

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

・ロン (雪) (田) (田) (田)

#### • Steps in Principal Component Analysis

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

$$\mathbf{S} = rac{1}{N} \mathbf{X} \mathbf{X}^ op$$
 (note:  $\mathbf{X}$  assumed  $D imes N$  here)

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

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

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

・ロト ・ 同 ト ・ 三 ト ・ 三 ・ つへの

#### • Steps in Principal Component Analysis

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

$$\mathbf{S} = rac{1}{N} \mathbf{X} \mathbf{X}^ op$$
 (note:  $\mathbf{X}$  assumed  $D imes N$  here)

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

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

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 **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 **Z** is  $N \times K$ 

・ロト ・ 同 ト ・ 三 ト ・ 三 ・ つへの

# PCA as Minimizing the Reconstruction Error

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

・ロン ・御 と く ヨ と く ヨ と 二 ヨー

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

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

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



イロト イポト イヨト イヨト

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



ヘロン 人間 とくほ とくほ とう

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



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

Machine Learning (CS771A)

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

• Reconstruction of  $x_n$  from  $z_n$  will be exact if we use all the D basis vectors

-

イロン 不同と 不同と 不同と

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

-

イロト 不得 トイヨト イヨト

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

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

・ロン (雪) (田) (田) (田)

- Reconstruction of  $x_n$  from  $z_n$  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  $x_n$  is

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

• We can now try "reconstructing"  $x_n$  from its embedding  $z_n$  as follows

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

イロン 不良 とくほど 不良 とうせい

- Reconstruction of  $x_n$  from  $z_n$  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  $x_n$  is

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

• We can now try "reconstructing"  $x_n$  from its embedding  $z_n$  as follows

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

• Total error or "loss" in reconstructing all the data points



・ロン (雪) (田) (田) (田)

• We want to find  $\boldsymbol{u}_1$  that minimizes the reconstruction error

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

-

イロト イポト イヨト イヨト

• We want to find  $\boldsymbol{u}_1$  that minimizes the reconstruction error

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

-

イロト イポト イヨト イヨト

• We want to find  $u_1$  that minimizes the reconstruction error

$$L(\boldsymbol{u}_{1}) = \sum_{n=1}^{N} ||\boldsymbol{x}_{n} - \boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}||^{2}$$
  
$$= \sum_{n=1}^{N} \{\boldsymbol{x}_{n}^{\top}\boldsymbol{x}_{n} + (\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n})^{\top}(\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}) - 2\boldsymbol{x}_{n}^{\top}\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}\}$$
  
$$= \sum_{n=1}^{N} -\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}\boldsymbol{x}_{n}^{\top}\boldsymbol{u}_{1} \qquad (\text{using }\boldsymbol{u}_{1}^{\top}\boldsymbol{u}_{1} = 1 \text{ and ignoring constants w.r.t. } \boldsymbol{u}_{1})$$

-

イロト 不得下 不良下 不良下

• We want to find  $\boldsymbol{u}_1$  that minimizes the reconstruction error

$$L(\boldsymbol{u}_{1}) = \sum_{n=1}^{N} ||\boldsymbol{x}_{n} - \boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}||^{2}$$
  
$$= \sum_{n=1}^{N} \{\boldsymbol{x}_{n}^{\top}\boldsymbol{x}_{n} + (\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n})^{\top}(\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}) - 2\boldsymbol{x}_{n}^{\top}\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}\}$$
  
$$= \sum_{n=1}^{N} -\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}\boldsymbol{x}_{n}^{\top}\boldsymbol{u}_{1} \qquad (\text{using }\boldsymbol{u}_{1}^{\top}\boldsymbol{u}_{1} = 1 \text{ and ignoring constants w.r.t. }\boldsymbol{u}_{1})$$

• Thus the problem is equivalent to the following maximization

-

イロン 不同と 不同と 不同と

• We want to find  $\boldsymbol{u}_1$  that minimizes the reconstruction error

$$L(\boldsymbol{u}_{1}) = \sum_{n=1}^{N} ||\boldsymbol{x}_{n} - \boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}||^{2}$$
  
$$= \sum_{n=1}^{N} \{\boldsymbol{x}_{n}^{\top}\boldsymbol{x}_{n} + (\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n})^{\top}(\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}) - 2\boldsymbol{x}_{n}^{\top}\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}\}$$
  
$$= \sum_{n=1}^{N} -\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}\boldsymbol{x}_{n}^{\top}\boldsymbol{u}_{1} \qquad (\text{using }\boldsymbol{u}_{1}^{\top}\boldsymbol{u}_{1} = 1 \text{ and ignoring constants w.r.t. }\boldsymbol{u}_{1})$$

• Thus the problem is equivalent to the following maximization

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

where **S** is the covariance matrix of the data (data assumed centered)

-

• We want to find  $\boldsymbol{u}_1$  that minimizes the reconstruction error

$$L(\boldsymbol{u}_{1}) = \sum_{n=1}^{N} ||\boldsymbol{x}_{n} - \boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}||^{2}$$
  
$$= \sum_{n=1}^{N} \{\boldsymbol{x}_{n}^{\top}\boldsymbol{x}_{n} + (\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n})^{\top}(\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}) - 2\boldsymbol{x}_{n}^{\top}\boldsymbol{u}_{1}\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}\}$$
  
$$= \sum_{n=1}^{N} -\boldsymbol{u}_{1}^{\top}\boldsymbol{x}_{n}\boldsymbol{x}_{n}^{\top}\boldsymbol{u}_{1} \qquad (\text{using }\boldsymbol{u}_{1}^{\top}\boldsymbol{u}_{1} = 1 \text{ and ignoring constants w.r.t. }\boldsymbol{u}_{1})$$

• Thus the problem is equivalent to the following maximization

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

where **S** is the covariance matrix of the data (data assumed centered)

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

Machine Learning (CS771A)

ヘロト 人間ト 人口ト 人口ト

• Eigenvalue  $\lambda_k$  measures the variance captured by the corresponding PC  $\boldsymbol{u}_k$ 

-

・ロト ・ 同ト ・ ヨト ・ ヨト

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



イロン 不得 とくほど 不良 とうほう

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



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

・ロト ・ 同ト ・ ヨト ・ ヨト … ヨー

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



- 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  $\boldsymbol{u}_k$
- The "left-over" variance will therefore be



• Another direct way is to look at the spectrum of the eigenvalues plot, or the plot of reconstruction error vs number of PC

k = K + 1



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

Machine Learning (CS771A)

• Note that PCA represents each  $x_n$  as  $x_n = Uz_n$ 

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

-

イロン 不同と 不同と 不同と

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

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

where **X** is  $D \times N$ , **U** is  $D \times K$  and **Z** is  $K \times N$ 



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

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

where **X** is  $D \times N$ , **U** is  $D \times K$  and **Z** is  $K \times N$ 



• The above approx. is equivalent to a low-rank matrix factorization of X

・ 同 ト ・ ヨ ト ・ ヨ ト

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

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

where **X** is  $D \times N$ , **U** is  $D \times K$  and **Z** is  $K \times N$ 



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

Machine Learning (CS771A)

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



3

イロト 不得下 不良下 不良下

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



• **U** is  $D \times K$  matrix with top K left singular vectors of **X** 

э.

イロン 不同と 不同と 不同と

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



- **U** is  $D \times K$  matrix with top K left singular vectors of **X**
- A is a  $K \times K$  diagonal matrix (with top K singular values)

ъ.

イロン 不同と イヨン イヨン

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



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

-

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



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

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



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

 $||\mathbf{X} - \mathbf{U} \wedge \mathbf{V}^{ op}||$ 

ъ.

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



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

 $|| \bm{X} - \bm{U} \bm{\Lambda} \bm{V}^\top ||$ 

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

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

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

$$oldsymbol{x}_n pprox \sum_{k=1}^K z_{nk}oldsymbol{u}_k$$
 or  $oldsymbol{X} pprox oldsymbol{UZ}$ 

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

-

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

$$oldsymbol{x}_n pprox \sum_{k=1}^K z_{nk}oldsymbol{u}_k$$
 or  $oldsymbol{X} pprox oldsymbol{UZ}$ 

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"

-

イロト 不得 トイヨト イヨト

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

$$oldsymbol{x}_n pprox \sum_{k=1}^K z_{nk}oldsymbol{u}_k$$
 or  $oldsymbol{X} pprox oldsymbol{\mathsf{UZ}}$ 

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"



• Each document in a collection as a comb. of a small no of "topics"

・ロト ・ 一 ・ ・ ヨト ・ ヨト

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

$$oldsymbol{x}_n pprox \sum_{k=1}^K z_{nk}oldsymbol{u}_k$$
 or  $oldsymbol{X} pprox oldsymbol{\mathsf{UZ}}$ 

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"



- 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"

・ロト ・ 一 ・ ・ ヨト ・ ヨト

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

$$oldsymbol{x}_n pprox \sum_{k=1}^K z_{nk}oldsymbol{u}_k$$
 or  $oldsymbol{X} pprox oldsymbol{\mathsf{UZ}}$ 

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"



- 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

Machine Learning (CS771A)

### **PCA: Example**



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



Each image's reconstructed version



イロン 不同と 不同と 不同と

э

### **PCA: Example**

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



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

Machine Learning (CS771A)

Linear Dimensionality Reduction: Principal Component Analysis

-

8 ▶ (< 38 →

• A linear projection method

3

イロト 不得す 不良す 不良す

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

-

イロン 不同と イヨン イヨン

- 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)

-

イロン 不同と イヨン イヨン

- 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)



< ロ > < 同 > < 回 > < 回 >

- 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 naïvely. Takes  $O(D^3)$  time
  - Many faster methods exists (e.g., Power Method)

3.1 (3.1)