

CS540 Introduction to Artificial Intelligence

Lecture 12

Young Wu

Based on lecture slides by Jerry Zhu and Yingyu Liang

June 27, 2019

High Dimensional Data

Motivation

- High dimensional data are training set with a lot of features.
- ① Document classification.
- ② MEG brain imaging.
- ③ Handwritten digits (or images in general).

Low Dimension Representation

Motivation

- Unsupervised learning techniques are used to find low dimensional representation.
- ① Visualization. 2 or 3 features
- ② Efficient storage.
- ③ Better generalization.
- ④ Noise removal.

Orthogonal Directions

Definition

- In Euclidean space (L_2 norm), a unit vector u_k has length 1.

$$\|u_k\|_2 = \underbrace{u_k^T u_k}_{\quad} = 1$$

- Two vectors $u_k, u_{k'}$ are orthogonal (or uncorrelated) if the dot product is 0.

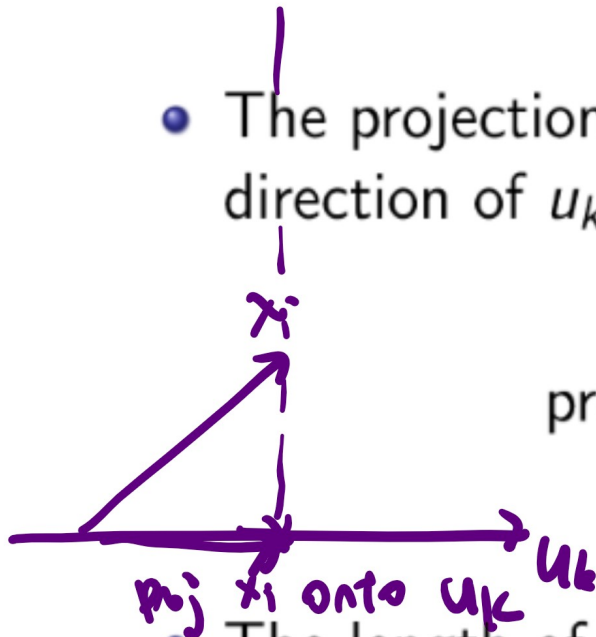
$$u_k \cdot u_{k'} = \underbrace{u_k^T}_{\quad} u_{k'} = 0$$



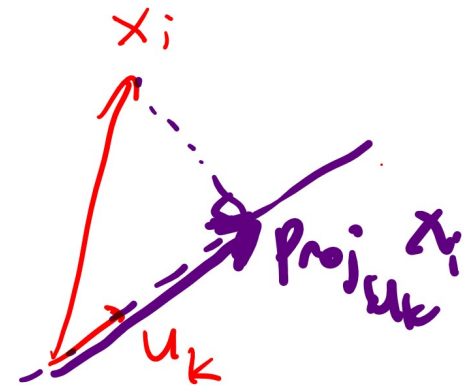
Projection

Definition

- The projection of x_i onto a unit vector u_k is the vector in the direction of u_k that is the closest to x_i .



$$\text{proj}_{u_k} x_i = \left(\frac{u_k^T x_i}{u_k^T u_k} \right) u_k = \underbrace{u_k^T x_i}_{\text{length}} \underbrace{u_k}_{\text{unit direction}}$$



- The length of the projection of x_i onto a unit vector u_k is $u_k^T x_i$.

$$\| \text{proj}_{u_k} x_i \|_2 = u_k^T x_i$$

Project Example, Part I

Quiz (Graded)

• What is the projection of $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ onto $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$?

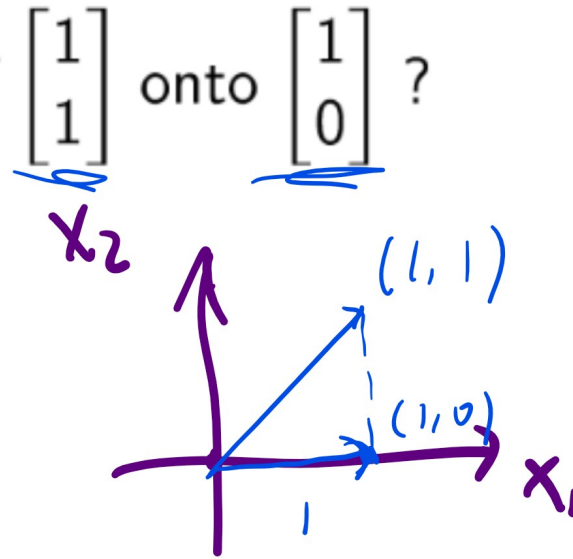
• A: 1

• B: $\frac{1}{\sqrt{2}}$

• C: $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

• D: $\begin{bmatrix} 1 \\ \frac{1}{\sqrt{2}} \\ 0 \end{bmatrix}$

• E: $\begin{bmatrix} \sqrt{2} \\ 0 \end{bmatrix}$



$$\begin{aligned}
 & U_k^T x_j \cdot U_k \\
 & \begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = (1 \cdot 1 + 0 \cdot 1) \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} \\
 & = 1 \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \\
 & \text{length direction}
 \end{aligned}$$

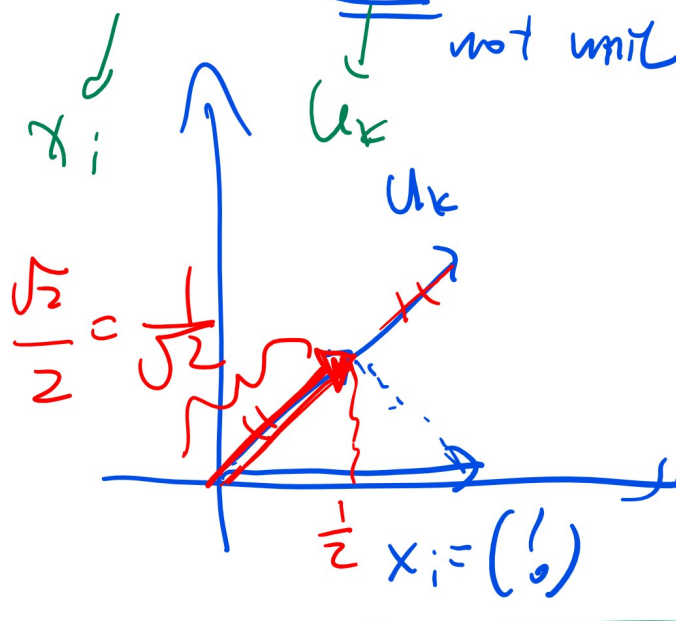
Project Example, Part II

Quiz (Graded)

• What is the projection of $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ onto $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$?

- ~~A: 1~~
- ~~B: $\frac{1}{\sqrt{2}}$~~
- ~~C: $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$~~
- D: $\begin{bmatrix} 1 \\ \sqrt{2} \\ 1 \\ \sqrt{2} \end{bmatrix}$
- **E: $\begin{bmatrix} 1 \\ \frac{1}{2} \\ 1 \\ \frac{1}{2} \end{bmatrix}$**

$$\left(\frac{1}{\sqrt{2}} \right)$$



$$\frac{\begin{pmatrix} 1 \\ 1 \end{pmatrix}}{\| \begin{pmatrix} 1 \\ 1 \end{pmatrix} \|_2} = \frac{\begin{pmatrix} 1 \\ 1 \end{pmatrix}}{\sqrt{(1,1) \begin{pmatrix} 1 \\ 1 \end{pmatrix}}}$$

$\sqrt{u_k^T u_k}$

① convert $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ to unit vector

$\hat{u}_k^T x_i \hat{u}_k$

② $\frac{u_k^T x_i}{u_k^T u_k} u_k$

$$\frac{(1,1) \begin{pmatrix} 1 \\ 0 \end{pmatrix}}{(1,1) \begin{pmatrix} 1 \\ 1 \end{pmatrix}} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} \\ \frac{1}{2} \end{pmatrix}$$

$$\hat{u}_k = \begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}$$

unit vector of u_k

Variance Definition

$$\hat{u}_k^T X_i \hat{u}_k$$

$$= \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right) \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix} = \begin{pmatrix} \frac{1}{2} \\ \frac{1}{2} \end{pmatrix}$$

only works for unit \hat{u}_k

- The sample variance of a data set $\{x_1, x_2, \dots, x_n\}$ is the sum of the squared distance from the mean.

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}$$

$n \times 1$ vector

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{im} \end{pmatrix}$$

$$\begin{pmatrix} x_{i1}^2 & x_{i1}x_{i2} & x_{i1}x_{i3} \\ x_{i2}x_{i1} & x_{i2}^2 & x_{i2}x_{i3} \\ \dots & \dots & \dots \\ x_{im}x_{i1} & x_{im}x_{i2} & x_{im}x_{i3} \end{pmatrix}$$

$n \times m$ matrix

$$\hat{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

$$\begin{pmatrix} x_{i1} \\ \vdots \\ x_{im} \end{pmatrix} \begin{pmatrix} x_{i1} & \dots & x_{im} \end{pmatrix}$$

~~***~~ NOT $(x_i - \hat{\mu})^T (x_i - \hat{\mu})$
Covariance matrix.

Normalization

Definition

- Normalize the data by subtracting the mean, then the variance expression can be simplified.

$$x_i = x_i - \mu$$

$$\hat{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n x_i x_i^T = \frac{1}{n-1} X^T X$$

Maximum Variance Directions

Definition

After a lot of math

compute variance along

- The goal is to find the direction that maximizes the projected variance.

$$u_k^T \lambda u_k$$

$$\lambda u_k^T u_k$$

$$\max_{u_k} u_k^T \hat{\Sigma} u_k \text{ such that } u_k^T u_k = 1$$

$$\Rightarrow \max_{u_k} u_k^T \hat{\Sigma} u_k - \lambda u_k^T u_k$$

$$\Rightarrow \hat{\Sigma} u_k = \lambda u_k$$

direction u_k .

a lot of math again,

max λ

PC
are

eigen vectors of covariance matrix.

find eigen values
eigen vectors

Eigenvalue

Definition

- The λ represents the projected variance.

$$u_k^T \hat{\Sigma} u_k = u_k^T \lambda u_k = \lambda$$

$\lambda_k \uparrow$ variance in $u_k \uparrow$

- The larger the variance, the larger the variability in direction u_k . There are m eigenvalues for a symmetric positive semidefinite matrix (for example, $X^T X$ is always symmetric PSD). Order the eigenvectors u_k by the size of their corresponding eigenvalues λ_k .

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k \geq \lambda_m$$

$k \ll m$

k principal components \rightarrow

\downarrow u_1 \downarrow u_2 \downarrow u_k

Eigenvalue Algorithm

Definition

- Solving eigenvalue using the definition (characteristic polynomial) is computationally inefficient.

$$\left(\hat{\Sigma} - \lambda_k I\right) u_k = 0 \Rightarrow \det \left(\hat{\Sigma} - \lambda_k I\right) = 0$$

- There are many fast eigenvalue algorithms that compute the spectral (eigen) decomposition for real symmetric matrices. Columns of Q are unit eigenvectors and diagonal elements of D are eigenvalues.

$$\begin{aligned}\hat{\Sigma} &= PDP^{-1}, D \text{ is diagonal} \\ &= QDQ^T, \text{ if } Q \text{ is orthogonal, i.e. } Q^T Q = I\end{aligned}$$

Spectral Decomposition Example

Quiz (Participation)

- Given the following spectral decomposition of $\hat{\Sigma}$, what is the first principal component?

$$\hat{\Sigma} = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix}^{-1}$$

Handwritten annotations: λ_1 , λ_2 , λ_3 above the diagonal of D ; v_1 above the first column of P^{-1} ; "largest λ " with an arrow pointing to the 3 in D ; P and P^{-1} labels under the matrices.

- A: $\begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$, B: $\begin{bmatrix} 1 \\ \sqrt{2} \\ 0 \\ 1 \\ \sqrt{2} \end{bmatrix}$, C: $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$, D: $\begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$, E: $\begin{bmatrix} 1 \\ -\sqrt{2} \\ 0 \\ 1 \\ \sqrt{2} \end{bmatrix}$
- Handwritten annotations: "unit vector" under B; a red box around B; an arrow from the box around B in the decomposition above pointing to B.

Principal Component Analysis

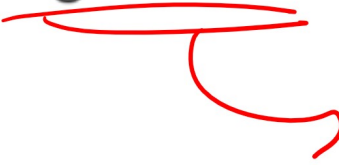
Algorithm

- Input: instances: $\{x_i\}_{i=1}^n$, the number of dimensions after reduction $K < m$.
- Output: K principal components.
- Find the largest K eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_K$.
- Return the corresponding unit orthogonal eigenvectors $u_1, u_2 \dots u_K$.

Reduced Feature Space

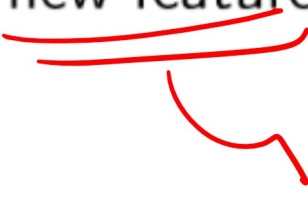
Discussion

- The original feature space is m dimensional.


$$(x_{i1}, x_{i2}, \dots, x_{im})^T$$

m

- The new feature space is K dimensional.


$$(u_1^T x_i, u_2^T x_i, \dots, u_K^T x_i)^T$$

K

- Other supervised learning algorithms can be applied on the new features.

Reduced Space Example

Quiz (Graded)

Q6

- 2017 Fall Final Q10

- If $u_1 = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \end{bmatrix}^T$ and $u_2 = \begin{bmatrix} -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \end{bmatrix}^T$. If one original data is $x = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}^T$. What is the new representation?

- A: $\begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ 1 \end{bmatrix}$, B: $\begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix}$, C: $\begin{bmatrix} \frac{1}{\sqrt{2}} \\ 2 \\ \frac{1}{\sqrt{2}} \end{bmatrix}$, D: $\begin{bmatrix} \frac{3}{\sqrt{2}} \\ 1 \\ \frac{1}{\sqrt{2}} \end{bmatrix}$, E: $\begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{3}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$

$$(u_1^T x, u_2^T x)$$

Number of Dimensions

Discussion

- There are a few ways to choose the number of principal components K .
- K can be selected given prior knowledge or requirement.
- K can be the number of non-zero eigenvalues.
- K can be the number of eigenvalues that are large (larger than some threshold).

remove $\lambda_k < 0.1$

$\lambda_k = \text{variance in } U_k \text{ direction}$
 < 0.1

Reconstruction Error

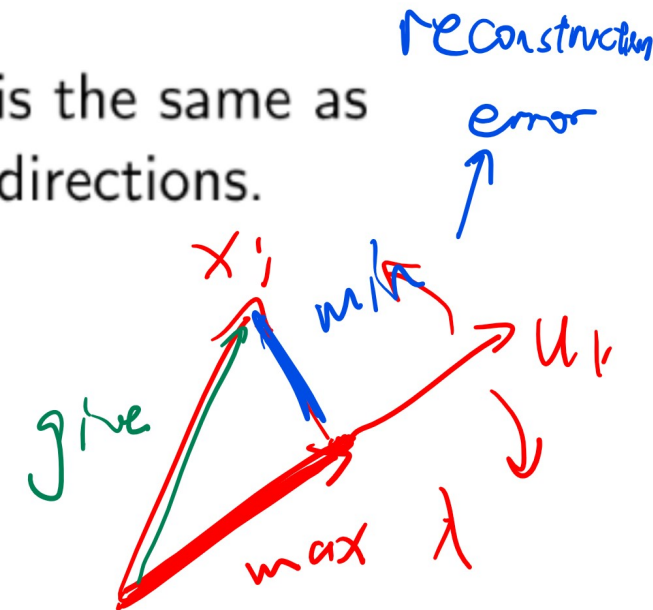
Discussion

- Reconstruction error is the squared error (distance) between the original data and its projection onto u_k .

$$\left\| x_i - \left(u_k^T x_i \right) u_k \right\|^2$$

- Finding the variance maximizing directions is the same as finding the reconstruction error minimizing directions.

$$\frac{1}{n} \sum_{i=1}^n \left\| x_i - \left(u_k^T x_i \right) u_k \right\|^2$$



Reconstruction Error Diagram

Discussion

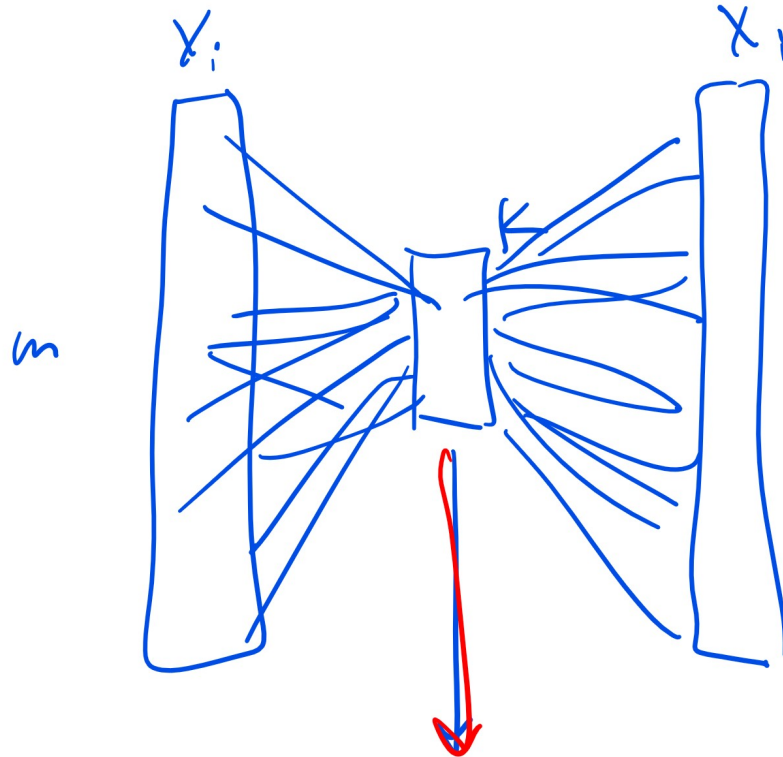
Autoencoder

Discussion

- A multi-layer neural network with the same input and output $y_i = x_i$ is called an autoencoder.
- The hidden layers have fewer units than the dimension of the input m .
- The hidden units form an encoding of the input with reduced dimensionality.

Autoencoder Diagram

Discussion



if $a_i = w^T x + b$

NO $\sigma(w^T x + b)$

Approximate PCA } non linear PCA
 non linear logistic

Eigenface

Discussion

- Eigenfaces are eigenvectors of face images (pixel intensities or HOG features).
- Every face can be written as a linear combination of eigenfaces. The coefficients determine specific faces.

$$x_i = \sum_{k=1}^m \left(\underbrace{u_k^T x_i}_{\text{new face}} \right) u_k \approx \sum_{k=1}^K \left(u_k^T x_i \right) u_k$$

Handwritten annotations:
 - A red arrow points from the '#' symbol above the first sum to the term $(u_k^T x_i)$.
 - A red arrow points from the text "new face" to the first sum.
 - A red arrow points from the text "unique identify a face" to the term $(u_1^T x_i) u_1$.
 - A red arrow points from the text "features" to the term $(u_2^T x_i) u_2$.
 - The terms $(u_1^T x_i) u_1$ and $(u_2^T x_i) u_2$ are circled in red.

- Eigenfaces and SVM can be combined to detect or recognize faces.

T-Distributed Stochastic Neighbor Embedding

Discussion

- t-distributed stochastic neighbor embedding is another non-linear dimensionality reduction method used mainly for visualization.
- Points in high dimensional spaces are embedded in 2 or 3-dimensional spaces to preserve the distance (neighbor) relationship between points.

Embedding Diagram

Discussion

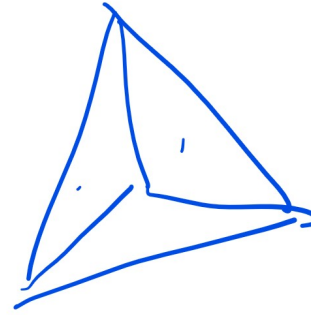
2D



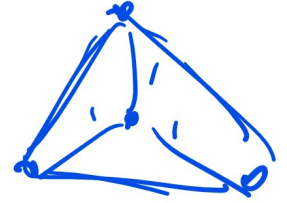
1D



3D



2D



embed

⇒ 2D → 1D keep distance relation