

# CS540 Introduction to Artificial Intelligence

## Lecture 16

Young Wu

Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

July 18, 2022

# Random Choice 1

## Quiz

P1, P2

P3 due tonight + 1 week

P4 can start, not submit.

} P6  
└

one week  
after final

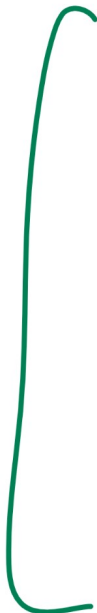
Q7

- Pick a random choice.
- A:
- B:
- C:
- D:
- E:

# Random Choice 2

## Quiz

Q8



- Pick the choice you think is the least popular.
- A :
- B :
- C :
- D :
- E :

Q9

# Random Choice 3

## Quiz

- Pick the choice based on the last digit of your ID.
- *A* : 0 – 1
- *B* : 2 – 3
- *C* : 4 – 5
- *D* : 6 – 7
- *E* : 8 – 9

Q10



# Unsupervised Learning

## Motivation

- Supervised learning:  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  .
  - Unsupervised learning:  $x_1, x_2, \dots, x_n$  .
  - There are a few common tasks without labels.
- 1 Clustering: separate instances into groups.
  - 2 Novelty (outlier) detection: find instances that are different.
  - 3 Dimensionality reduction: represent each instance with a lower dimensional feature vector while maintaining key characteristics.

# 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

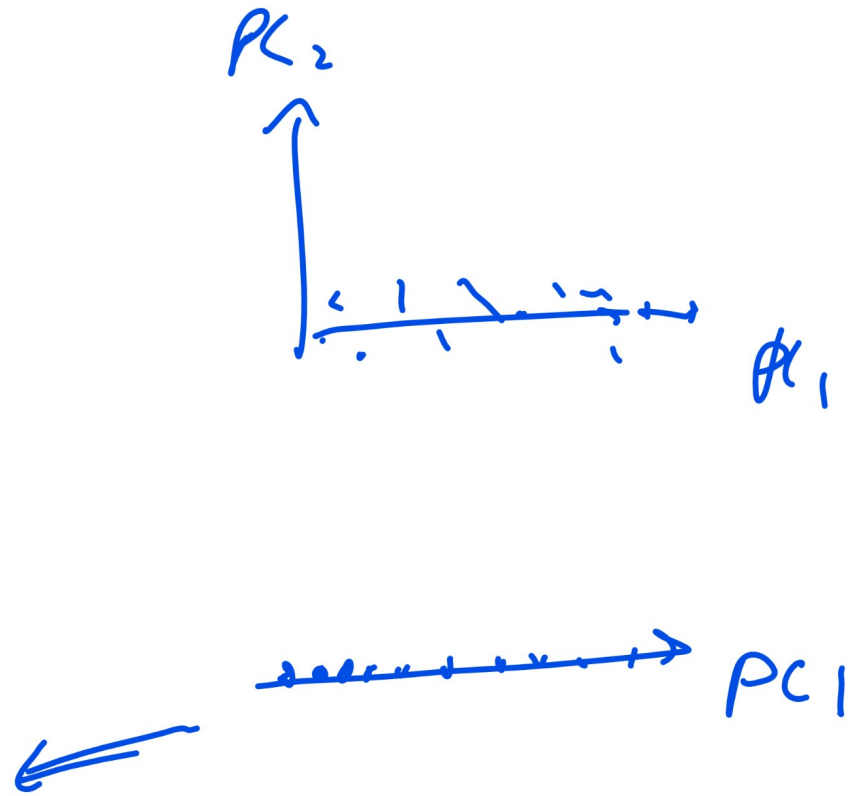
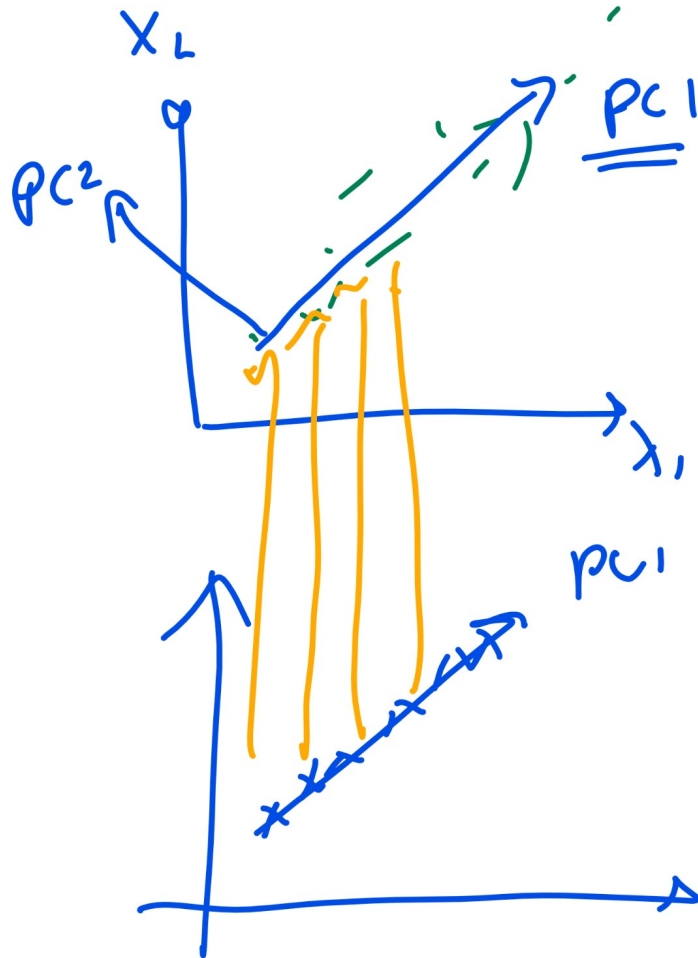
## Motivation

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



# Dimension Reduction Demo

## Motivation



# Projection

## Definition

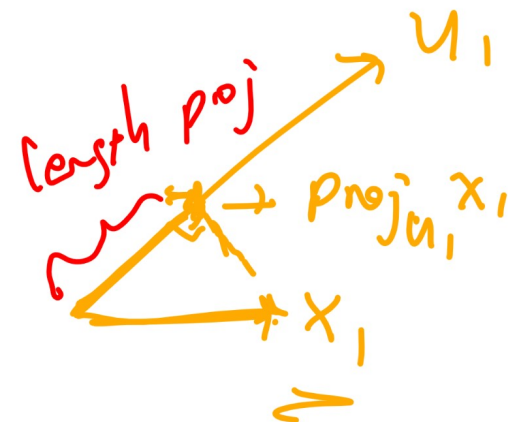
- The projection of  $x_i$  onto a unit vector  $u_k$  is the vector in the direction of  $u_k$  that is the closest to  $x_i$ .

$$\text{proj}_{u_k} x_i = \left( \frac{u_k^T x_i}{u_k^T u_k} \right) u_k = \frac{u_k^T x_i}{\|u_k\|^2} u_k$$

if  $u_k$  is not unit vector  $(1 = u_{k1}^2 + u_{k2}^2 + \dots + u_{kn}^2)$

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

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



# Variance

## Definition

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

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

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

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

$(x_i - \hat{\mu})^T (x_i - \hat{\mu})$   
inner product  
dot product

outer-product.





# Projection Example 1

Quiz

$$\frac{1}{2} \left[ (\sqrt{2} - \bar{\mu})^2 + \left( \frac{\sqrt{2}}{2} - \bar{\mu} \right)^2 \right]$$

where  $\bar{\mu} = \frac{1}{2}(\sqrt{2} + \frac{\sqrt{2}}{2})$

↓  
proj  
variance

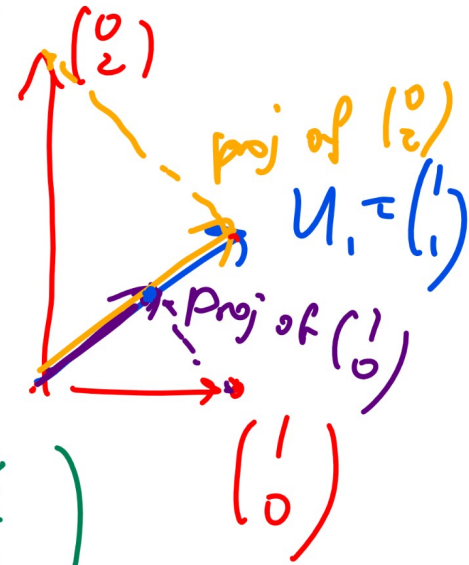
$$\text{proj}_{\begin{pmatrix} 1 \\ 1 \end{pmatrix}} \begin{pmatrix} 0 \\ 2 \end{pmatrix} = \frac{\begin{pmatrix} 1 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 2 \end{pmatrix}}{\begin{pmatrix} 1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix}} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \frac{2}{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

not unit

- What is the projection of  $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$  and  $\begin{bmatrix} 0 \\ 2 \end{bmatrix}$  onto  $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$  and what is the projected variance?

Unit vector  $\begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}$

$\sqrt{2}, \frac{\sqrt{2}}{2}, \dots$



$$\text{proj}_{\begin{pmatrix} 1 \\ 1 \end{pmatrix}} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix} = \left( \frac{1}{\sqrt{2}} \cdot 1 + \frac{1}{\sqrt{2}} \cdot 0 \right) \begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix} = \begin{pmatrix} \frac{1}{2} \\ \frac{1}{2} \end{pmatrix}$$

# Projection Example 3

## Quiz

• What is the projection of

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

onto

$$\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Why? Q11

$$\sqrt{1^2 + 1^2 + 1^2} = \sqrt{3}$$

$$\frac{u^T x}{u^T u} u = \text{proj}_u x$$

$$\begin{pmatrix} \frac{1}{\sqrt{3}} \\ \frac{1}{\sqrt{3}} \\ \frac{1}{\sqrt{3}} \end{pmatrix}$$

$$\begin{aligned} &\sqrt{\left(\frac{1}{\sqrt{3}}\right)^2 + \left(\frac{1}{\sqrt{3}}\right)^2 + \left(\frac{1}{\sqrt{3}}\right)^2} \\ &= \sqrt{\frac{1}{3} + \frac{1}{3} + \frac{1}{3}} \\ &= \sqrt{1} = 1 \end{aligned}$$

$$\frac{6}{3} \cdot \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = 2 \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

•  $A: [2 \ 2 \ 2]^T$

•  $B: [3 \ 3 \ 3]^T$

•  $C: [4 \ 4 \ 4]^T$

•  $D: [6 \ 6 \ 6]^T$

• E: I don't understand.



# Projection Example 4

Quiz

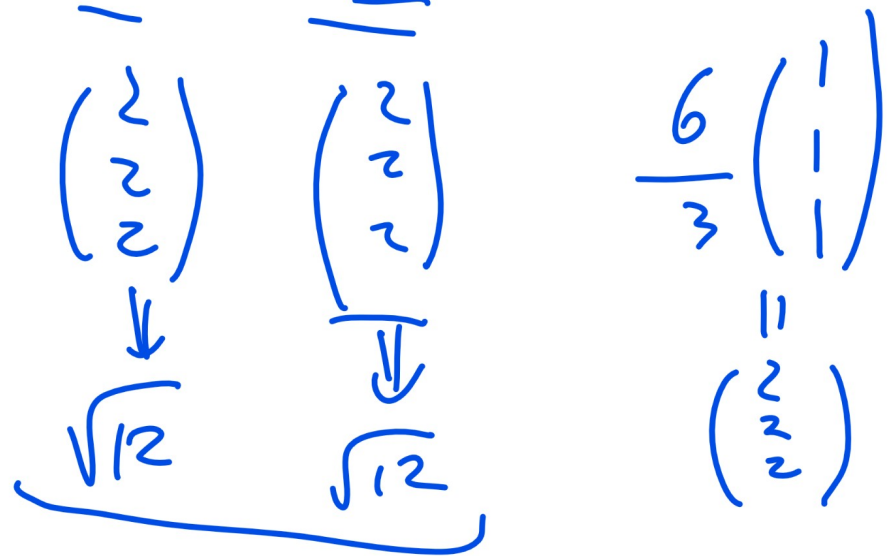
$$\left\{ \frac{1}{2} \left[ (\sqrt{12} - \sqrt{12})^2 + (\sqrt{12} - \sqrt{12})^2 \right] \right\} = 0$$

~~5-1~~

Q12

• What is the projection <sup>ed</sup> variance of  $\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$  and  $\begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}$  onto  $\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$  ?

- A : 0
- B : 12
- C : 24
- D : 48
- E : I don't understand.



# Maximum Variance Directions

## Definition

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

$$\begin{aligned} & \max_{u_k} u_k^T \hat{\Sigma} u_k \text{ such that } u_k^T u_k = 1 \\ & \Rightarrow \max_{u_k} u_k^T \hat{\Sigma} u_k - \lambda u_k^T u_k \\ & \Rightarrow \hat{\Sigma} u_k = \lambda u_k \end{aligned}$$

*derivative*


*eigenvalue.*

*eigenvector*

# Eigenvalue

## Definition

- The  $\lambda$  represents the projected variance.

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


- 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  $\lambda_k$ .

$$\underline{\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m}$$

# 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.

$$\hat{\Sigma} = PDP^{-1}, D \text{ is diagonal}$$

$$= QDQ^T, \text{ if } Q \text{ is orthogonal, i.e. } Q^T Q = I$$

# Spectral Decomposition Example 1 <sup>200</sup> → $K=3$

Quiz

$$PC_1 = \begin{pmatrix} \frac{1}{\sqrt{2}} \\ 0 \\ \frac{1}{\sqrt{2}} \end{pmatrix} \quad PC_2 = \begin{pmatrix} \frac{1}{\sqrt{2}} \\ 0 \\ -\frac{1}{\sqrt{2}} \end{pmatrix} \quad PC_3 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

- Given the following spectral decomposition of  $\hat{\Sigma}$ , what are the first two principal components?

$\hat{\Sigma} = \begin{bmatrix} 1 & 0 & 1 \\ \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \\ 0 & 1 & 0 \\ 1 & 0 & 1 \\ -\frac{1}{\sqrt{2}} & 0 & -\frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & 0 & -\frac{1}{\sqrt{2}} \\ 0 & 1 & 0 \\ \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \end{bmatrix}$

$\underbrace{\begin{bmatrix} 1 & 0 & 1 \\ \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \\ 0 & 1 & 0 \\ 1 & 0 & 1 \\ -\frac{1}{\sqrt{2}} & 0 & -\frac{1}{\sqrt{2}} \end{bmatrix}}_P$   $\underbrace{\begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}}_{\text{diag}}$   $\underbrace{\begin{bmatrix} \frac{1}{\sqrt{2}} & 0 & -\frac{1}{\sqrt{2}} \\ 0 & 1 & 0 \\ \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \end{bmatrix}}_{P^{-1}}$

$\underbrace{\begin{bmatrix} 1 & 0 & 1 \\ \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \\ 0 & 1 & 0 \\ 1 & 0 & 1 \\ -\frac{1}{\sqrt{2}} & 0 & -\frac{1}{\sqrt{2}} \end{bmatrix}}_{\text{eigen vectors}}$   $\underbrace{\begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}}_{\text{eigen values}}$

$PC_2 = \begin{pmatrix} \frac{1}{\sqrt{2}} \\ 0 \\ -\frac{1}{\sqrt{2}} \end{pmatrix}$   $PC_1 = \begin{pmatrix} \frac{1}{\sqrt{2}} \\ 0 \\ \frac{1}{\sqrt{2}} \end{pmatrix}$   $PC_3 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$



# Spectral Decomposition Example 2

## Quiz

- Given the following  $\hat{\Sigma}$ , what are the first two principal components? Q1

unit vector

$e_i v = \text{variance}$

$\hat{\Sigma} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 3 \end{bmatrix}$

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

- A:  $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ , B:  $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ , C:  $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ , D:  $\begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$ , E:  $\begin{bmatrix} 0 \\ 0 \\ 3 \end{bmatrix}$

PC<sub>1</sub>      PC<sub>2</sub>

unit vector



# Reduced Feature Space

## Discussion

- The original feature space is  $m$  dimensional.

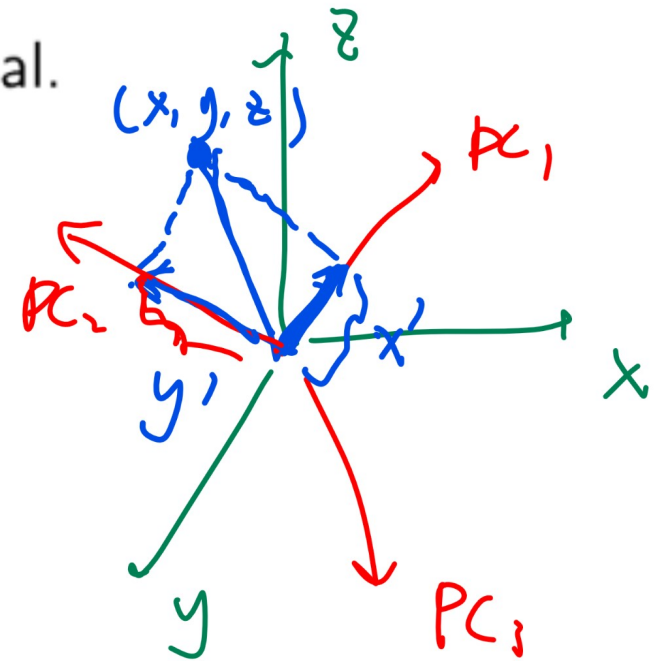
$(x, y, z)$   $\swarrow$   $(x_{i1}, x_{i2}, \dots, x_{im})^T$

- The new feature space is  $K$  dimensional.

$(x', y')$   $\swarrow$   $(u_1^T x_i, u_2^T x_i, \dots, u_K^T x_i)^T$

length of projection on  $T$

$\underbrace{\hspace{1cm}}_{PC_1}$   $\underbrace{\hspace{1cm}}_{PC_2}$   $\underbrace{\hspace{1cm}}_{PC_K}$



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



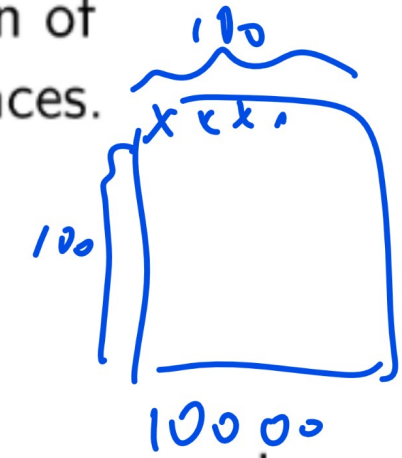
# Eigenface

## 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.

$$x_i = \sum_{k=1}^m (u_k^T x_i) u_k \approx \sum_{k=1}^K (u_k^T x_i) u_k$$

Handwritten annotations:  $K=m$  above the first sum,  $K < m$  above the second sum, a blue double arrow under  $u_k^T x_i$ , and a red arrow pointing to  $x_i$ .



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

$K = 112$

$$(x, y, z) \Rightarrow (\underbrace{x'}_{PC1}, \underbrace{y'}_{PC2}, \dots)$$

# Reduced Space Example 1

## Quiz

• If  $u_1 = \begin{bmatrix} 1 \\ \sqrt{2} \\ 0 \\ 1 \\ \sqrt{2} \end{bmatrix}$  and  $u_2 = \begin{bmatrix} 1 \\ \sqrt{2} \\ 0 \\ 1 \\ -\sqrt{2} \end{bmatrix}$ . If one original item is  $x = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$  (circled in red), what is its new representation and the reconstructed vector using only the two principal components? (circled in blue)

*Handwritten notes:*  
 $\Rightarrow u_3 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$   
 rc construction  
 approx  $x$   
 $(2\sqrt{2}, -\sqrt{2})$  (boxed in red)  
 $2\sqrt{2} \cdot u_1 + -\sqrt{2} \cdot u_2 = \begin{pmatrix} 2 \\ 2 \\ 0 \end{pmatrix} + \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix}$  (circled in blue)

$(x'_1, x'_2)$

length of proj

$$1 \cdot \frac{1}{\sqrt{2}} + 2 \cdot 0 + 3 \cdot \frac{1}{\sqrt{2}} = 2\sqrt{2}$$

$$u_1^T x, u_2^T x = 1 \cdot \frac{1}{\sqrt{2}} + 2 \cdot 0 + 3 \cdot (-\frac{1}{\sqrt{2}}) = -\sqrt{2}$$



# Reduced Space Example 2

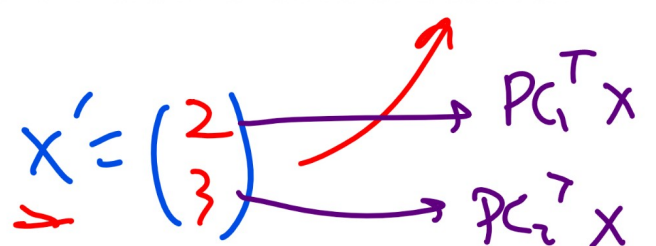
## Quiz

Q2

•  $\hat{\Sigma} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 3 \end{bmatrix}$ . If one original data is  $x = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$ . What is the reconstructed vector using only the first two principal components?

- A:  $\begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}$ , B:  $\begin{bmatrix} 2 \\ 3 \\ 0 \end{bmatrix}$ , C:  $\begin{bmatrix} 0 \\ 2 \\ 3 \end{bmatrix}$ , D:  $\begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix}$ , E: I don't understand.

$PC_1 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$       $PC_2 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$



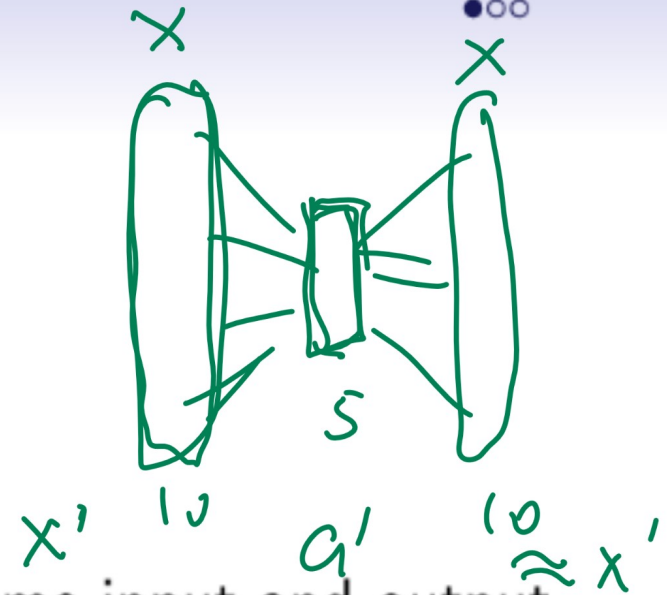
$2 \cdot PC_1 + 3 \cdot PC_2 = \begin{pmatrix} 0 \\ 2 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 3 \end{pmatrix} = \begin{pmatrix} 0 \\ 2 \\ 3 \end{pmatrix}$

## Autoencoder

## Discussion

non-linear "PCA"  $\Leftrightarrow a = g(w^T x + b)$

PCA  $\Leftrightarrow a = w^T x + b$



- 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.





# Kernel PCA

## Discussion

- A kernel can be applied before finding the principal components.

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

back  
6:30

feature map

SVM

SVM  
 $\varphi = \infty$   
dim

- The principal components can be found without explicitly computing  $\varphi(x_i)$ , similar to the kernel trick for support vector machines.
- Kernel PCA is a non-linear dimensionality reduction method.

K