

## ML w8 Unsupervised learning

Dataset don't have "Y" (label).

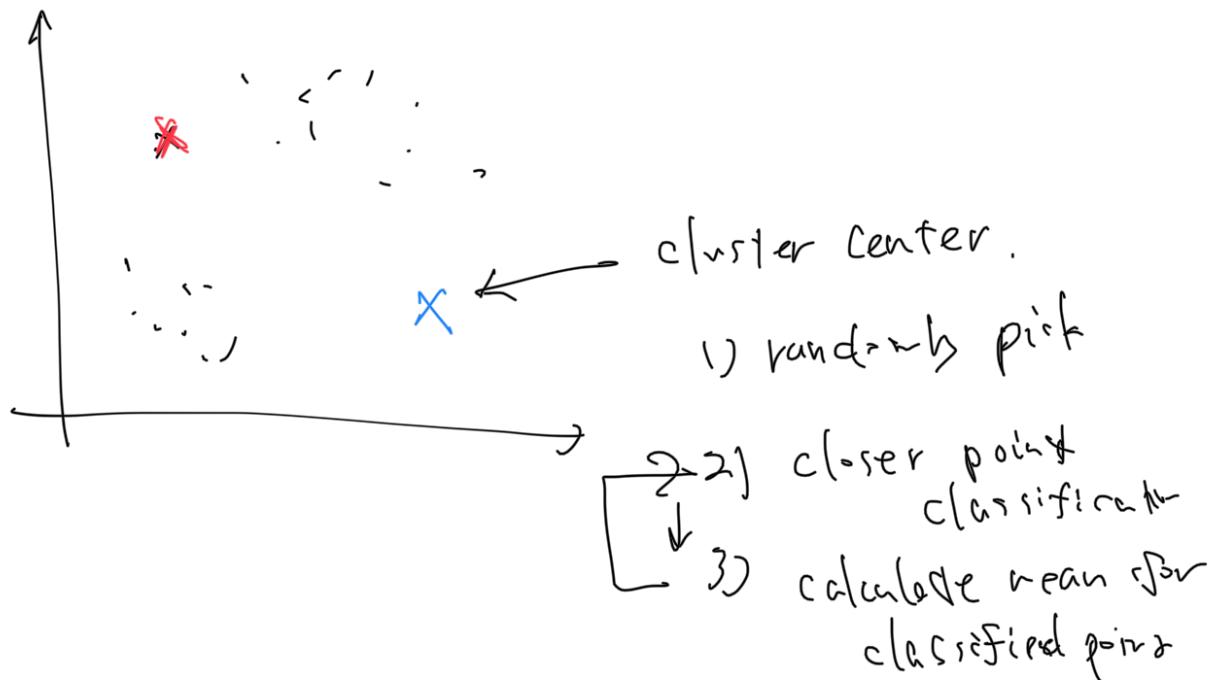
Training set  $\{x^{(1)}, x^{(2)}, x^{(3)} \dots x^{(n)}\}$

example.

- Market segmentation
- SN analysis.
- Organize computing clusters
- Astronomical data analysis.

---

K-Means Algorithm.



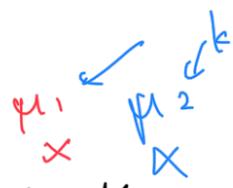
K-means algorithm

Input:

1. n 1 n - 1 ... 1

-  $K$  (number of Clusters)

- Training set



Randomly initialize  $k$  cluster centroids  $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$

cluster assignment [ for  $i=1 \dots m$

step       $c^{(i)} := \text{index} (\text{from } 1 \text{ to } K) \text{ of cluster}$

closest to  $x^{(i)}$

$$\min_{k \in c^{(i)}} \|x^{(i)} - \mu_k\|$$

square  
(distance)

$K$  : total number centroided ]

$k$  : index

$f_{ik} = 1 \text{ to } k:$

$\mu_k$  : average (mean) of points assigned to cluster  $k$

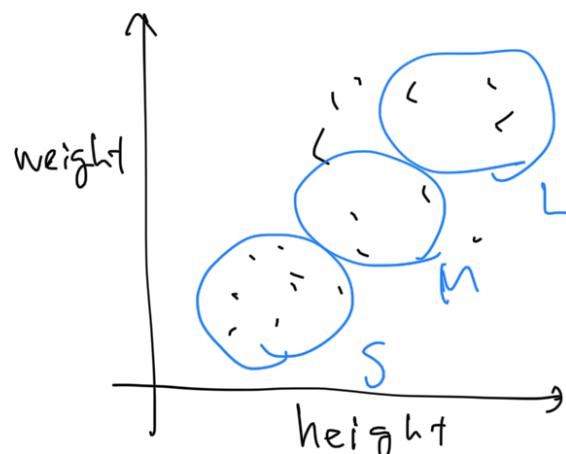
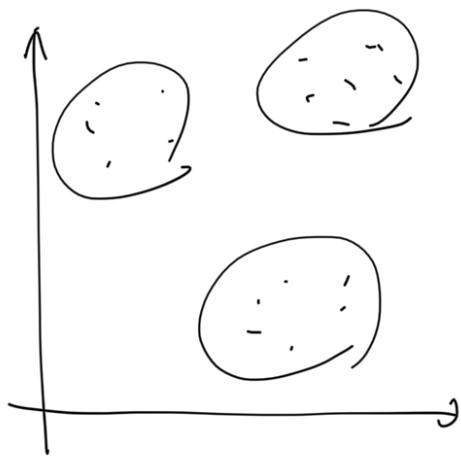
$$\text{e.g. } x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)} \rightarrow c^{(1)} = \{1\} = \{x^{(1)}\}, c^{(5)} = \{5\}$$

}

$$\mu_1 = \frac{1}{4} [x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}] \in \mathbb{R}^n$$

K means for non-separated cluster.

$S, M, L:$



## Optimization Objectives.

→  $c^{(i)}$  = index of cluster ( $1, 2, \dots, k$ ) to which example  $x^{(i)}$  is currently assigned

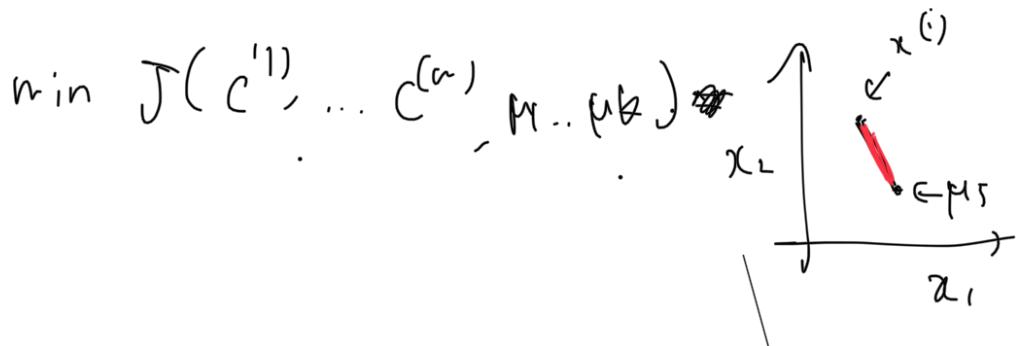
$\mu_k$  = cluster centroid  $k$  ( $\mu_k \in \mathbb{R}^n$ )

$\uparrow$  index of centroid  $k$ .

$$x^{(i)} \rightarrow 5 \quad \text{where} \quad c^{(i)} = 5$$

$$\mu_{c^{(i)}} = \underline{\mu_5}$$

$$\rightarrow J(c^{(1)}, \dots, c^{(n)}, \mu_1, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$



Randomly initialize  $K$  cluster centroids  $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$

Repeat {  
 for  $i = 1 \dots m$        $\downarrow$  minimize  $J(\dots, c^{(1)}, c^{(2)}, \dots, c^{(n)})$   
 (holding  $\mu_1 \dots \mu_k$ )

$c^{(i)}$  = index (from  $1 \dots k$ ) of cluster centroid

close to  $x^{(i)}$

for  $k = 1 \dots K$

$\mu_k$ : Average (mean) of points assigned to cluster  $k$

}

:

Random Initialization

随机数



Should have  $K < m$  e.g.  $\cancel{K^2}$



$$\mu_1 = x^{(i)}$$

$$\mu_2 = x^{(j)}$$

f, r getting global optim, make random initialization several times.

For  $i=1 \rightarrow 100 \{$

→ Randomly initialise k-means.

Ran k-means, Get  $C^{(1)} \dots C^{(n)}$ ,  $\mu_1 \dots \mu_k$

→ Compute Cost  $J$

}

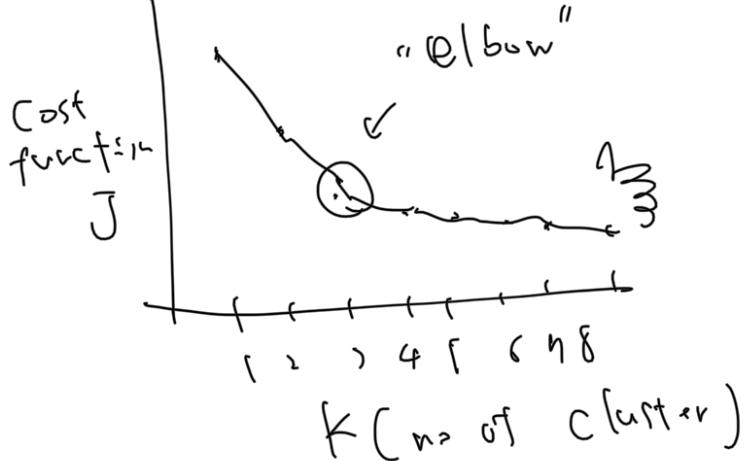
Pick clustering gave lowest cost  $J$ .

$K=2 \Rightarrow \leftarrow \text{good job.}$

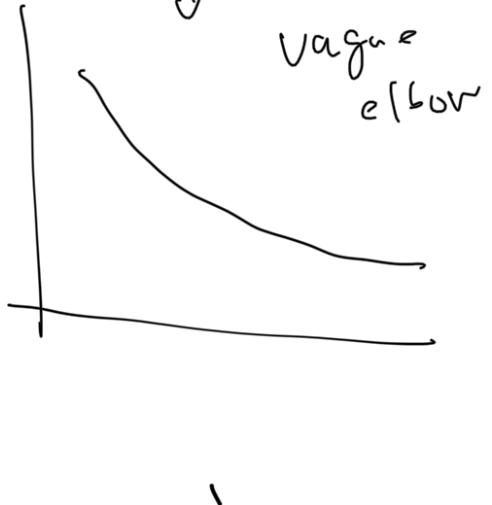
## Choosing the number of clusters

What is the right value of  $k$ ?

Elbow method

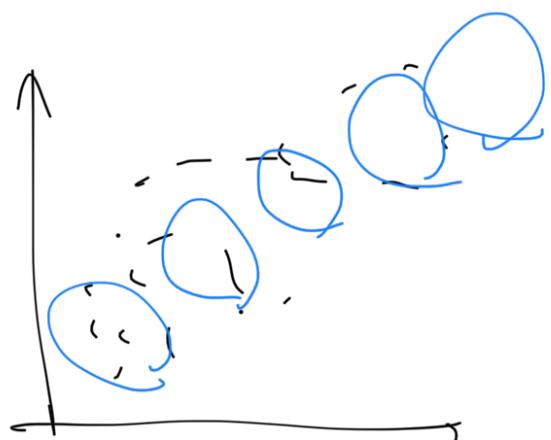
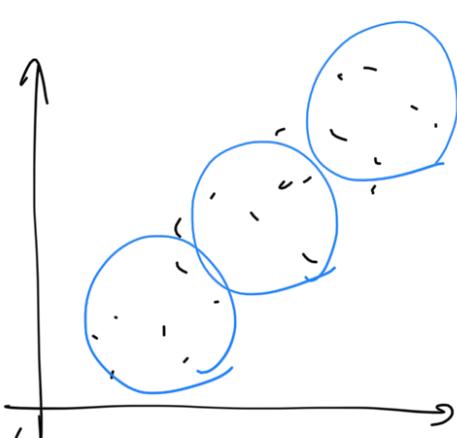


Realistic ...  
↓

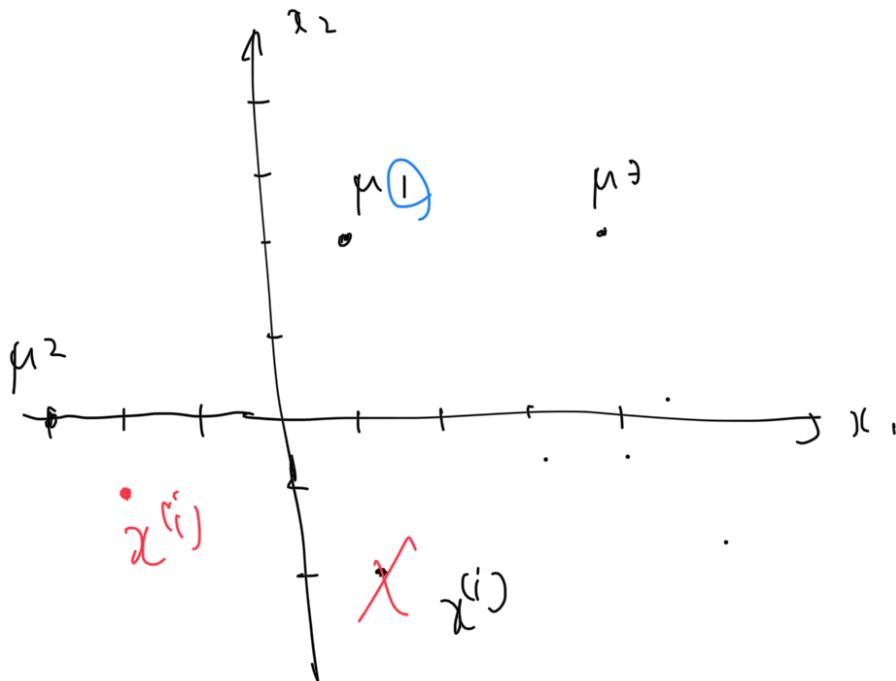


$k = 3$  (S, M, L)

$k = 5$  (XS, S, M, L, XL)



"How well it perform for that [cater \$ purpose].  
downstream  
由源，评估结果是否



# Unsupervised Learning

合計点数5

1. For which of the following tasks might K-means clustering be a suitable algorithm? Select all that apply.

1点

- Given a set of news articles from many different news websites, find out what are the main topics covered.
- From the user usage patterns on a website, figure out what different groups of users exist.
- Given historical weather records, predict if tomorrow's weather will be sunny or rainy.
- Given many emails, you want to determine if they are Spam or Non-Spam emails.

2. Suppose we have three cluster centroids  $\mu_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ ,  $\mu_2 = \begin{bmatrix} -3 \\ 0 \end{bmatrix}$  and  $\mu_3 = \begin{bmatrix} 4 \\ 2 \end{bmatrix}$ . Furthermore, we have a training example  $x^{(i)} = \begin{bmatrix} -2 \\ 1 \end{bmatrix}$ . After a cluster assignment step, what will  $c^{(i)}$  be?

1点

$c^{(i)} = 2$

$c^{(i)}$  is not assigned

$c^{(i)} = 1$

$c^{(i)} = 3$

2. Suppose we have three cluster centroids  $\mu_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ ,  $\mu_2 = \begin{bmatrix} -3 \\ 0 \end{bmatrix}$  and  $\mu_3 = \begin{bmatrix} 4 \\ 2 \end{bmatrix}$ . Furthermore, we have a training example  $x^{(i)} = \begin{bmatrix} -2 \\ 1 \end{bmatrix}$ . After a cluster assignment step, what will  $c^{(i)}$  be?

0 / 1点

$c^{(i)} = 2$

$c^{(i)}$  is not assigned

$c^{(i)} = 1$

$c^{(i)} = 3$

**⊗ 不正解**

$x^{(i)}$  is closest to  $\mu_2$ , so  $c^{(i)} = 2$ , not 1

3. K-means is an iterative algorithm, and two of the following steps are repeatedly carried out in its inner-loop. Which two?

1点

- The cluster assignment step, where the parameters  $c^{(i)}$  are updated.
- Randomly initialize the cluster centroids.
- Test on the cross-validation set.
- Move the cluster centroids, where the centroids  $\mu_k$  are updated.

4. Suppose you have an unlabeled dataset  $\{x^{(1)}, \dots, x^{(m)}\}$ . You run K-means with 50 different random

1点

initializations, and obtain 50 different clusterings of the data. What is the recommended way for choosing which one of these 50 clusterings to use?

- Plot the data and the cluster centroids, and pick the clustering that gives the most "coherent" cluster centroids.
- Manually examine the clusterings, and pick the best one.
- Compute the distortion function  $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$ , and pick the one that minimizes this.
- Use the elbow method.

5. Which of the following statements are true? Select all that apply.

1点

- The standard way of initializing K-means is setting  $\mu_1 = \dots = \mu_k$  to be equal to a vector of zeros.
- Since K-Means is an unsupervised learning algorithm, it cannot overfit the data, and thus it is always better to have as large a number of clusters as is computationally feasible.
- For some datasets, the "right" or "correct" value of K (the number of clusters) can be ambiguous, and hard even for a human expert looking carefully at the data to decide.
- If we are worried about K-means getting stuck in bad local optima, one way to ameliorate (reduce) this problem is if we try using multiple random initializations.

## Dimensionality Reduction.

- Reduce data from 2D to 1D

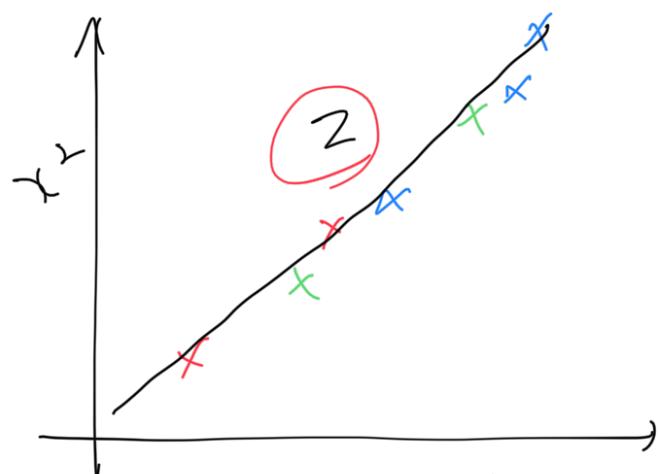
$\nearrow \uparrow$

inch centimeter

$\nwarrow \nearrow$

Round error

2D  $\rightarrow$  1D



Reduce from 2D to 1D.

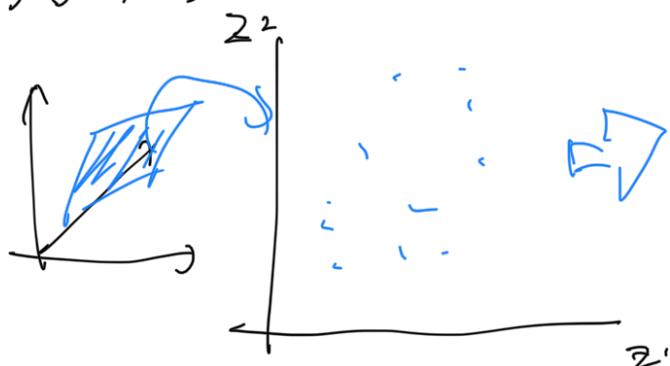
$x_1$

$z_1$

$x_2$

$z_2$

3D  $\rightarrow$  2D



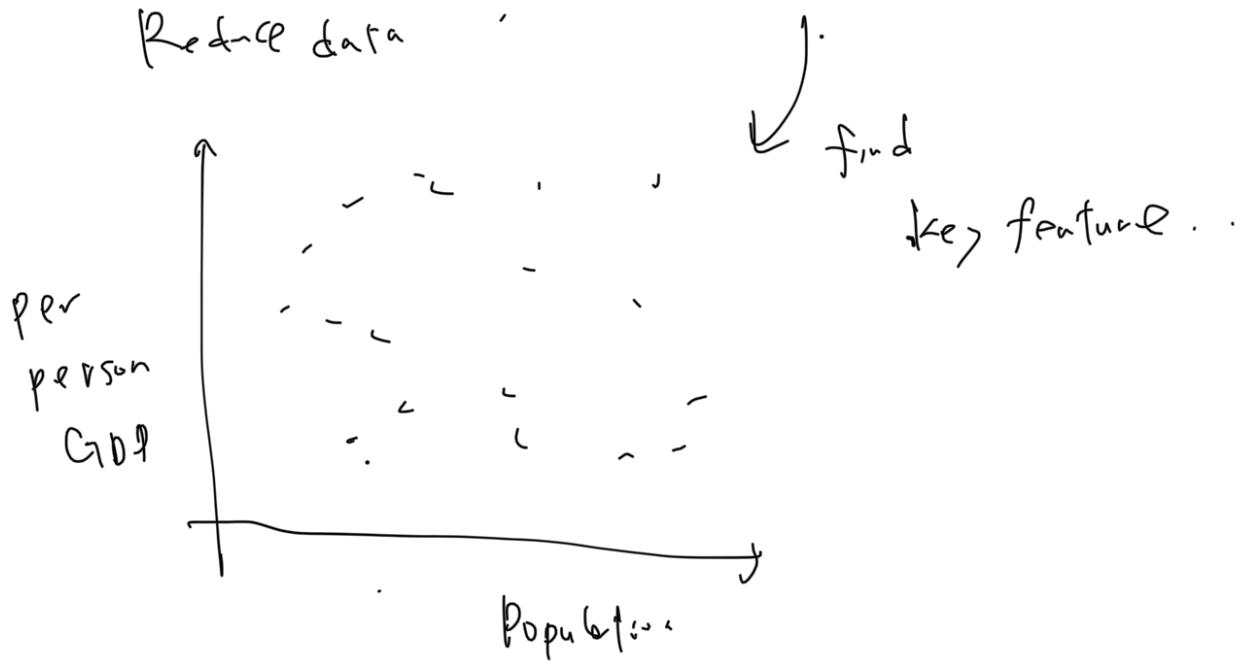
$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

$$z^{(i)} = \begin{bmatrix} z_1^{(i)} \\ z_2^{(i)} \end{bmatrix}$$

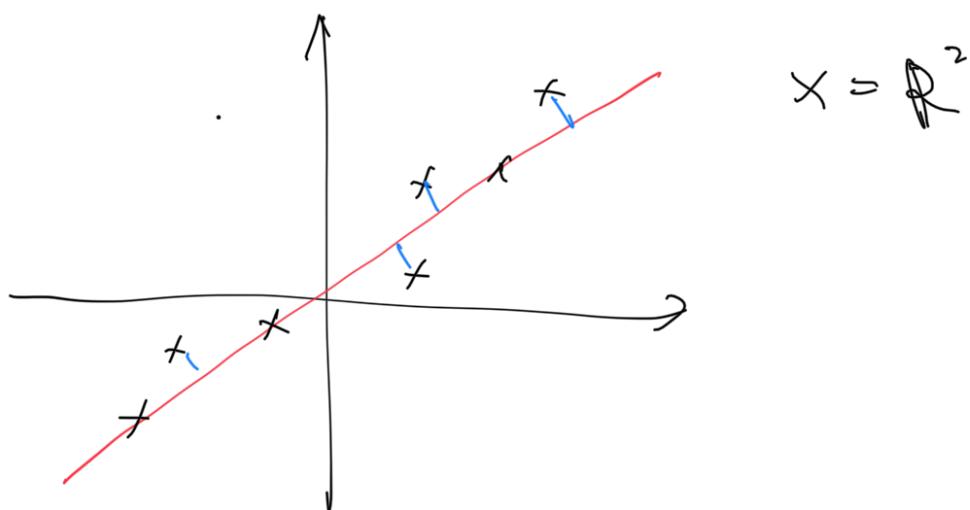
## Data Visualization

$X \in \mathbb{R}^{50}$  50 features

$z^{(i)} \in \mathbb{R}^2 \rightarrow$  plot 2D or 3D data



## PCA : Principal Component Analysis



Reduce from 2-dimensional  $\rightarrow$  1 dimension  
 find a direction (a vector  $u^{(1)} \in \mathbb{R}^n$ )  
 onto which to project the data so as to minimize  
 the projector error.

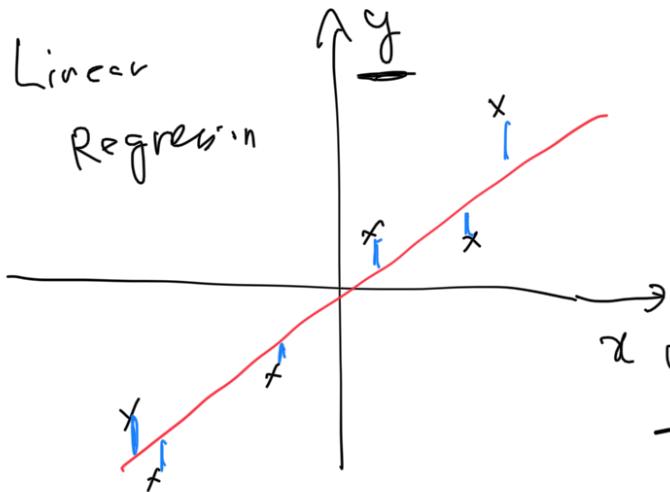
Reduce  $n$ -dimensional to k-dimension  
 vectors .. (1) (2) (3) ... (k)

onto which to project the data, so as to . . .

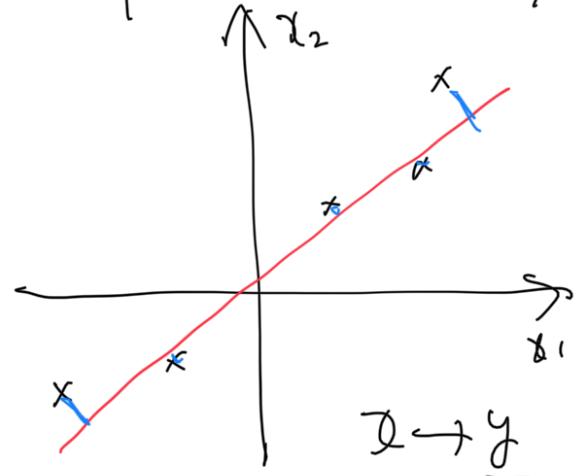
→ find k direction vectors.

PCA is not linear regression

Linear  
Regression



PCA: without  $\Sigma$



every feature is  
treated equally

## PCA - Data pre-processing

Training set:  $x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(n)}$

Pre processing features scaling, mean normalization

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

Replace each  $x_j^{(i)}$  with  $x_j^{(i)} - \mu_j$ .

If different features on large scale difference  
scale features to have comparable range of values.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j} \quad (\sigma_j) \leftarrow \text{range of values.}$$

$$x^{(i)} \in \mathbb{R}^n \rightarrow z^{(i)} \in \mathbb{R}^2$$

$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

"Covariance matrix"

$$\Sigma = \frac{1}{m} \sum_{i=1}^n (x^{(i)}) (x^{(i)})^\top$$

n x n matrix  
 特異值分解  
 svd |  
 (singular value decomposition)  
 共變方差矩阵  
 ↓  
 eig (sigma)

$$[U, \Sigma, V] = \text{svd}(\Sigma);$$

$$U = \begin{bmatrix} u^{(1)}, u^{(2)}, \dots, u^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

n x n matrix  
 ↓  
 k

$$x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$$

$$\underline{z} \stackrel{(4)}{=} \left[ u^{(1)} \ u^{(2)} \ \dots \ u^{(k)} \right]^T x \stackrel{(1)}{=} \begin{bmatrix} -u^{(1)\top} \\ -u^{(2)\top} \\ \vdots \\ -u^{(k)\top} \end{bmatrix} x$$

$n \times k$   
U reduced

↑

$$\rightarrow \underline{z} = \mathbb{R}^k \quad (\text{k dimensional vector})$$

$$\text{sigma} = (1/m) \times \underline{x}' \times \underline{x}_j$$

$$\rightarrow [U, S, V] = \text{svd}(\text{sigma});$$

$$\rightarrow U_{\text{reduce}} = U(:, 1:k);$$

$$\rightarrow \underline{z} = U_{\text{reduce}}' \times \underline{x};$$

---

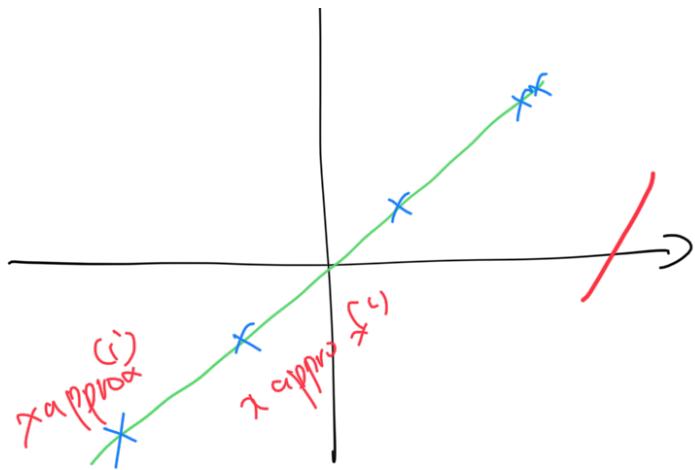

$$\underline{z} = \left[ \underline{u}^{(1)} \ \underline{u}^{(2)} \ \dots \ \underline{u}^{(k)} \right]^T x = \begin{bmatrix} -(u^{(1)})^\top \\ -(u^{(2)})^\top \\ \vdots \\ -(u^{(k)})^\top \end{bmatrix} x$$

$$\underline{z}_j = (u^{(j)})^\top x$$


---

Reconstruction from compressed representation





$$z \in \mathbb{R} \rightarrow x \in \mathbb{R}^2$$

$$\approx z$$

$$x_{\text{approx}} \stackrel{(1)}{=} u_{\text{reduce}} \cdot z^{(1)}$$

$\uparrow$   
 $\mathbb{R}^n$

$\underbrace{\quad}_{n \times k}$        $\underbrace{\quad}_{k \times 1}$   
 $\underbrace{\quad}_{n \times 1}$

Choosing the number of principal component

Choosing " $k$ "

$$\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$$

Average squared projection error  
 $< 0.01$

$$\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

total variation in data.

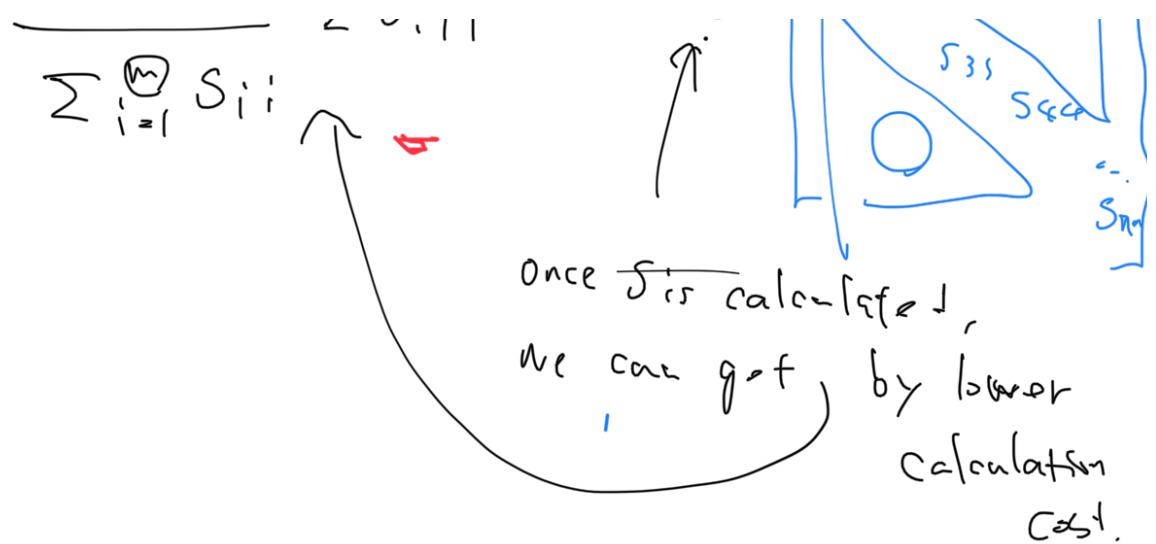
↳ "99% of variance is retained"

(Implementation)

$$[U, S, V] = \text{svd}(\text{Sigma})$$

p.s.t.  $\oplus$  smallest values of  $k$ :  
 $\sum_{i=1}^k s_{ii} > 99$

$$S = \begin{bmatrix} s_{11} & & \\ & \ddots & \\ & & s_{kk} \end{bmatrix}$$



Advice for applying PCA

Supervised learning speed up

$$\rightarrow (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots (x^{(m)}, y^{(m)})$$

Extract input:

$$\text{Unlabeled data set: } x^{(1)} x^{(2)} \dots x^{(m)} \in \mathbb{R}^{1000}$$

Firstly,

New training set

$$(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots (z^{(m)}, y^{(m)}) \in \mathbb{R}^{1000}$$

Then, cross validation and  
test sets.

$\downarrow$  PCA  $\downarrow$

$$h_\theta(z) = \frac{1}{1 + e^{-\theta^T z}}$$

$$x \rightarrow z$$

Application of PCA

Compression

- reduce memory/disk needed to store data
- speed up learning algorithm

Visualization

PCA is for speeding up algorithm

Bad use of PCA : to prevent overfitting

BAD!!



Use regularization instead

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

thus  $\uparrow$  know "w<sub>j</sub>" value.

## Design of ML System

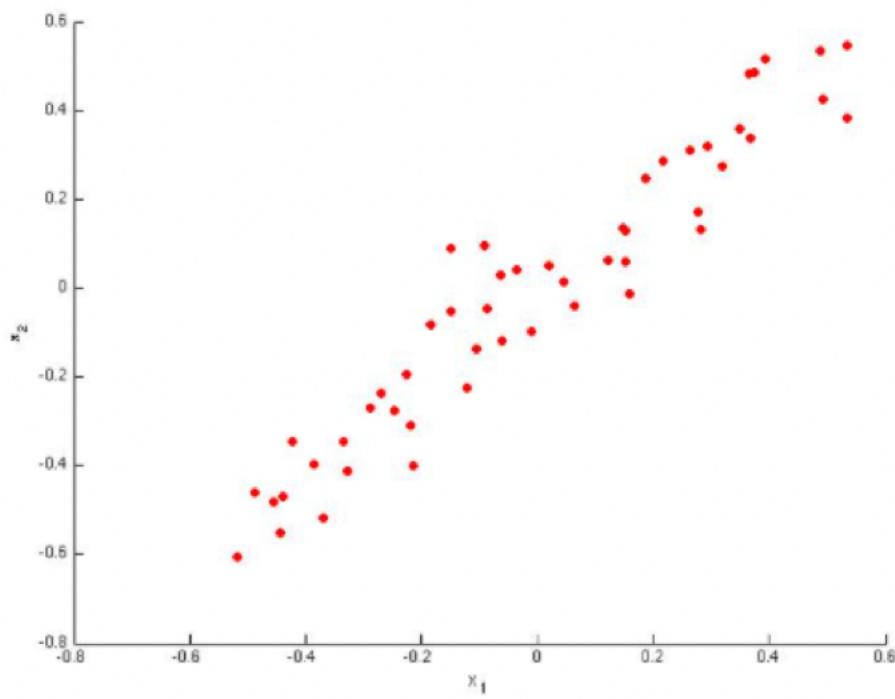
Ask: How about doing the whole thing without  
using PCA?

# Principal Component Analysis

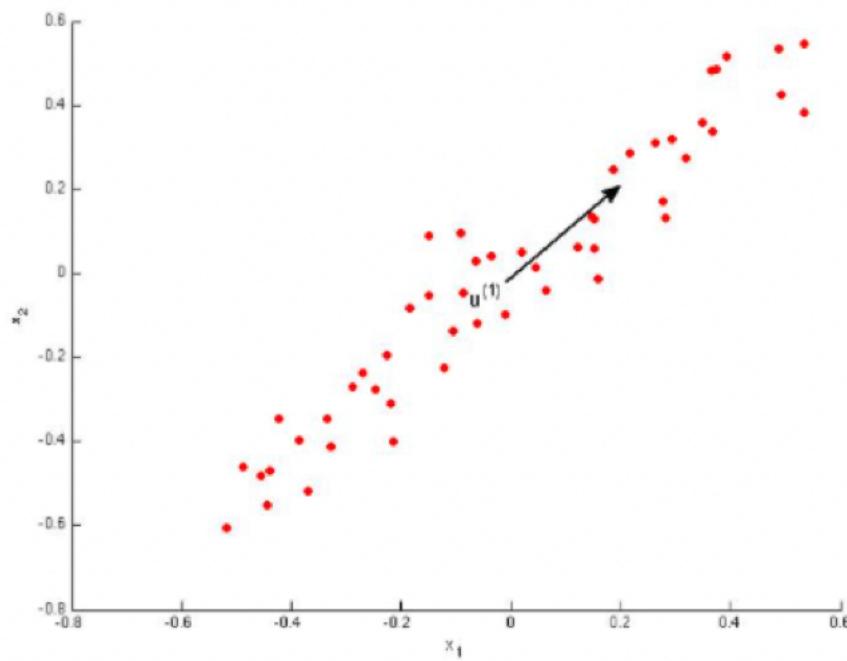
合計点数 5

1. Consider the following 2D dataset:

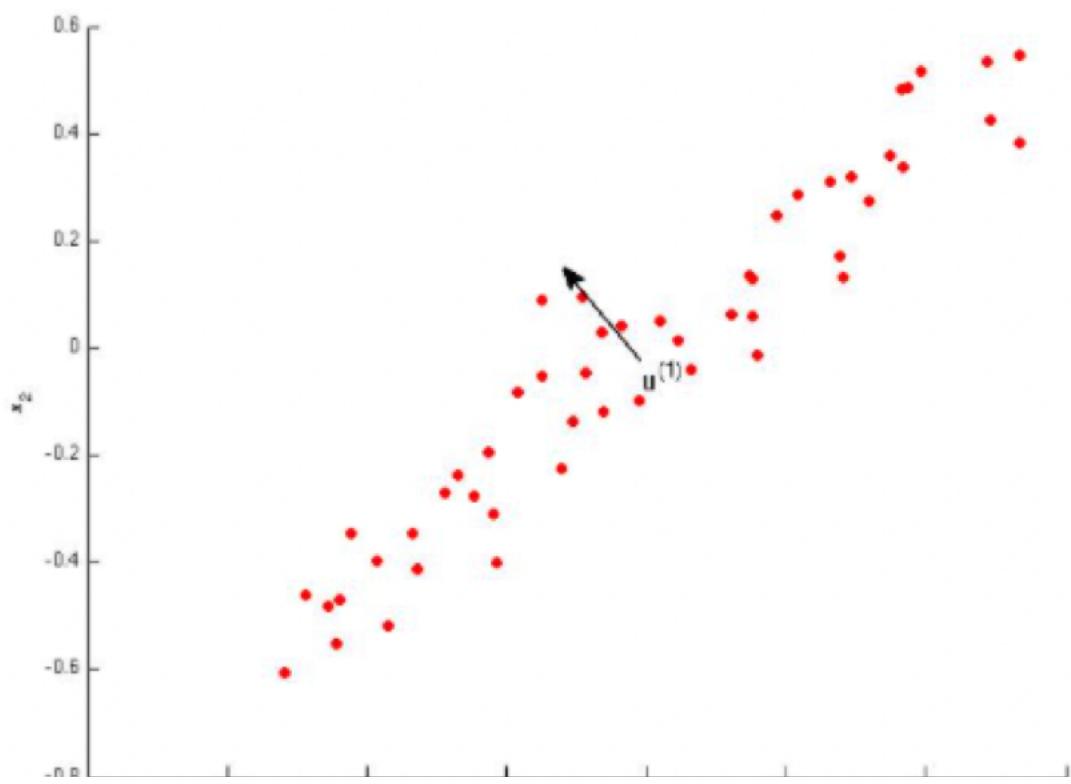
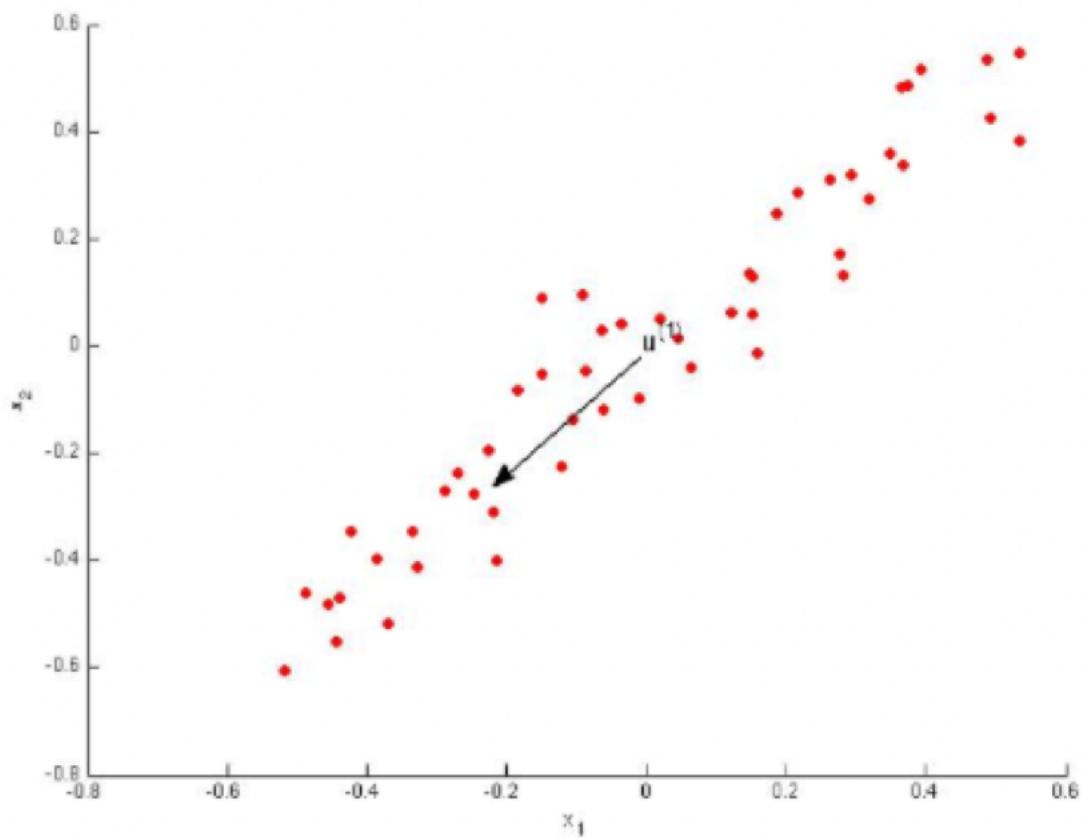
1点

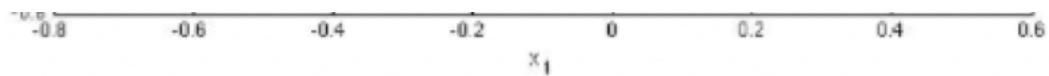


Which of the following figures correspond to possible values that PCA may return for  $u^{(1)}$  (the first eigenvector / first principal component)? Check all that apply (you may have to check more than one figure).

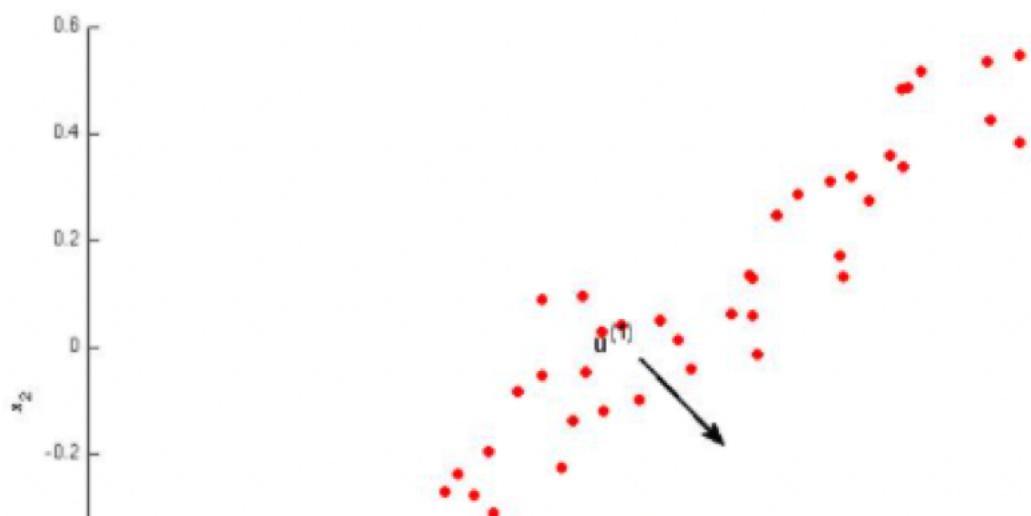








□



2. Which of the following is a reasonable way to select the number of principal components  $k$ ?

(Recall that  $n$  is the dimensionality of the input data and  $m$  is the number of input examples.)

- Use the elbow method.
  - Choose  $k$  to be the largest value so that at least 99% of the variance is retained
  - Choose  $k$  to be 99% of  $m$  (i.e.,  $k = 0.99 * m$ , rounded to the nearest integer).
  - Choose  $k$  to be the smallest value so that at least 99% of the variance is retained.
3. Suppose someone tells you that they ran PCA in such a way that "95% of the variance was retained." What is an equivalent statement to this?

- $\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \bar{x}_{\text{approx}}^{(i)}\|^2} \geq 0.95$
- $\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \bar{x}_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.05$
- $\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \bar{x}_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \geq 0.95$
- $\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \bar{x}_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \geq 0.05$

4. Which of the following statements are true? Check all that apply.

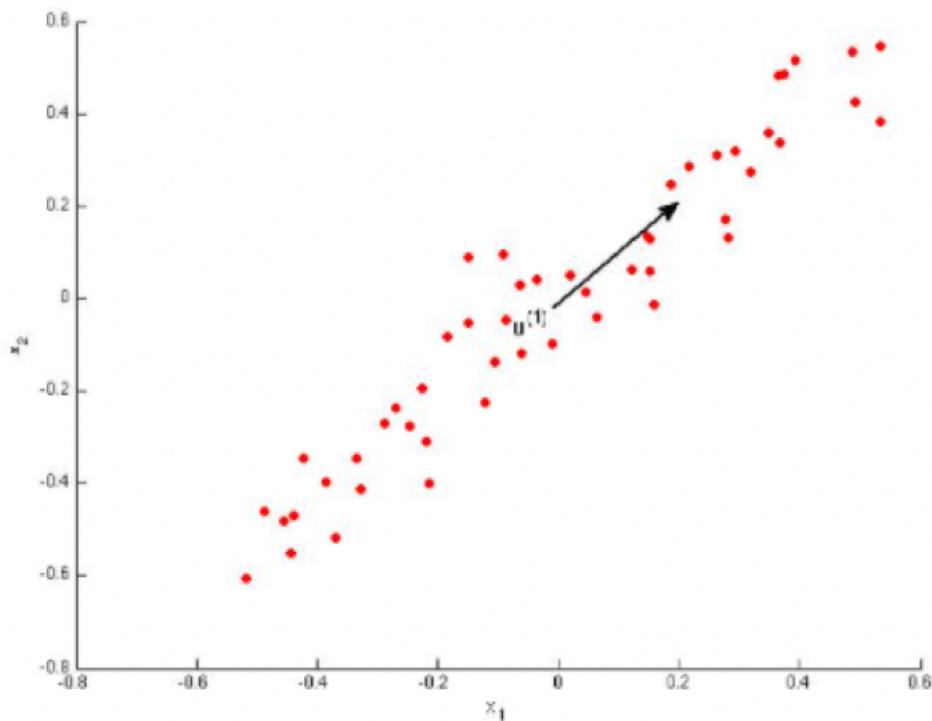
- Given only  $z^{(i)}$  and  $U_{\text{reduce}}$ , there is no way to reconstruct any reasonable approximation to  $x^{(i)}$ .
- Even if all the input features are on very similar scales, we should still perform mean normalization (so that each feature has zero mean) before running PCA.
- PCA is susceptible to local optima; trying multiple random initializations may help.
- Given input data  $x \in \mathbb{R}^n$ , it makes sense to run PCA only with values of  $k$  that satisfy  $k \leq n$ . (In particular, running it with  $k = n$  is possible but not helpful, and  $k > n$  does not make sense.)

5. Which of the following are recommended applications of PCA? Select all that apply.

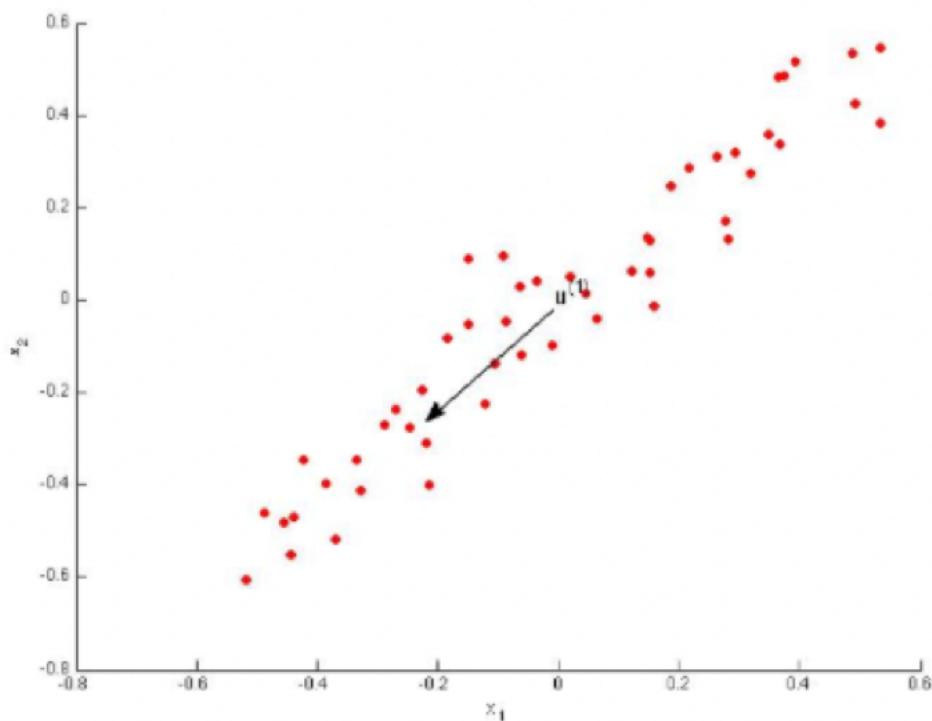
1点

- Data compression: Reduce the dimension of your input data  $x^{(i)}$ , which will be used in a supervised learning algorithm (i.e., use PCA so that your supervised learning algorithm runs faster).
- Data visualization: To take 2D data, and find a different way of plotting it in 2D (using k=2).
- Data compression: Reduce the dimension of your data, so that it takes up less memory / disk space.
- As a replacement for (or alternative to) linear regression: For most learning applications, PCA and linear regression give substantially similar results.

—  
Answer



The maximal variance is along the  $y = x$  line, so this option is correct.



✓ 正解

The maximal variance is along the  $y = x$  line, so the negative vector along that line is correct for the first principal component.

2. Which of the following is a reasonable way to select the number of principal components  $k$ ?

1 / 1点

(Recall that  $n$  is the dimensionality of the input data and  $m$  is the number of input examples.)

- Use the elbow method.
- Choose  $k$  to be the largest value so that at least 99% of the variance is retained
- Choose  $k$  to be 99% of  $m$  (i.e.,  $k = 0.99 * m$ , rounded to the nearest integer).
- Choose  $k$  to be the smallest value so that at least 99% of the variance is retained.



This is correct, as it maintains the structure of the data while maximally reducing its dimension.

3. Suppose someone tells you that they ran PCA in such a way that "95% of the variance was retained." What is an equivalent statement to this?

1 / 1点

- $\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \bar{x}_{\text{approx}}^{(i)}\|^2} \geq 0.95$
- $\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \bar{x}_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.05$
- $\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \bar{x}_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \geq 0.95$
- $\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \bar{x}_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \geq 0.05$



This is the correct formula.

4. Which of the following statements are true? Check all that apply.

1 / 1点

- Given only  $z^{(i)}$  and  $U_{\text{reduce}}$ , there is no way to reconstruct any reasonable approximation to  $x^{(i)}$ .
- Even if all the input features are on very similar scales, we should still perform mean normalization (so that each feature has zero mean) before running PCA.



If you do not perform mean normalization, PCA will rotate the data in a possibly undesired way.

- PCA is susceptible to local optima; trying multiple random initializations may help.
- Given input data  $x \in \mathbb{R}^n$ , it makes sense to run PCA only with values of  $k$  that satisfy  $k \leq n$ . (In particular, running it with  $k = n$  is possible but not helpful, and  $k > n$  does not make sense.)

Ⓐ 正解

The reasoning given is correct: with  $k = n$ , there is no compression, so PCA has no use.

5. Which of the following are recommended applications of PCA? Select all that apply.

- Data compression: Reduce the dimension of your input data  $x^{(i)}$ , which will be used in a supervised learning algorithm (i.e., use PCA so that your supervised learning algorithm runs faster).
- Data visualization: To take 2D data, and find a different way of plotting it in 2D (using k=2).

ⓧ これを選択しないでください

You should use PCA to visualize data with dimension higher than 3, not data that you can already visualize.

- Data compression: Reduce the dimension of your data, so that it takes up less memory / disk space.

Ⓑ 正解

If memory or disk space is limited, PCA allows you to save space in exchange for losing a little of the data's information. This can be a reasonable tradeoff.

- As a replacement for (or alternative to) linear regression: For most learning applications, PCA and linear regression give substantially similar results.