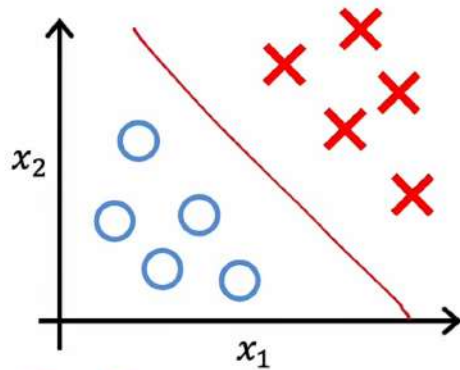


# Unsupervised Learning

## Clustering 聚类

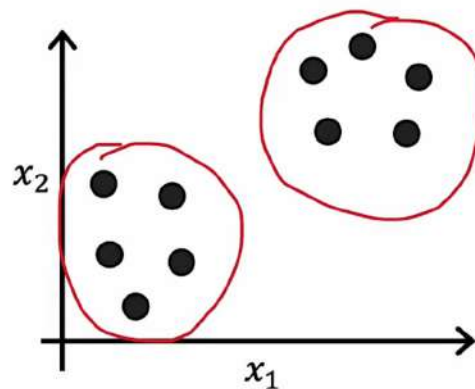
What is clustering

### Supervised learning



Training set:  $\{(\underline{x^{(1)}}), \underline{y^{(1)}}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$  ?

### Unsupervised learning



Clustering

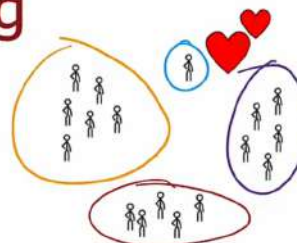
Training set:  $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$

### Applications of clustering

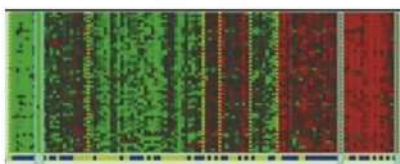


Grouping similar news

- Growing skills
- Develop career
- Stay updated with AI, understand how it affects your field of work



Market segmentation



DNA analysis



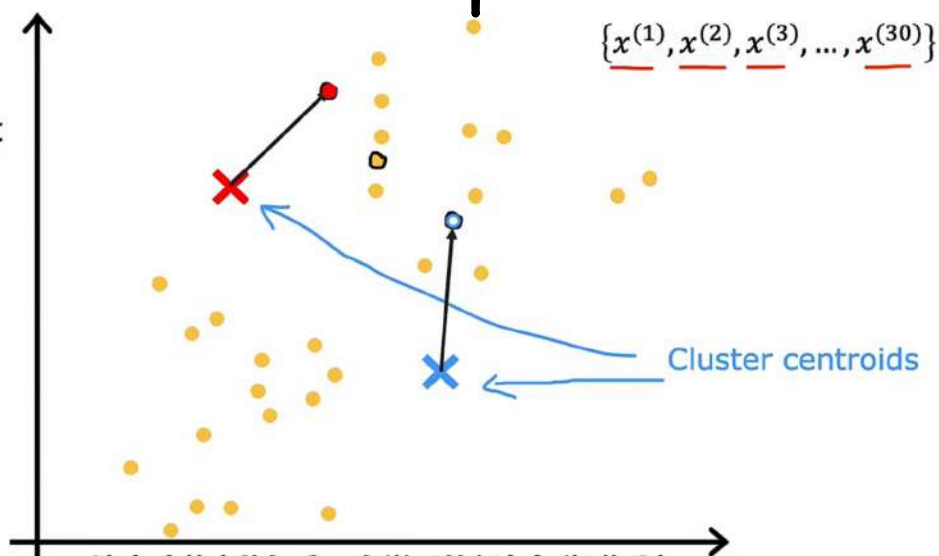
Image credit: NASA/JPL-Caltech/E. Churchwell (Univ. of Wisconsin, Madison)

Astronomical data analysis

# K-means intuition Two steps

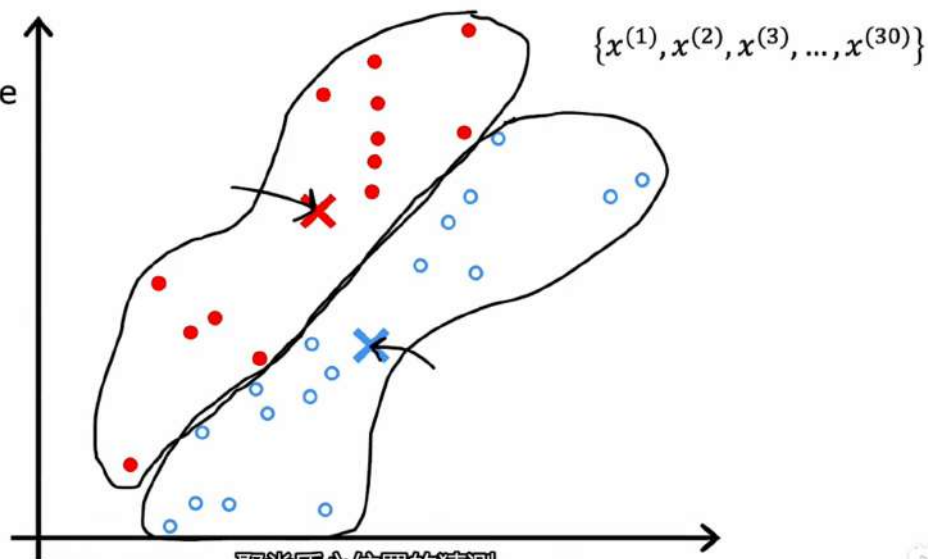
## Step 1:

Assign each point to its closest centroid



## Step 2:

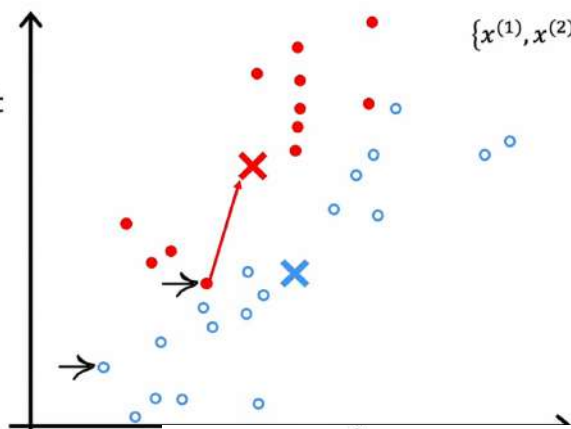
Recompute the centroids



聚类质心位置的猜测

## Step 1:

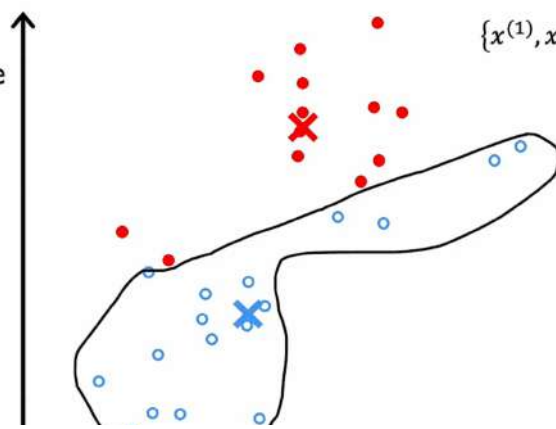
Assign each point to its closest centroid



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

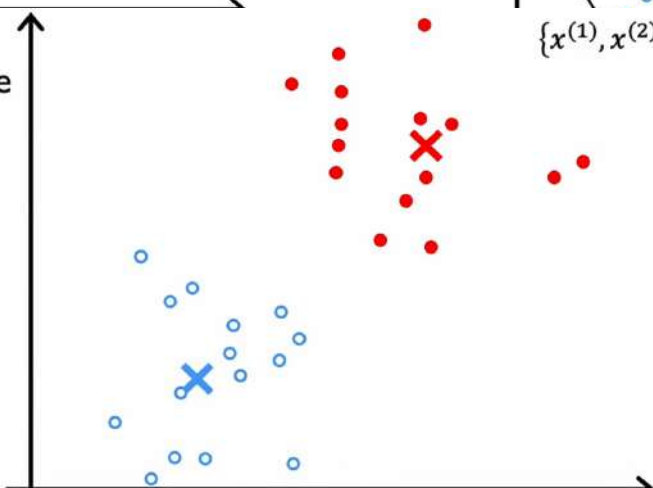
## Step 2:

Recompute the centroids



## Step 2:

Recompute the centroids



# K-means algorithm

## K-means algorithm

Randomly initialize  $K$  cluster centroids  $\mu_1, \mu_2, \dots, \mu_K$

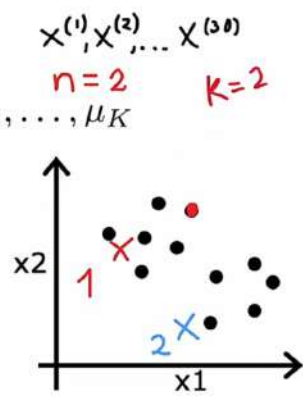
Repeat {

# Assign points to cluster centroids

for  $i = 1$  to  $m$

$c^{(i)} :=$  index (from 1 to  $K$ ) of cluster centroid closest to  $x^{(i)}$

$\min_k \|x^{(i)} - \mu_k\|^2$



## K-means algorithm

Randomly initialize  $K$  cluster centroids  $\mu_1, \mu_2, \dots, \mu_K$

Repeat {

# Assign points to cluster centroids

for  $i = 1$  to  $m$

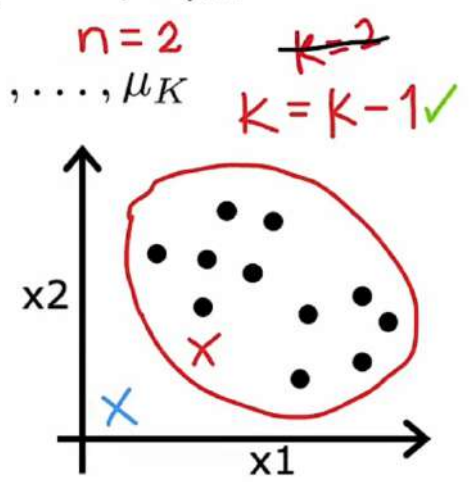
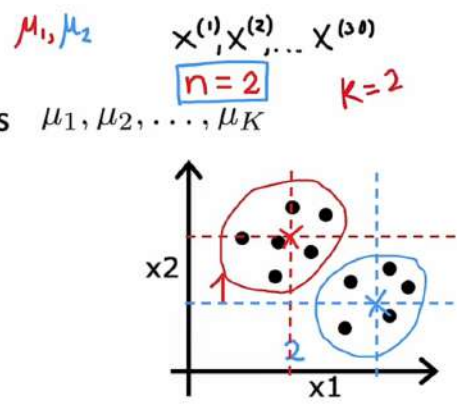
$c^{(i)} :=$  index (from 1 to  $K$ ) of cluster centroid closest to  $x^{(i)}$

# Move cluster centroids

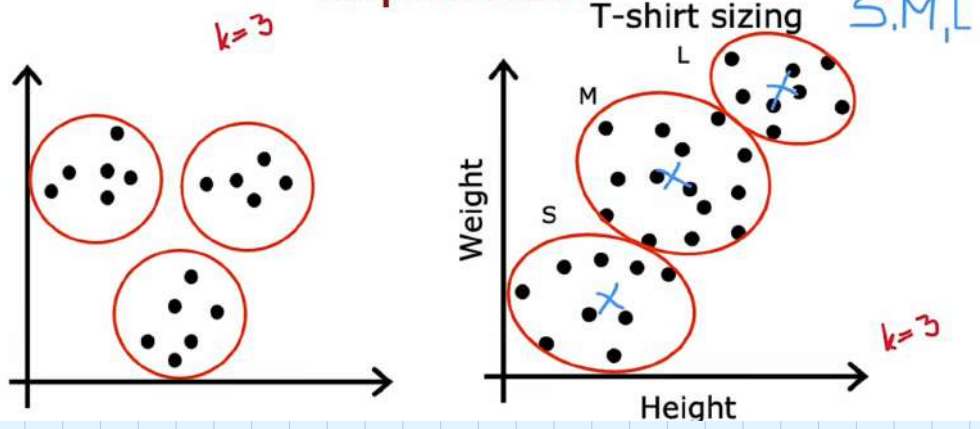
for  $k = 1$  to  $K$

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

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



## K-means for clusters that are not well separated



遇到这问题  
直接用 K

在分类还是不行

# Optimization objective 优化目标

## K-means optimization objective

$c^{(i)}$  = index of cluster (1, 2, ..., K) to which example  $x^{(i)}$  is currently assigned

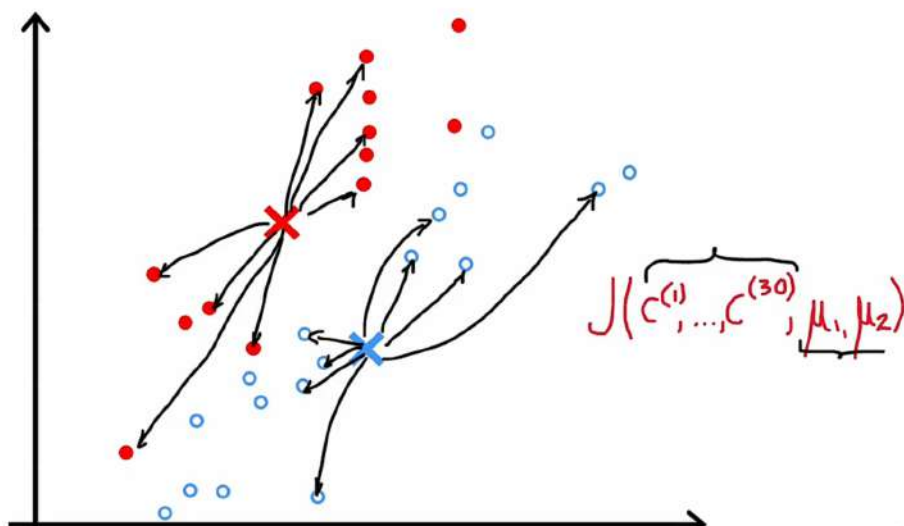
$\mu_k$  = cluster centroid  $k$

$\mu_{c^{(i)}}$  = cluster centroid of cluster to which example  $x^{(i)}$  has been assigned

Cost function

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

Distortion



## Cost function for K-means

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

Repeat {

# Assign points to cluster centroids

for  $i = 1$  to  $m$

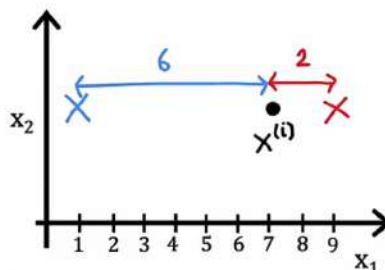
$c^{(i)} :=$  index of cluster centroid closest to  $x^{(i)}$

# Move cluster centroids

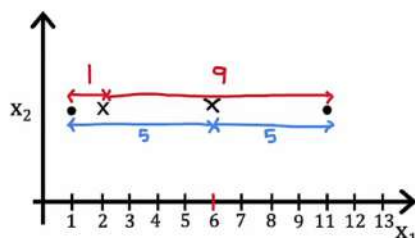
for  $k = 1$  to  $K$

$\mu_k :=$  average of points in cluster  $k$

}



## Moving the centroid



$$\frac{1}{2}(1^2 + 9^2) = \frac{1}{2}(1+81) = 41$$

$$\frac{1}{2}(1+81) = 41$$

$$\frac{1}{2}(5^2 + 5^2) = 25$$



# Initializing K-means

## K-means algorithm

Step 0: Randomly initialize  $K$  cluster centroids  $\mu_1, \mu_1, \dots, \mu_k$

Repeat {

Step 1: Assign points to cluster centroids

Step 2: Move cluster centroids

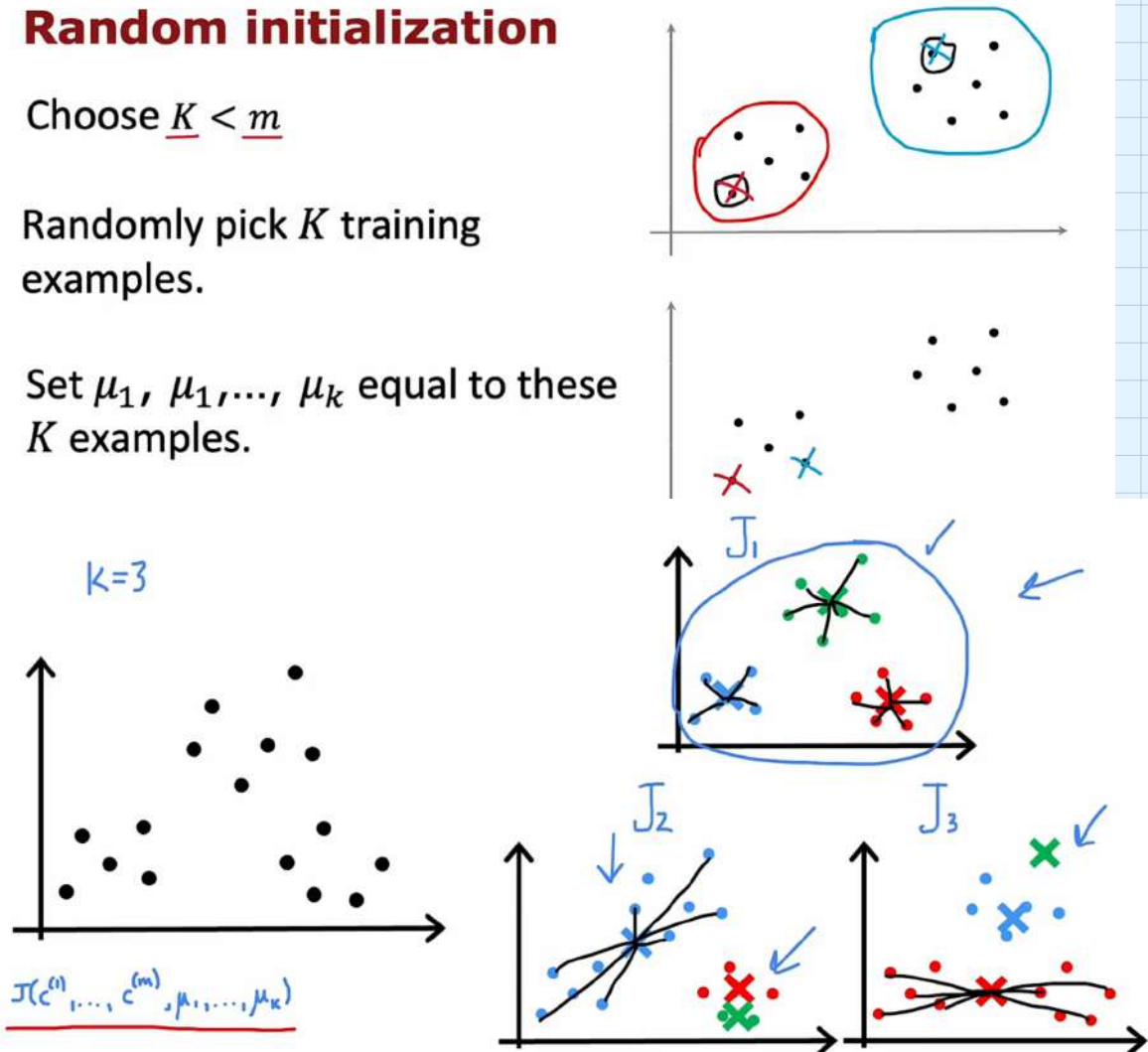
}

## Random initialization

Choose  $K < m$

Randomly pick  $K$  training examples.

Set  $\mu_1, \mu_1, \dots, \mu_k$  equal to these  $K$  examples.



多個  
找  
最佳  
局部  
解

## Random initialization

For  $i = 1$  to 100 { 50-1000

Randomly initialize K-means.  $k$  random examples

Run K-means. Get  $c^{(1)}, \dots, c^{(m)}, \mu_1, \mu_1, \dots, \mu_k$  ←

Computer cost function (distortion)

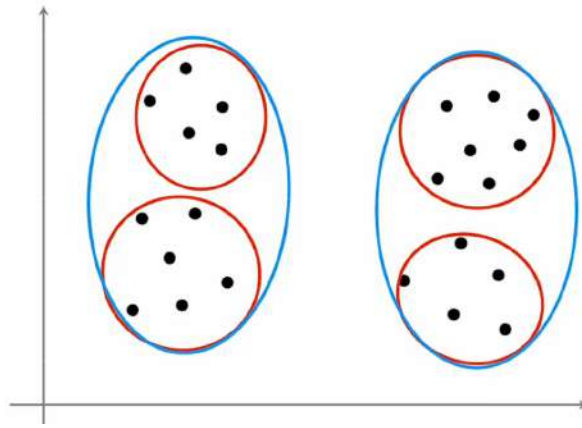
$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \mu_1, \dots, \mu_k)$  ←

}

Pick set of clusters that gave lowest cost ( $J$ )

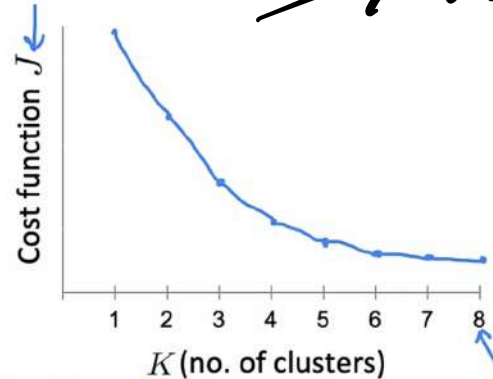
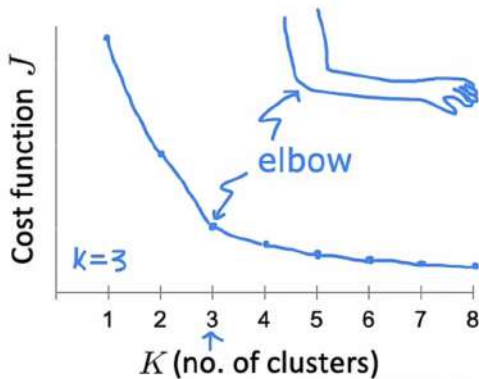
# Choosing the number of clusters

What is the right value of  $K$ ?



## Choosing the value of $K$

Elbow method:

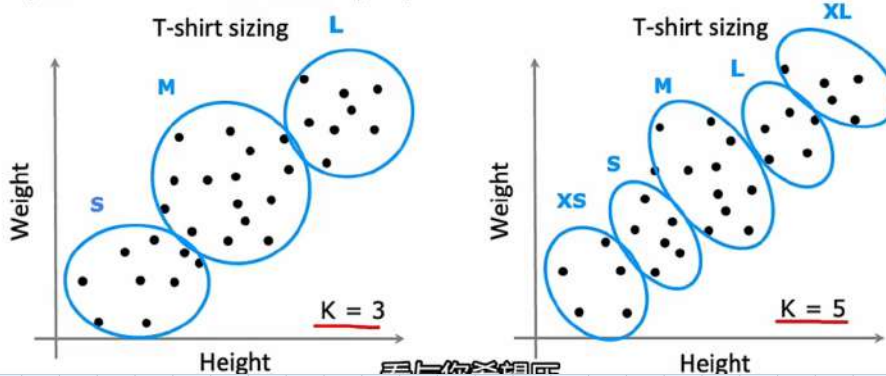


→ 不推荐

## Choosing the value of $K$

Sometimes, you're running K-means to get clusters to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose.

E.g.



→ 根据需求选择K，  
两个都跑才评估

# Second algorithm ~ Anomaly Detection

Finding unusual events

异常检测

## Anomaly detection example

Aircraft engine features:

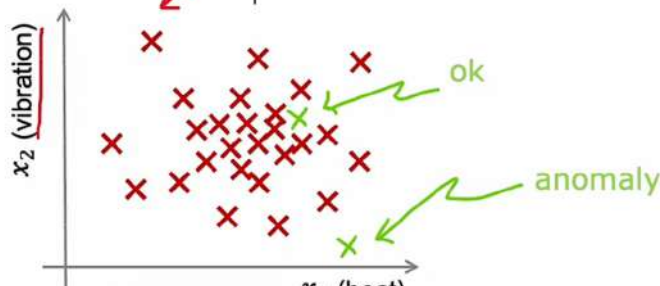
$x_1$  = heat generated

$x_2$  = vibration intensity

...

Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

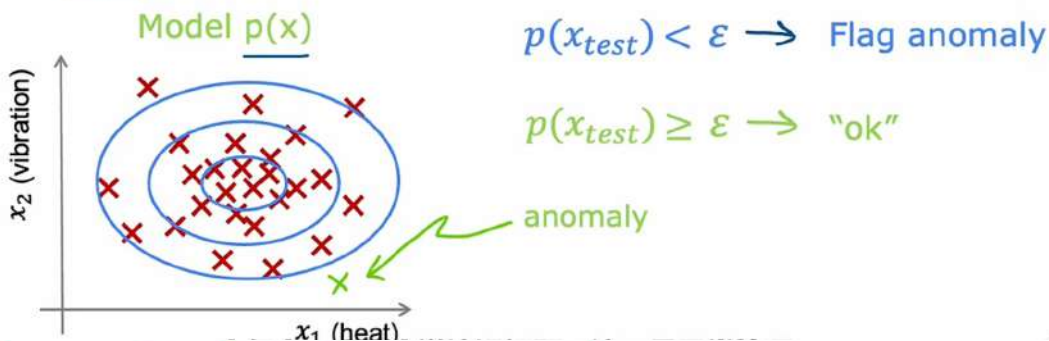
New engine:  $x_{test}$



## Density estimation

Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

Is  $x_{test}$  anomalous?



## Anomaly detection example

Fraud detection:

$\rightarrow x^{(i)}$  = features of user  $i$ 's activities

Model  $p(x)$  from data.

Identify unusual users by checking which have  $p(x) < \epsilon$

Monitoring computers in a data center.

$x^{(i)}$  = features of machine  $i$

$x_1$  = memory use,  $x_2$  = number of disk accesses/sec,

$x_3$  = CPU load,  $x_4$  = CPU load/network traffic.

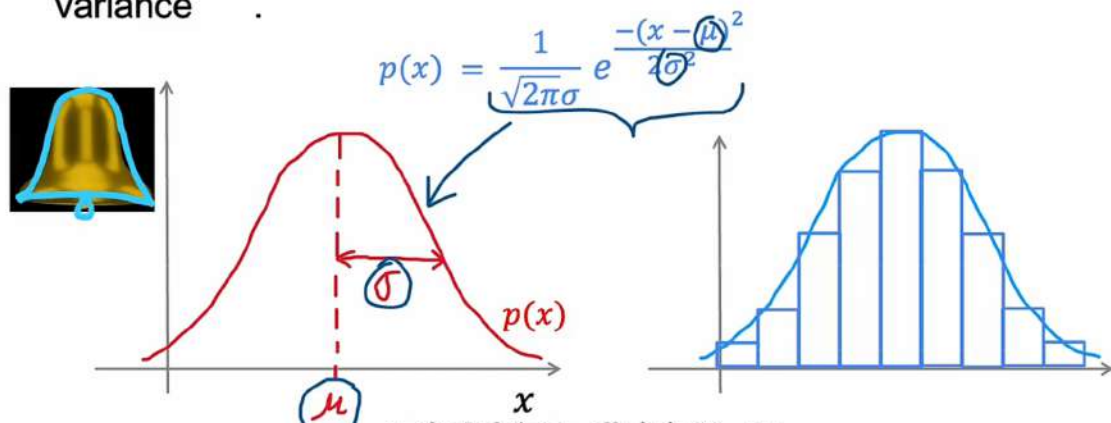
...

可以做欺诈  
防

# Gaussian (Normal) Distribution

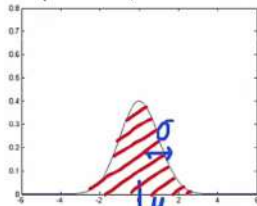
## Gaussian (Normal) distribution

Say  $x$  is a number. If  $x$  is a distributed Gaussian with mean  $\mu$ ,  $\sigma^2$  variance

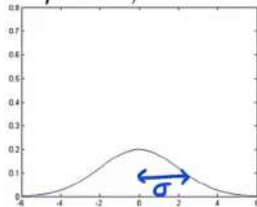


## Gaussian distribution example

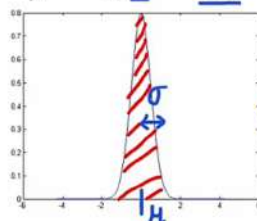
$$\mu = 0, \sigma = 1$$



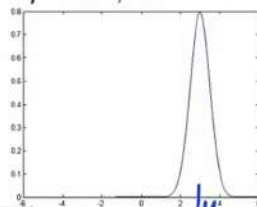
$$\mu = 0, \sigma = 2 \quad \sigma^2 = 4$$



$$\mu = 0, \sigma = 0.5 \quad \sigma^2 = 0.25$$

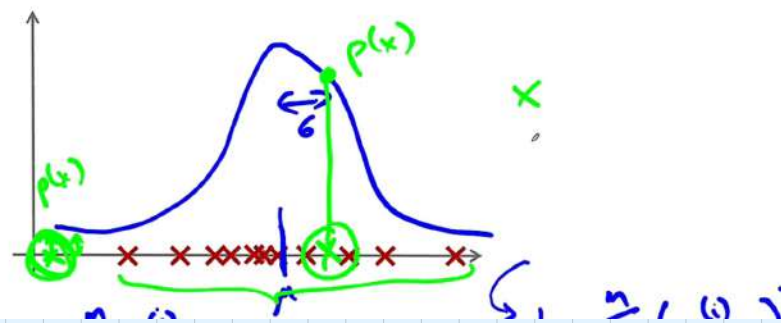
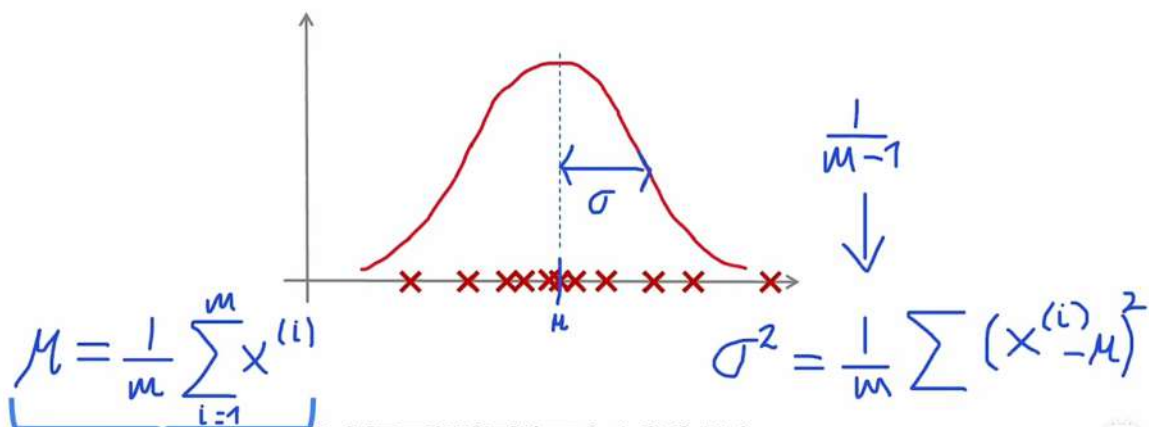


$$\mu = 3, \sigma = 0.5$$



## Parameter estimation

Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$





# Algorithm

## Density estimation

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

Each example  $x_i$  has  $n$  features

$$x = \begin{bmatrix} x_1 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$p(x) = p(x_1; \mu_1, \sigma_1^2) * p(x_2; \mu_2, \sigma_2^2) * p(x_3; \mu_3, \sigma_3^2) *$$

$$= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

$$\sum \rightarrow + \quad \prod \rightarrow \times$$

$$p(x_1 = \text{high temp}) = 1/10$$

$$p(x_2 = \text{high vibra}) = 1/20$$

$$p(x_1, x_2) = p(x_1) * p(x_2)$$

$$= \frac{1}{10} * \frac{1}{20} = \frac{1}{200}$$

## Anomaly detection algorithm

1. Choose  $n$  features  $x_i$  that you think might be indicative of anomalous examples.

2. Fit parameters  $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

Vectorized formula

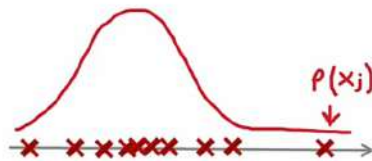
$$\vec{\mu} = \frac{1}{m} \sum_{i=1}^m \vec{x}^{(i)} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

## Anomaly detection algorithm

1. Choose  $n$  features  $x_i$  that you think might be indicative of anomalous examples.

2. Fit parameters  $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

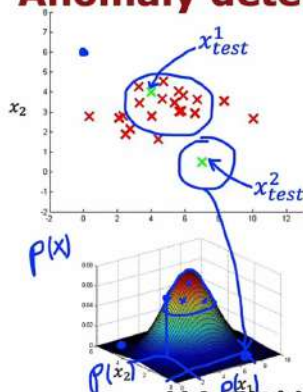
$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$



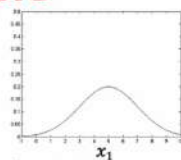
3. Given new example  $x$ , compute  $p(x)$ :

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

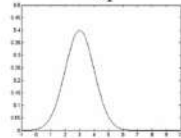
## Anomaly detection example



$$\begin{aligned} \mu_1 &= 5, \sigma_1^2 = 2 \\ \mu_2 &= 3, \sigma_2^2 = 1 \end{aligned}$$



$$p(x_1; \mu_1, \sigma_1^2)$$



$$p(x_2; \mu_2, \sigma_2^2)$$

$$\epsilon = 0.02$$

$$p(x_{test}^{(1)}) = 0.0426$$

"ok"

$$p(x_{test}^{(2)}) = 0.0021$$

"anomaly"

# Developing and evaluating an anomaly detection system

## The importance of real-number evaluation

When developing a learning algorithm (choosing features, etc.), making decisions is much easier if we have a way of evaluating our learning algorithm.

Assume we have some labeled data, of anomalous and non-anomalous examples. ( $y = 0$  if normal,  $y = 1$  if anomalous).

Training set:  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$  (assume normal examples/not anomalous)  $y=0$

Cross validation set:  $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$  } Include a few anomalous examples

Test set:  $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$  }  $y=1$

## Aircraft engines monitoring example

10000 good (normal) engines  $y=0$  2-50  
→ 20 flawed engines (anomalous)  $y=1$

(Training set: 6000 good engines  $y=0$  Train algorithm  
(CV: 2000 good engines ( $y=0$ ), 10 anomalous ( $y=1$ )  $\epsilon$   $x_j$   
Test: 2000 good engines ( $y=0$ ), 10 anomalous ( $y=1$ )

Alternative:

→ Training set: 6000 good engines 2  
→ CV: 4000 good engines ( $y=0$ ), 20 anomalous ( $y=1$ )  
No test set  $\epsilon, x_i$

## Algorithm evaluation

→ Fit model  $p(x)$  on training set  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

On a cross validation/test example  $x$ , predict

$$y = \begin{cases} 1 & \text{if } p(x) < \epsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \epsilon \text{ (normal)} \end{cases}$$

10  
2000

Possible evaluation metrics:

- True positive, false positive, false negative, true negative
- Precision/Recall
- $F_1$ -score

-  $F_1$ -score

Can also use cross validation set to choose parameter  $\epsilon$

Andrew

# Anomaly detection vs Supervised learning

## Anomaly detection vs. Supervised learning

Very small number of positive examples ( $y = 1$ ). (0-20 is common).  
Large number of negative examples ( $y = 0$ ).  
examples.  $p(x)$

Many different "types" of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalous examples we've seen so far.

Fraud

Large number of positive and negative examples.

20 positive examples

Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set.

Spam

新欺诈类型, 则可能更适用。

防欺诈

## Anomaly detection

- Fraud detection
- Manufacturing - Finding new previously unseen defects in manufacturing. (e.g. aircraft engines)
- Monitoring machines in a data center
- ⋮

## vs. Supervised learning

- Email spam classification
- Manufacturing - Finding known, previously seen defects scratches  $y = 1$
- Weather prediction (sunny/rainy/etc.)
- Diseases classification
- ⋮

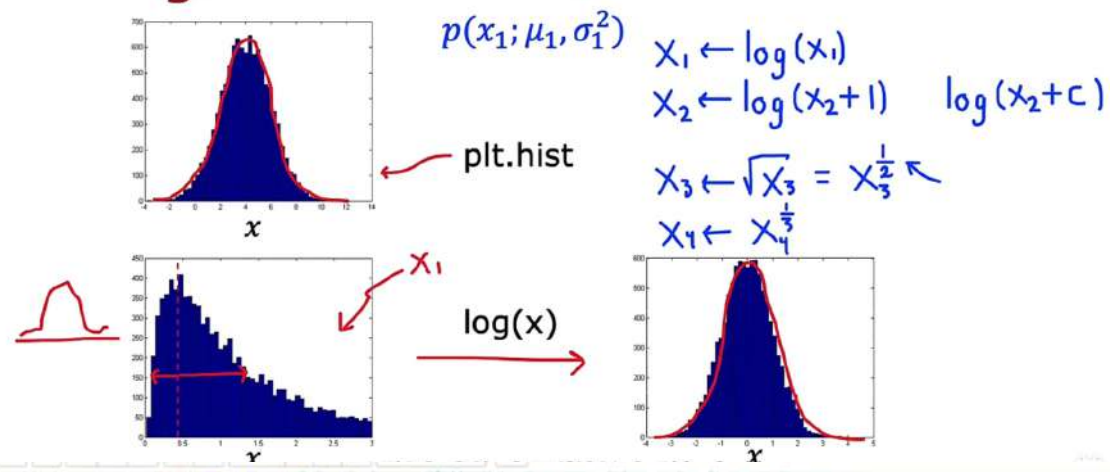
↓  
检测新东西

↑  
similar old and future



# Choosing what features to use

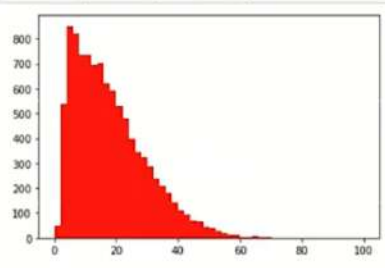
## Non-gaussian features



```
random = random - min(random) #Shift the set so the minimum value is equal to zero.
random = random / max(random) #Standardize all the vlues between 0 and 1.
random = random * maxValue    #Multiply the standardized values by the maximum value.

x = random

#Plot histogram to check skewness
plt.hist(x, bins=50, color='r');
```

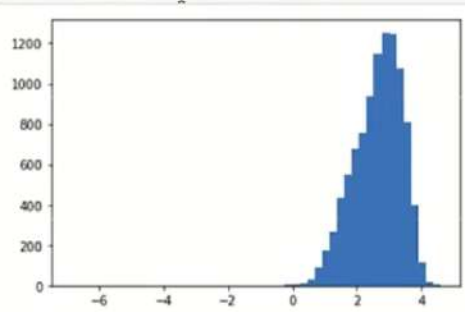


```
In [5]: plt.hist(x**0.7, bins=50);
```



再次发出嘶嘶声的一克垃圾箱，在这种情况下，它可能看起来像这样。  
in which case it might look like this.

```
In [11]: plt.hist(np.log(x+0.001), bins=50);
```



所以你会得到看起来像这样的嘶嘶声。

And so you get the blood cream that looks like a

```
In [16]: plt.hist(np.log(x+7), bins=50);
```



因此，通过在 Jupiter notebook  
So, by trying things out in Jupiter notet

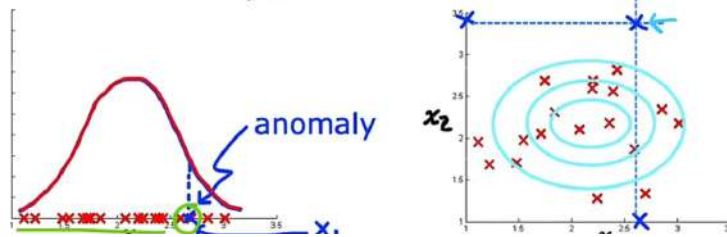


# Error analysis for anomaly detection

Want  $p(x) \gg \epsilon$  large for normal examples  $x$ .  
 $p(x) \ll \epsilon$  small for anomalous examples  $x$ .

Most common problem:

$p(x)$  is comparable (say, both large) for normal and anomalous examples



## Monitoring computers in a data center

Choose features that might take on unusually large or small values in the event of an anomaly.

$x_1$  = memory use of computer

$x_2$  = number of disk accesses/sec

$x_3$  = CPU load ←

$x_4$  = network traffic ←

$$x_5 = \frac{\text{CPU load}}{\text{network traffic}}$$

$$x_6 = \frac{(\text{CPU load})^2}{\text{network traffic}}$$

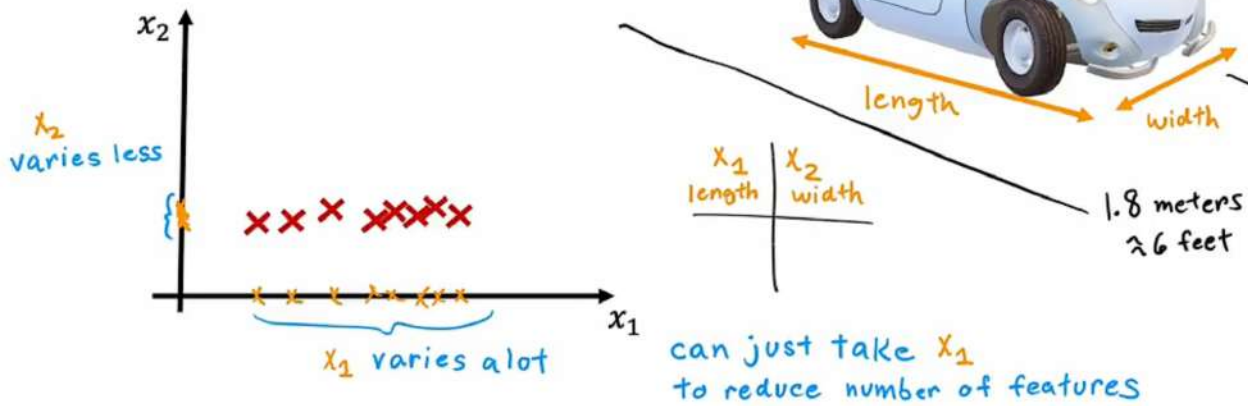
Deciding feature choice based on  $p(x)$

– large for normal examples, and becomes small for anomaly in the cross validation set

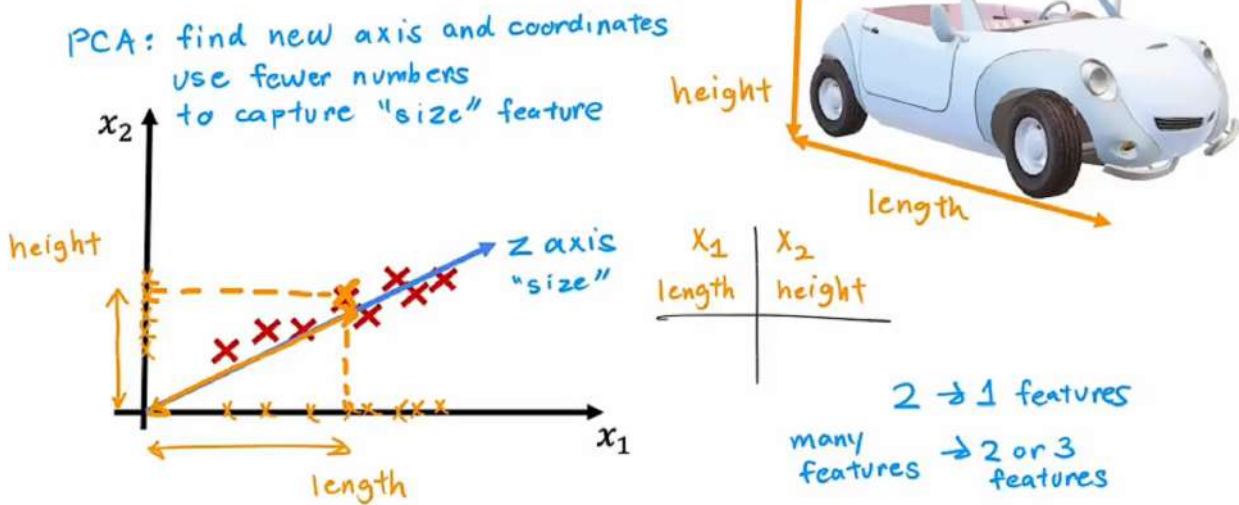
## 主成分分析 (PCA算法)

↳ 常用率可视化

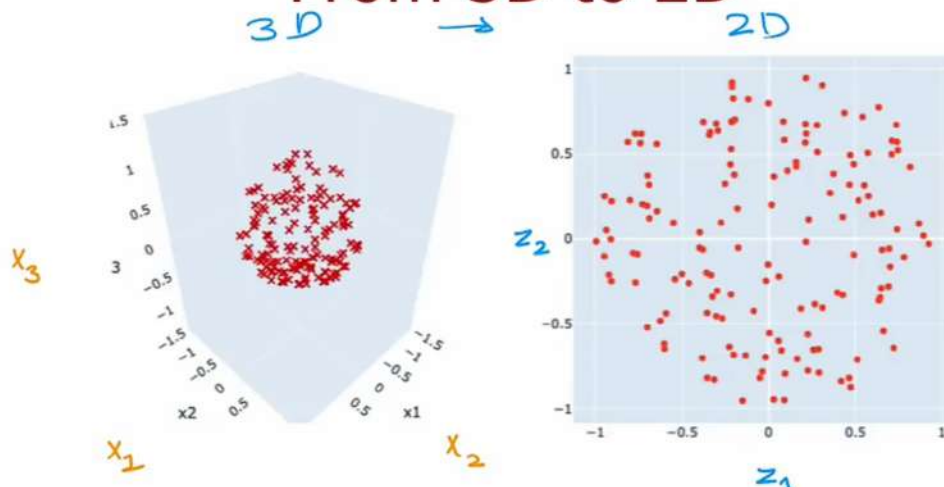
### Car measurements



### Size



### From 3D to 2D



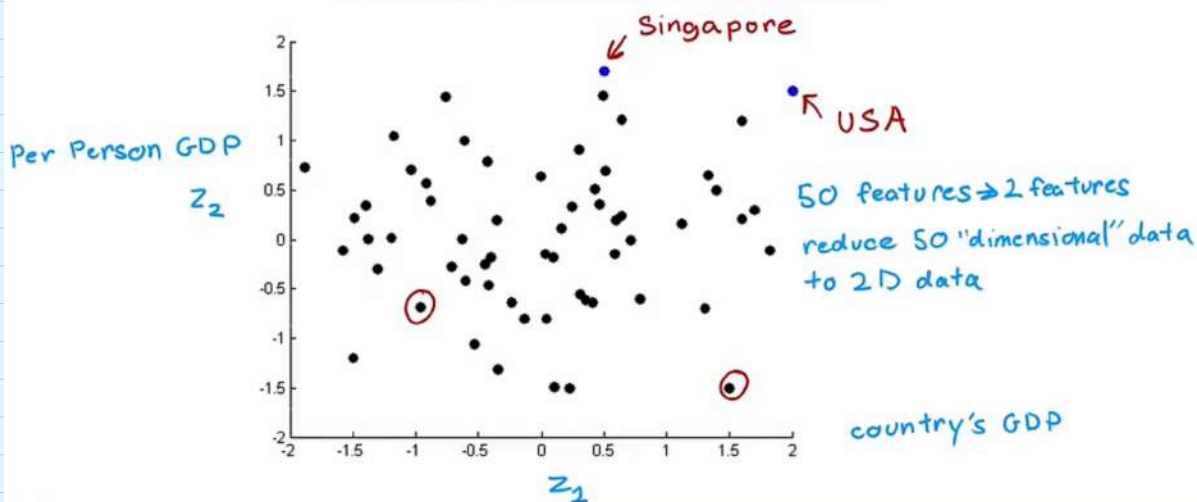
50 features

	$x_1$	$x_2$	$x_3$	$x_4$	...	
Country	GDP (trillions of US\$)	Per capita GDP (thousands of Intl. \$)	Human Development Index	Life expectancy	Poverty Index (Gini as percentage)	Mean household income (thousands of US\$)
Canada	1.577	39.17	0.908	80.7	32.6	67.293
China	5.878	7.54	0.687	73	46.9	10.22
India	1.632	3.41	0.547	64.7	36.8	0.735
Russia	1.48	19.84	0.755	65.5	39.9	0.72

Country	$z_1$	$z_2$
Canada	1.6	1.2
China	1.7	0.3
India	1.6	0.2
Russia	1.4	0.5

## Data visualization



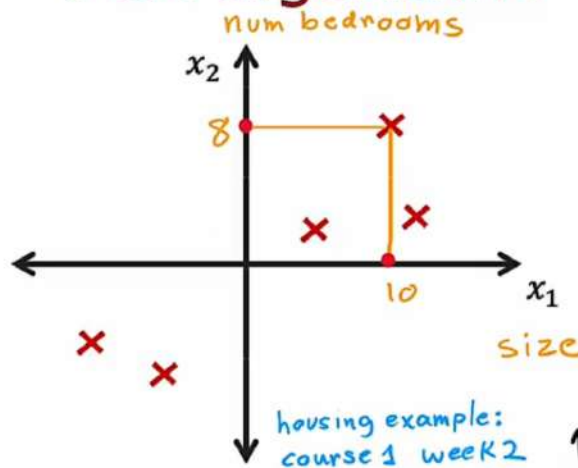
## PCA algorithm

### PCA algorithm

coordinates

$x_1=10$   $x_2=8$

Can we choose a different axis?



Preprocess features

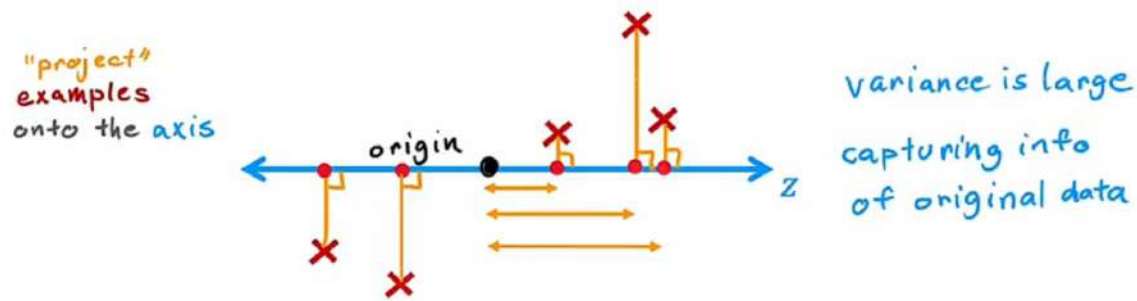
Normalized to have zero mean



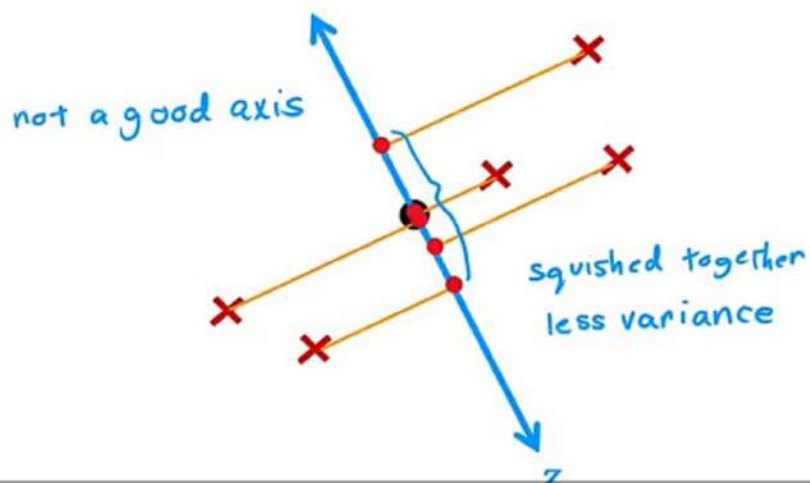
feature scaling



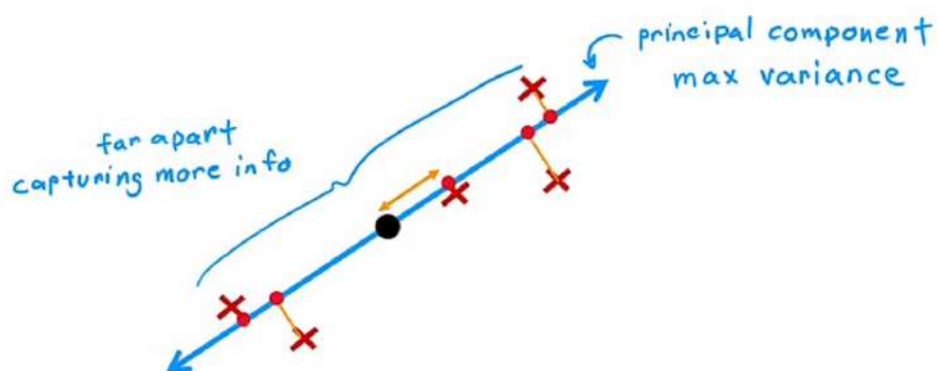
## Choose an axis



## Choose an axis



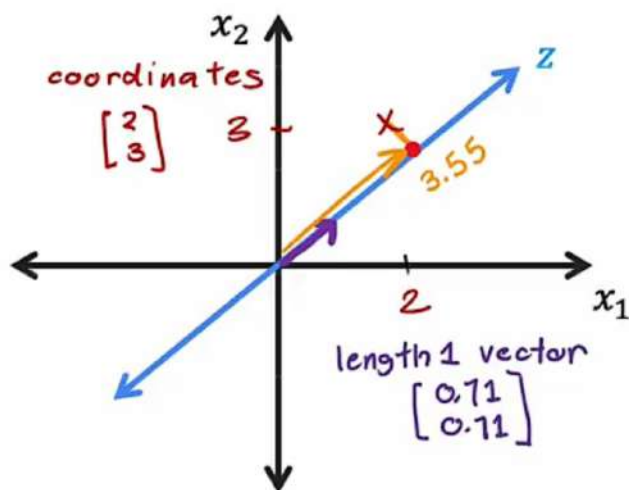
## Choose an axis



最好的选择



# Coordinate on the new axis

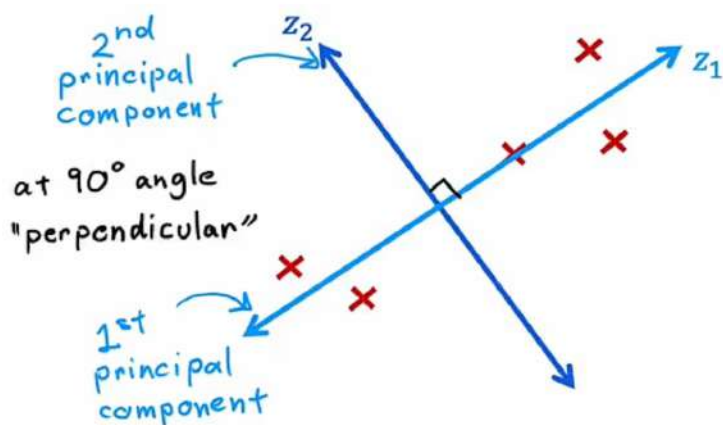


dot product

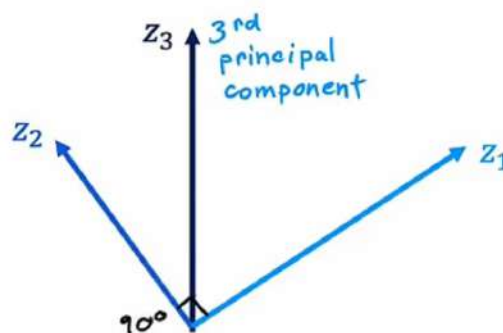
$$\begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix}$$

$$2 \times 0.71 + 3 \times 0.71 = 3.55$$

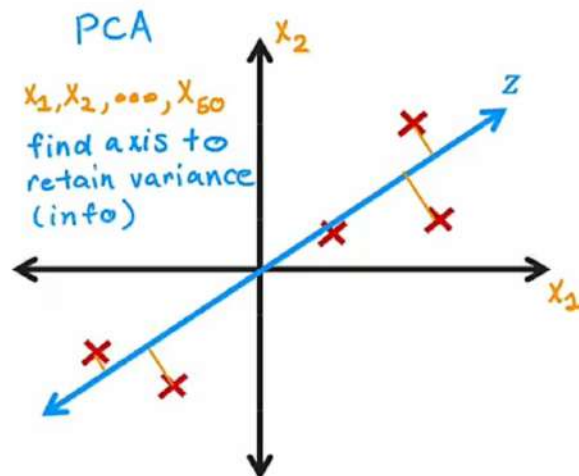
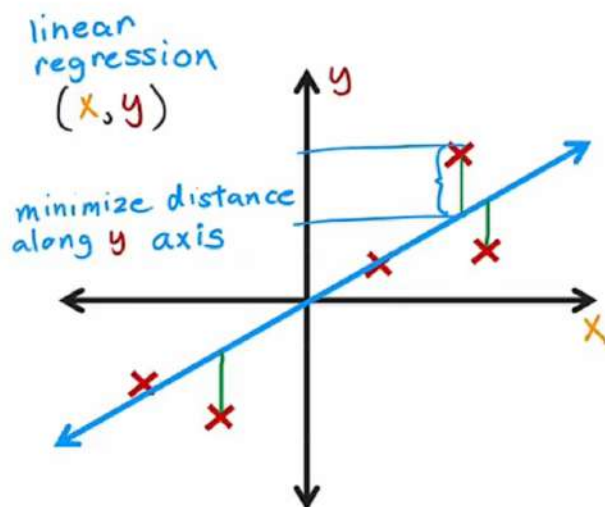
## More principal components



50 features  $\rightarrow$  3 principal components

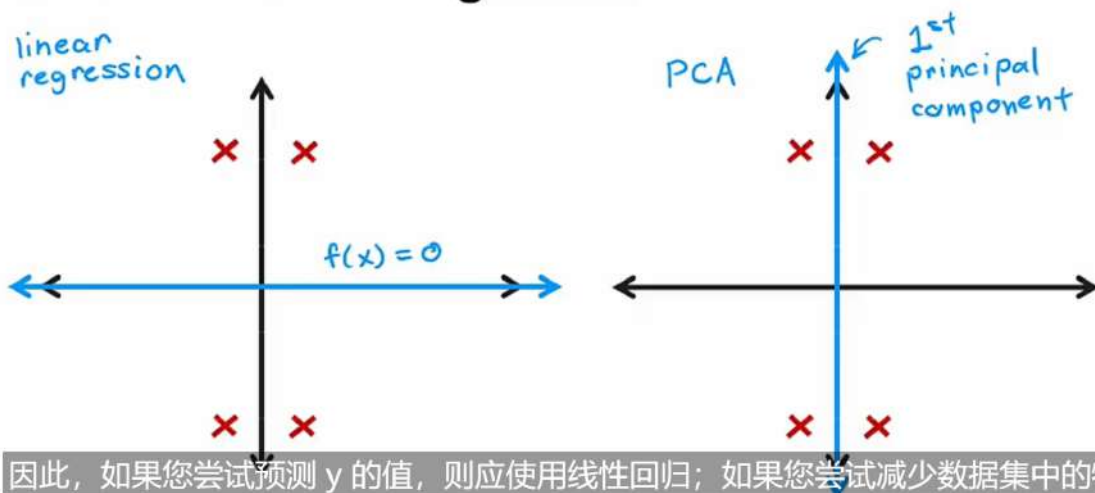


## PCA is not linear regression



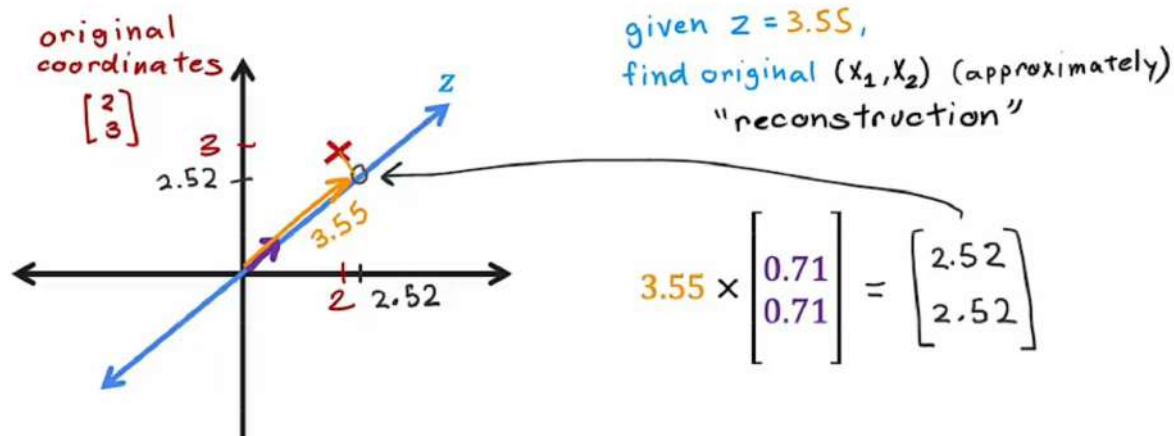
不一样

# PCA is not linear regression



因此，如果您尝试预测  $y$  的值，则应使用线性回归；如果您尝试减少数据集中的特征数量，例如将其可视化，则应使用 PCA。

## Approximation to the original data



## PCA In code

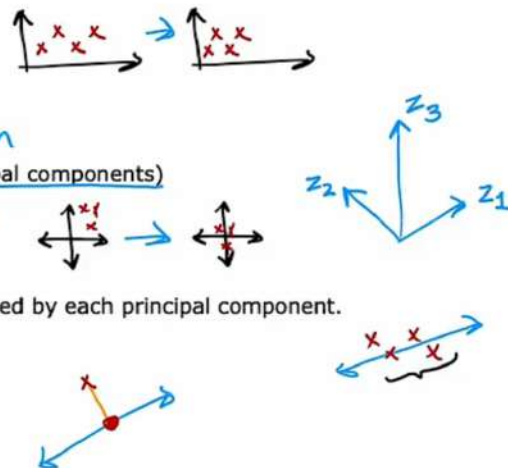
### PCA in scikit-learn

Optional pre-processing: Perform feature scaling

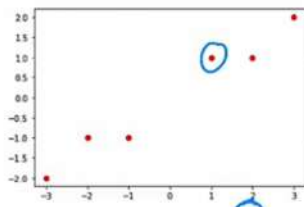
1. "fit" the data to obtain 2 (or 3) new axes (principal components)  
fit includes mean normalization

2. Optionally examine how much variance is explained by each principal component.  
info explained\_variance\_ratio

3. Transform (project) the data onto the new axes  
transform



## Example

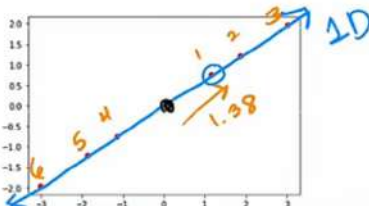


```
X = np.array([[1, 1], [2, 1], [3, 2],
               [-1, -1], [-2, -1], [-3, -2]])
```

2D

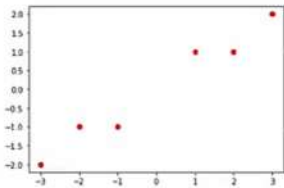
```
pca_1 = PCA(n_components=1)
pca_1.fit(X)
pca_1.explained_variance_ratio_ 0.992
X_trans_1 = pca_1.transform(X)
X_reduced_1 = pca_1.inverse_transform(X_trans_1)
```

```
array([
1 [ 1.38340578],
2 [ 2.22189802],
3 [ 3.6053038 ],
4 [-1.38340578],
5 [-2.22189802],
6 [-3.6053038 ]])
```



1D

## Example

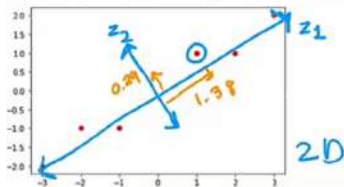


2D

```
X = np.array([[1, 1], [2, 1], [3, 2],
               [-1, -1], [-2, -1], [-3, -2]])
```

```
pca_2 = PCA(n_components=2)
pca_2.fit(X)
pca_2.explained_variance_ratio_ 0.992 0.008
X_trans_2 = pca_2.transform(X)
X_reduced_2 = pca_2.inverse_transform(X_trans_2)
```

```
array([
z1 z2
→ [ 1.38340578, 0.2935787 ],
[ 2.22189802, -0.25133484],
[ 3.6053038 , 0.04224385],
[-1.38340578, -0.2935787 ],
[-2.22189802, 0.25133484],
[-3.6053038 , -0.04224385]])
```



2D

## Applications of PCA

★ Visualization reduce to 2 or 3 features

Less frequently used for:

- Data compression (to reduce storage or transmission costs) 50 → 10
- Speeding up training of a supervised learning model

$n = 1000 \rightarrow 100$