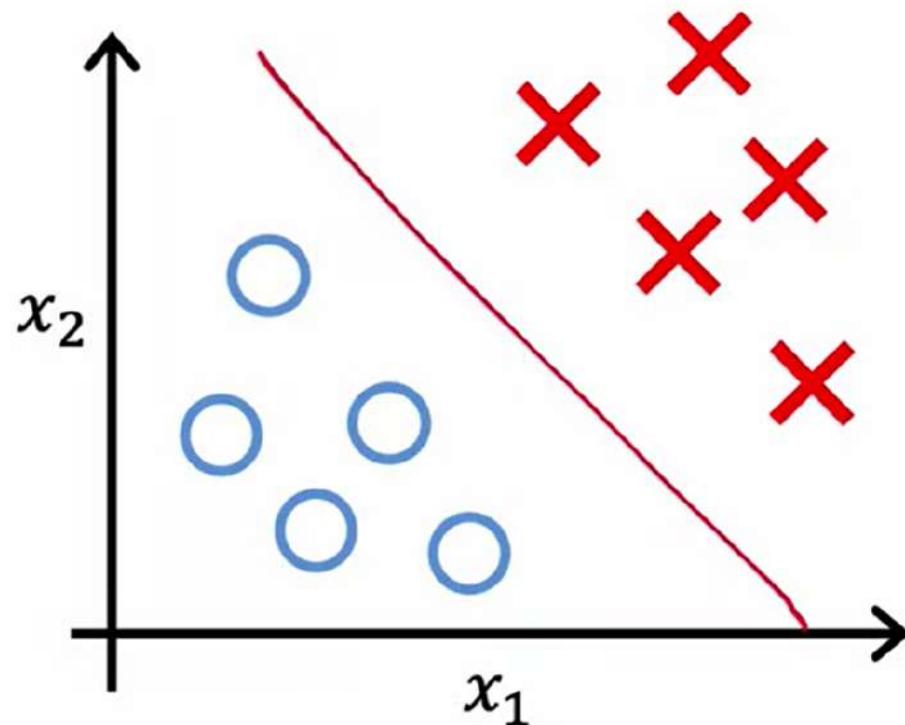


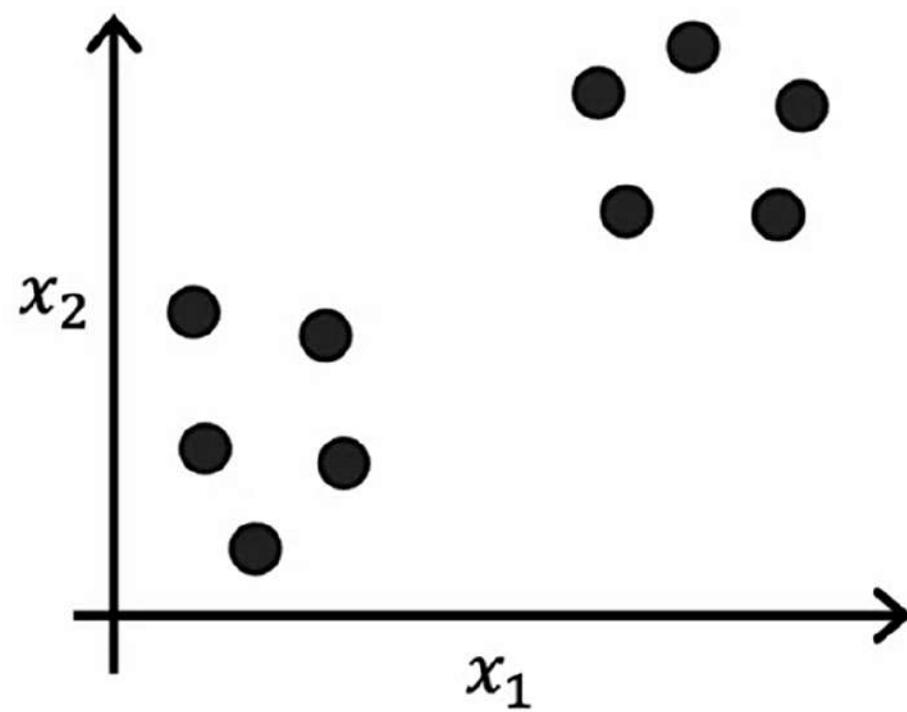
# Supervised learning



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

In contrast, in unsupervised learning,

# Unsupervised learning



Clustering

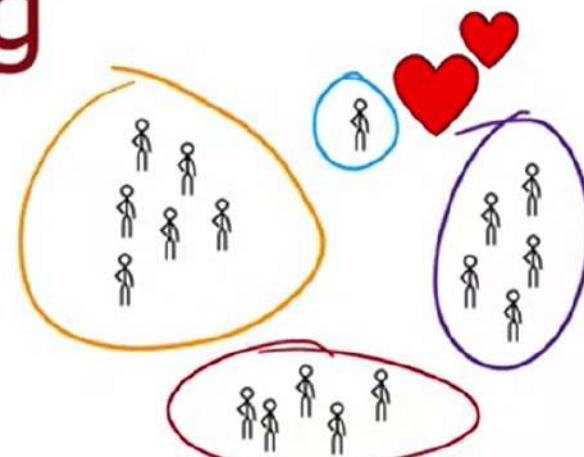
Training set:  $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(N)}\}$

# Applications of clustering

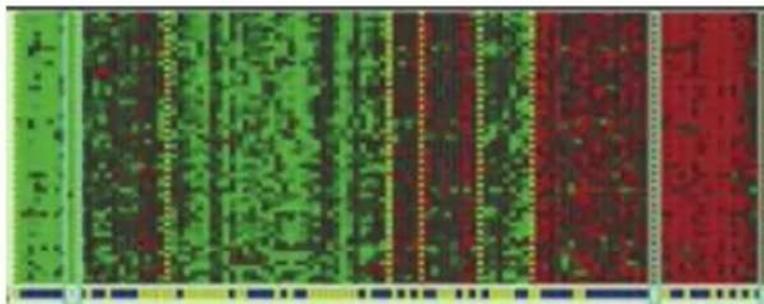


Grouping similar news

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



Market segmentation



DNA analysis

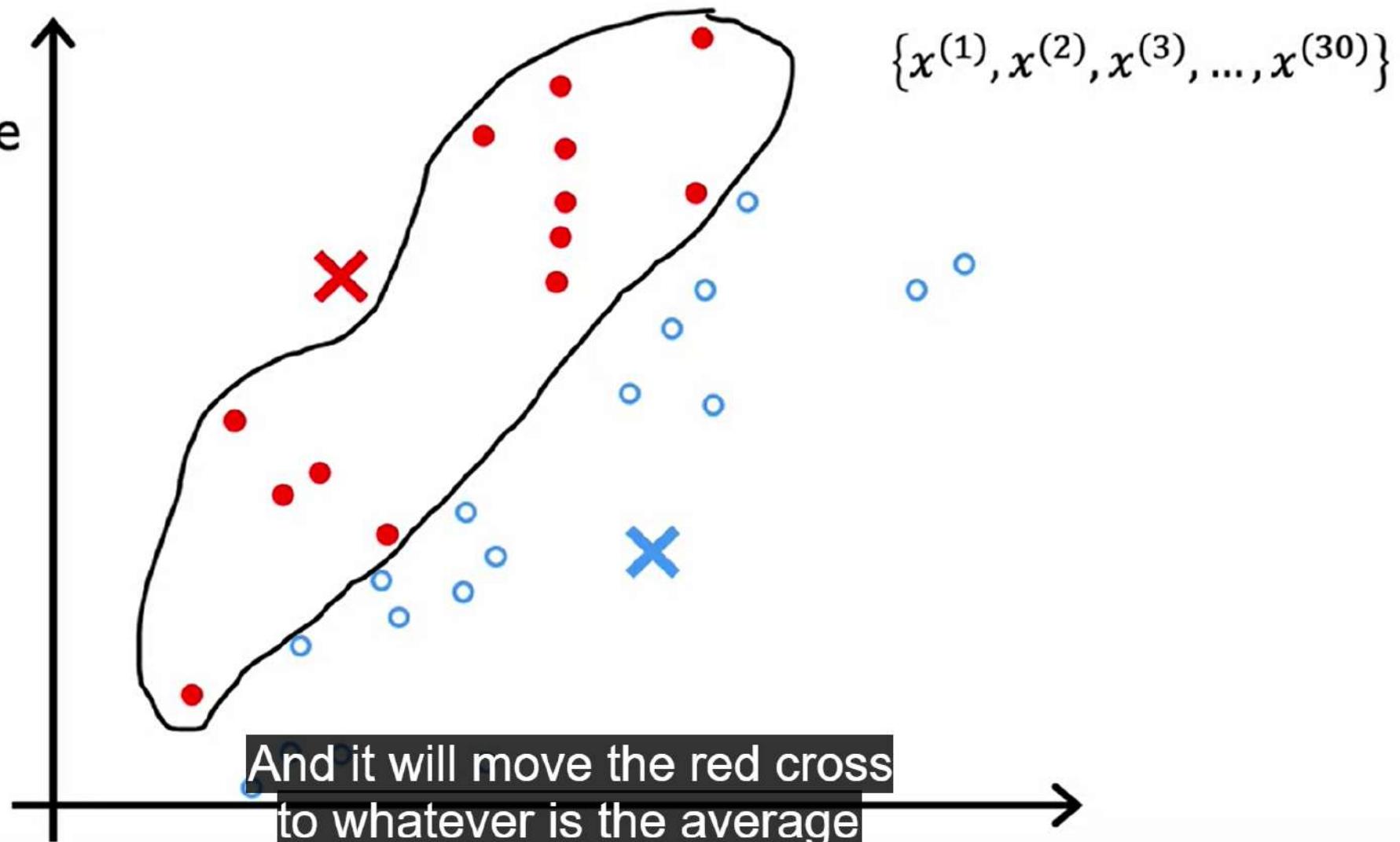
One of the applications I found



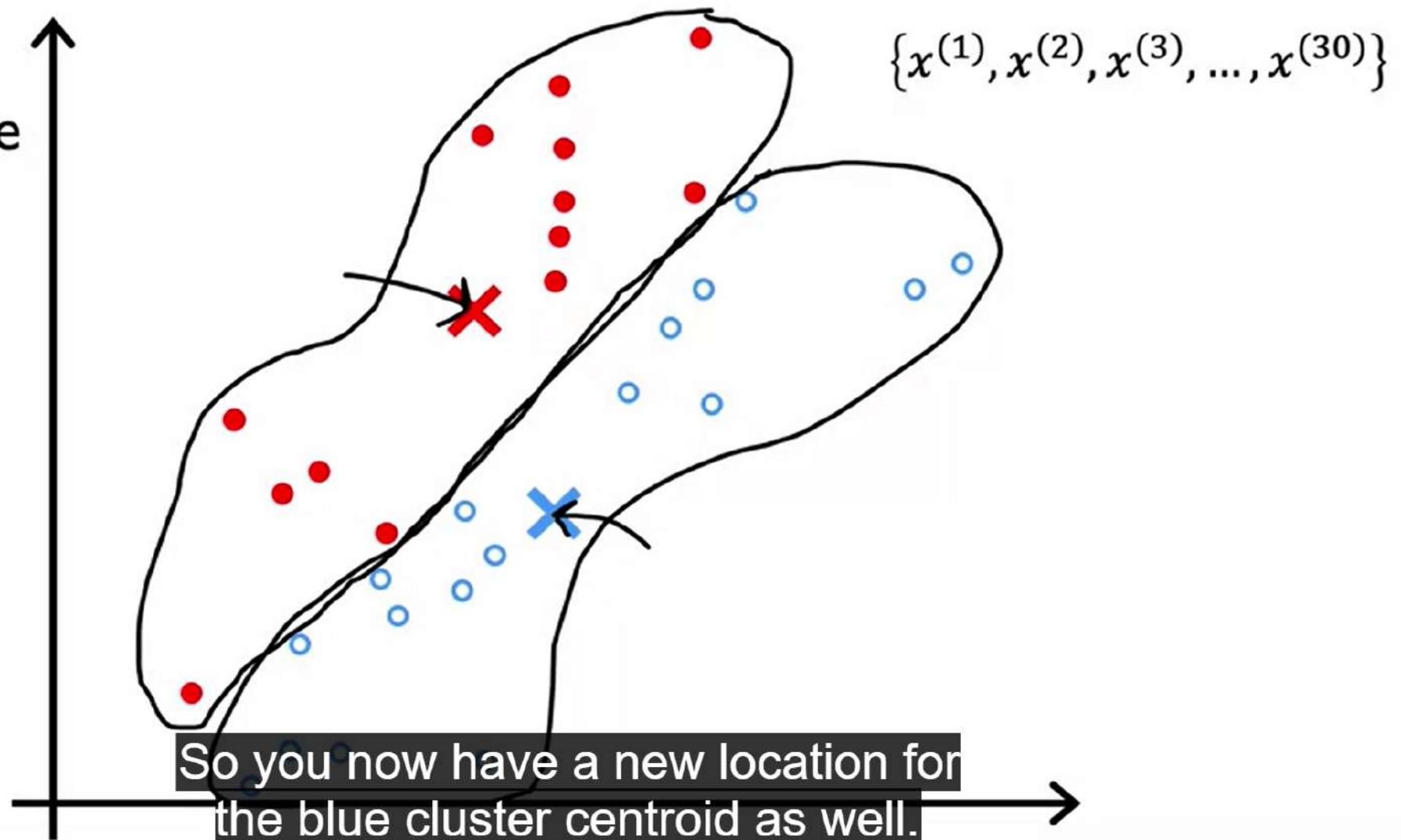
Image credit: NASA/JPL-Caltech/E. Churchwell (Univ. of Wisconsin, Madison)

Astronomical data analysis

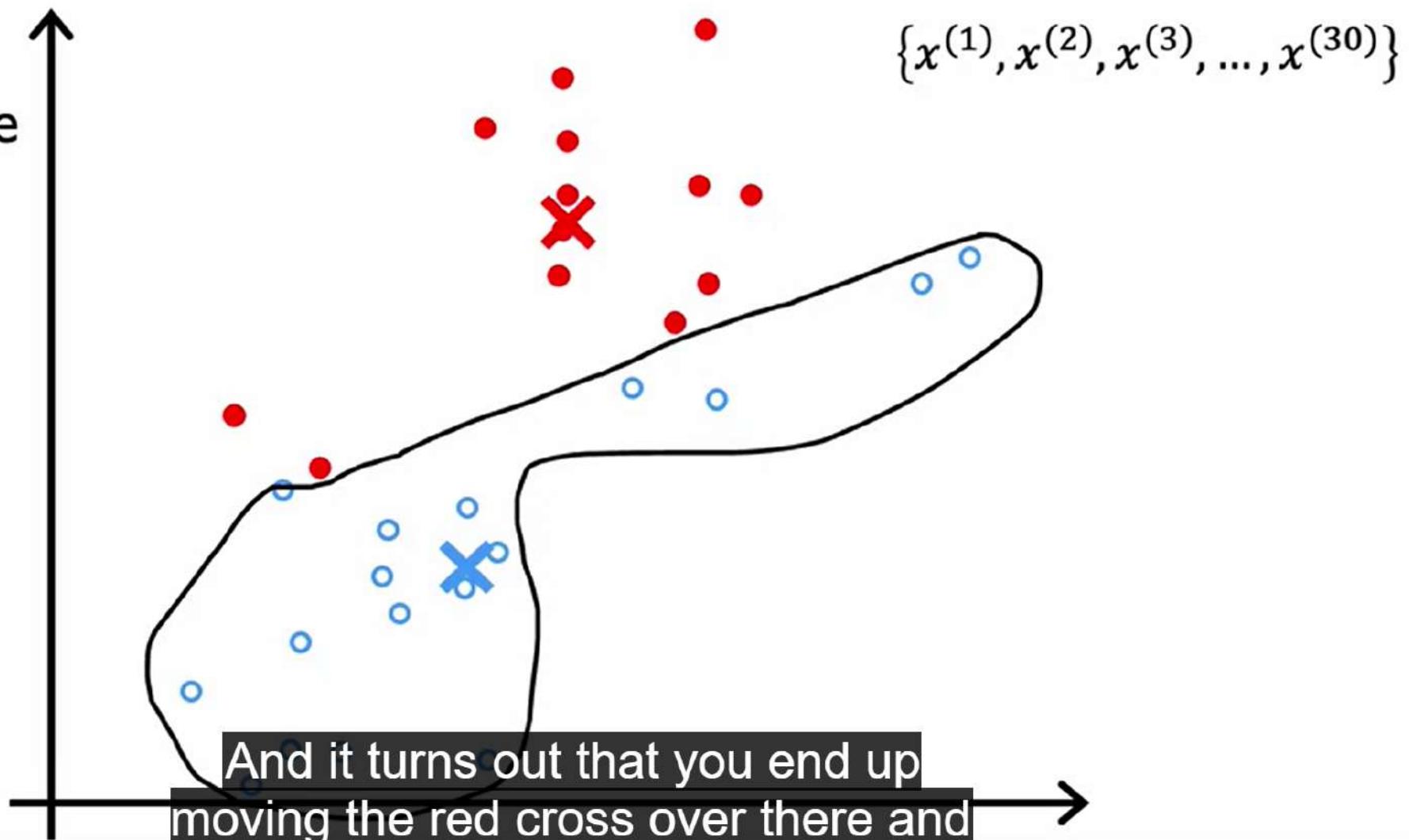
**Step 2:**  
Recompute  
the  
centroids



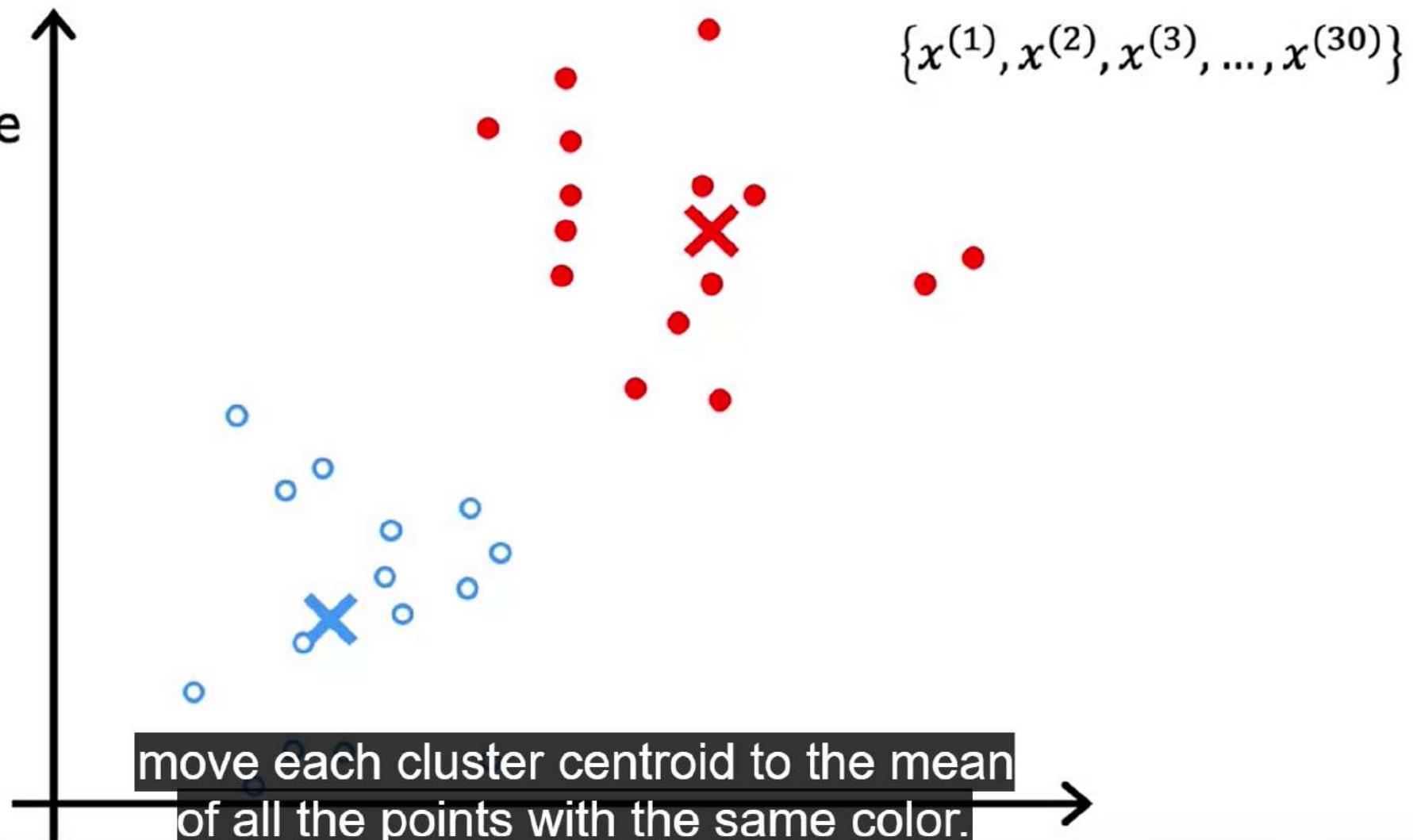
**Step 2:**  
Recompute  
the  
centroids



**Step 2:**  
Recompute  
the  
centroids



**Step 2:**  
Recompute  
the  
centroids



## 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  $k$

}

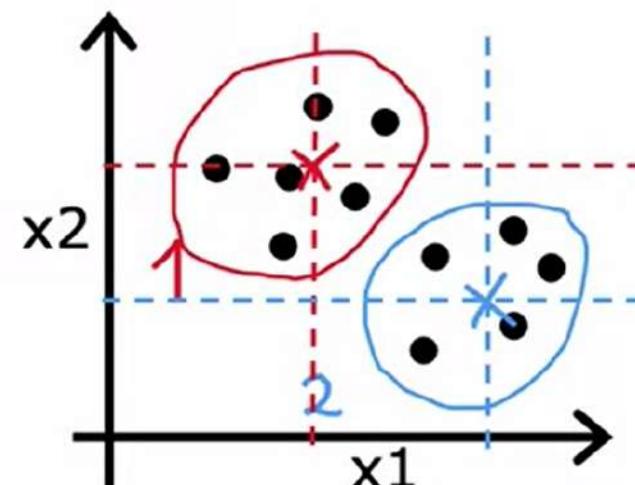
$$\mu_1 = \frac{1}{4} [x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}]$$

$\mu_1, \mu_2$

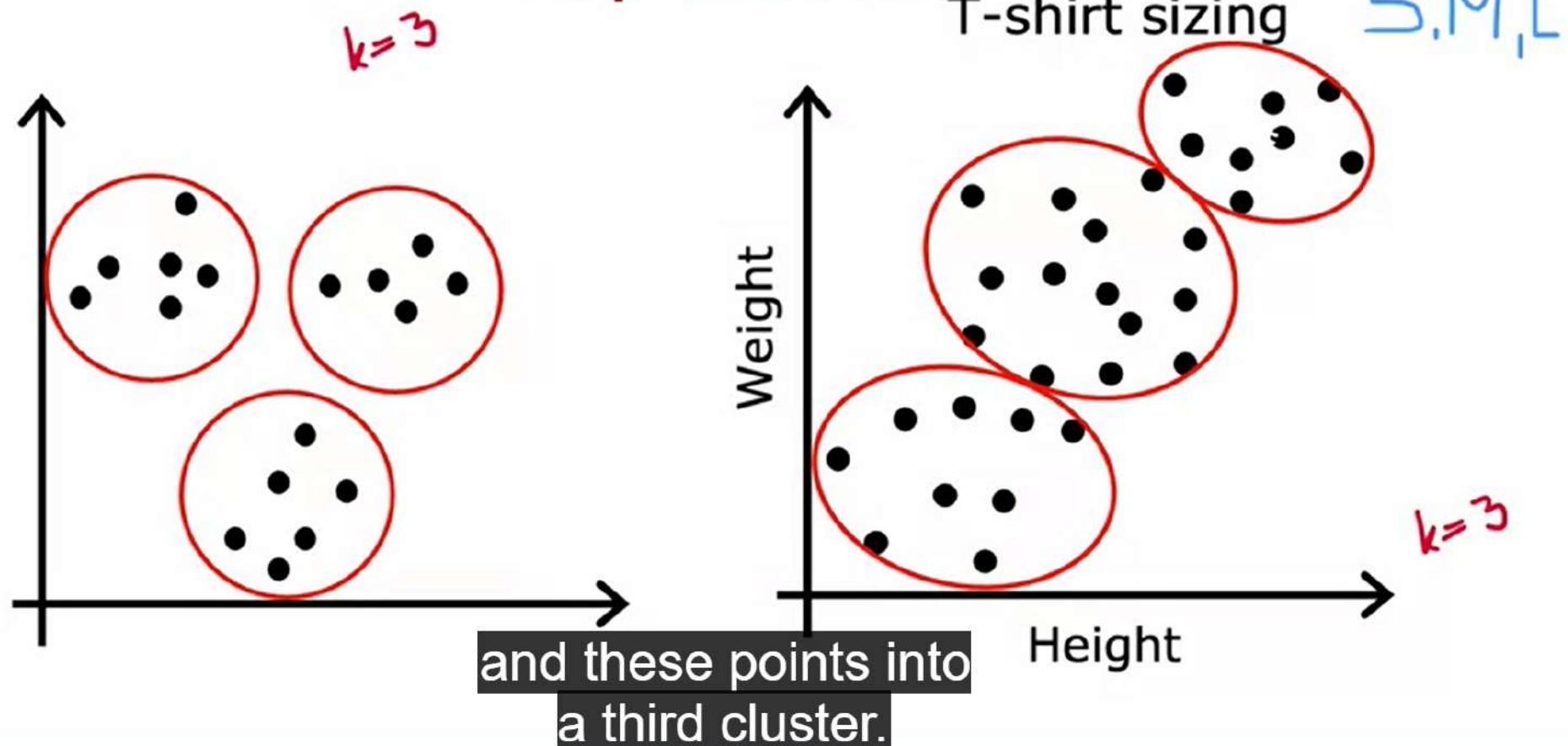
$x^{(1)}, x^{(2)}, \dots, x^{(10)}$

$n=2$

$K=2$



# K-means for clusters that are not well separated



## K-means optimization objective

$c^{(i)}$  = index of cluster ( $1, 2, \dots, K$ ) 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

$x^{(10)}$   $c^{(10)}$   $\mu_c^{(10)}$

### Cost function

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

$$\min_{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

These are the locations of all

the clusters centroid is defined as

## 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  $x^{(i)}$

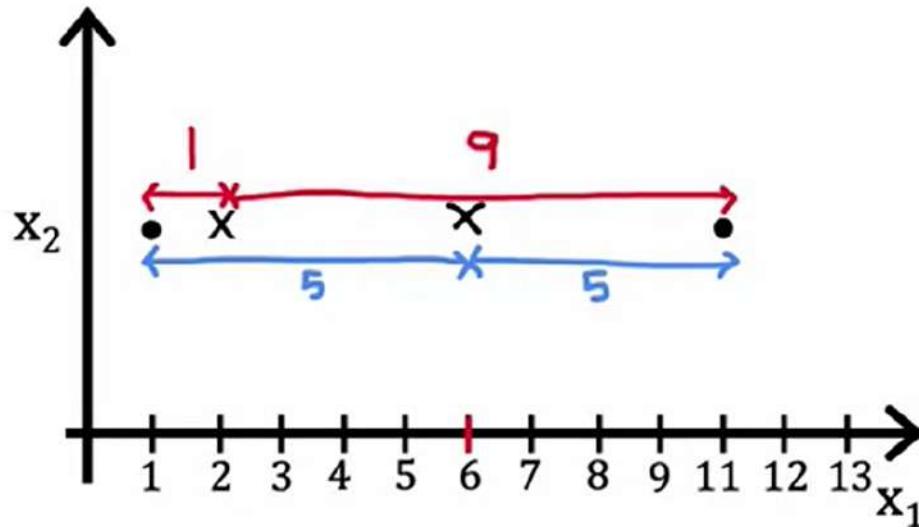
# Move cluster centroids

for  $k = 1$  to  $K$

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

} the distortion as much as possible.

# Moving the centroid



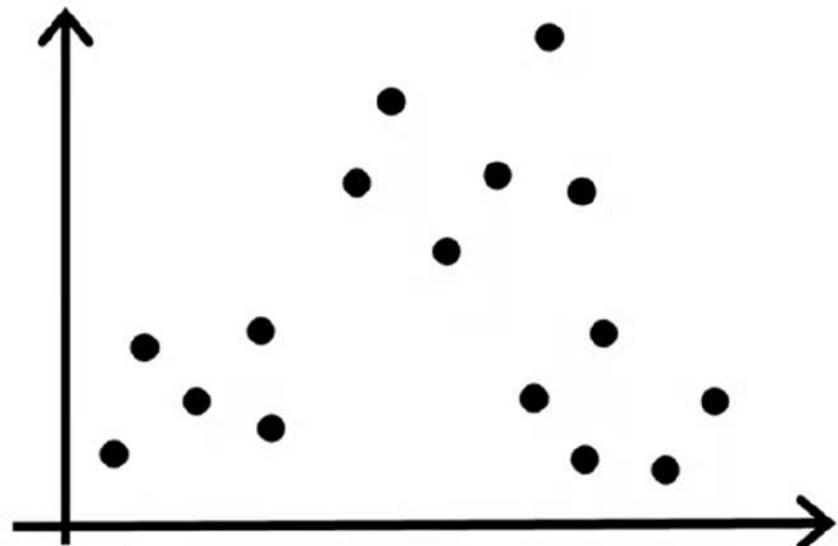
$$\frac{1}{2}(1^2 + 9^2) = \frac{1}{2}(1+81) = 41$$

$$\frac{1}{2}(1 + 11) = 6$$

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

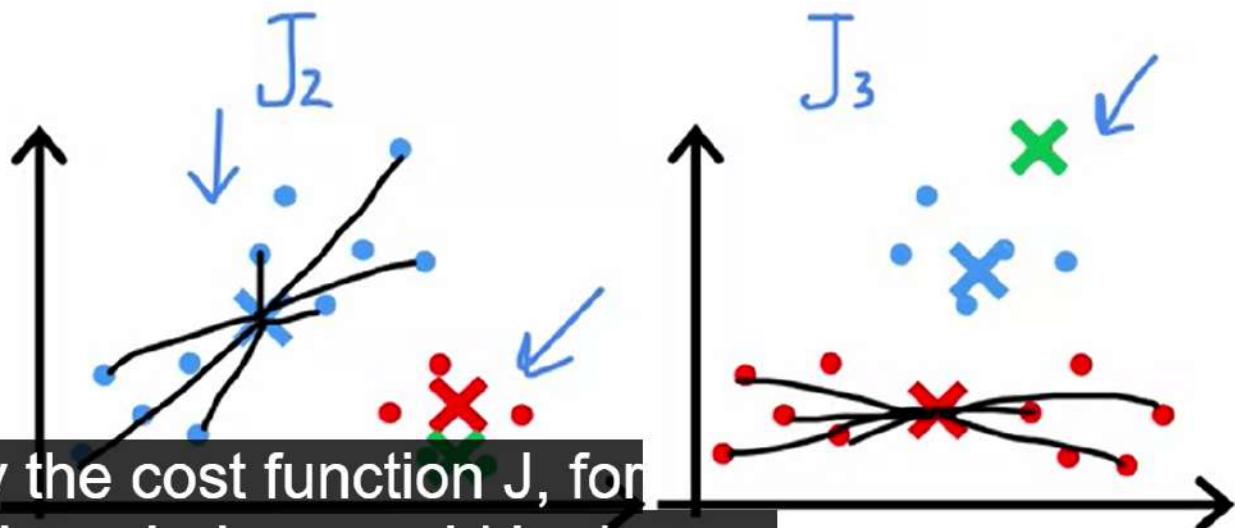
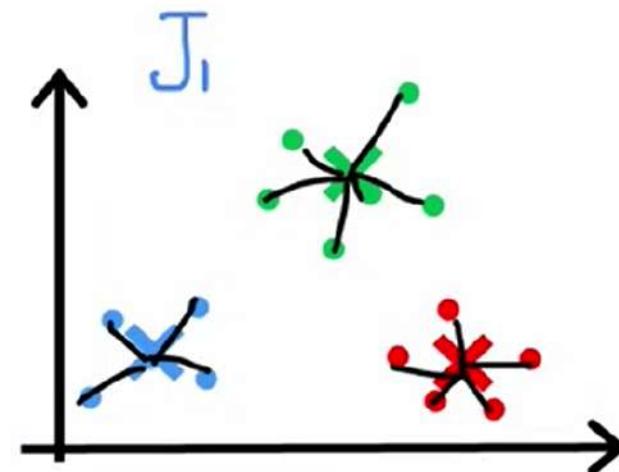
This average location in the middle  
of these two training examples,

$k=3$



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

which is why the cost function  $J$ , for these examples down below would be larger.



## Random initialization

For  $i = 1$  to 100 {

    Randomly initialize K-means.

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

    Computer cost function (distortion)

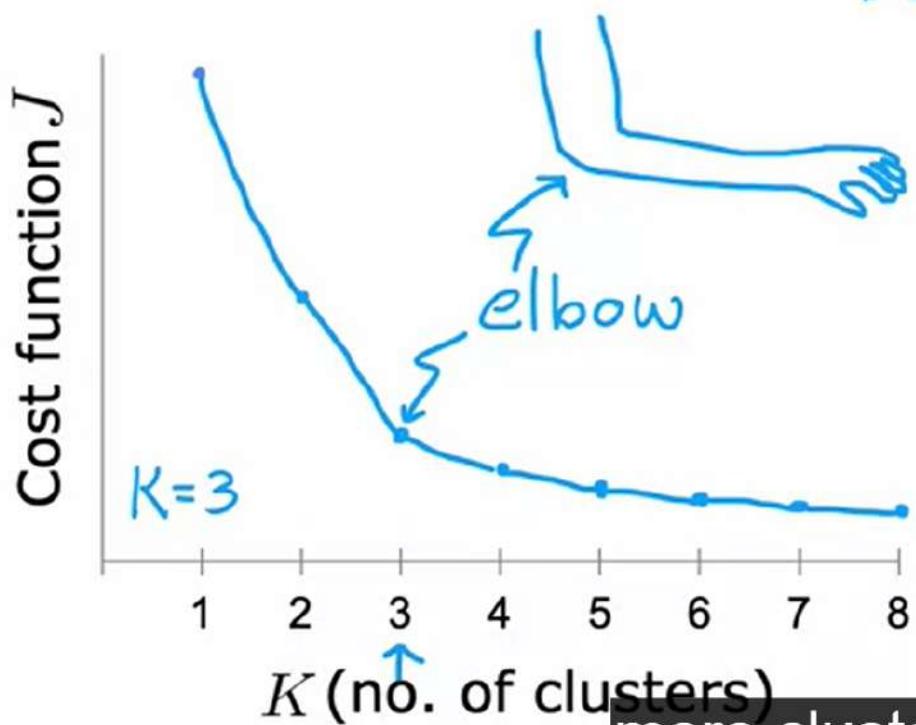
$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \mu_1, \dots, \mu_k) \leftarrow$

}

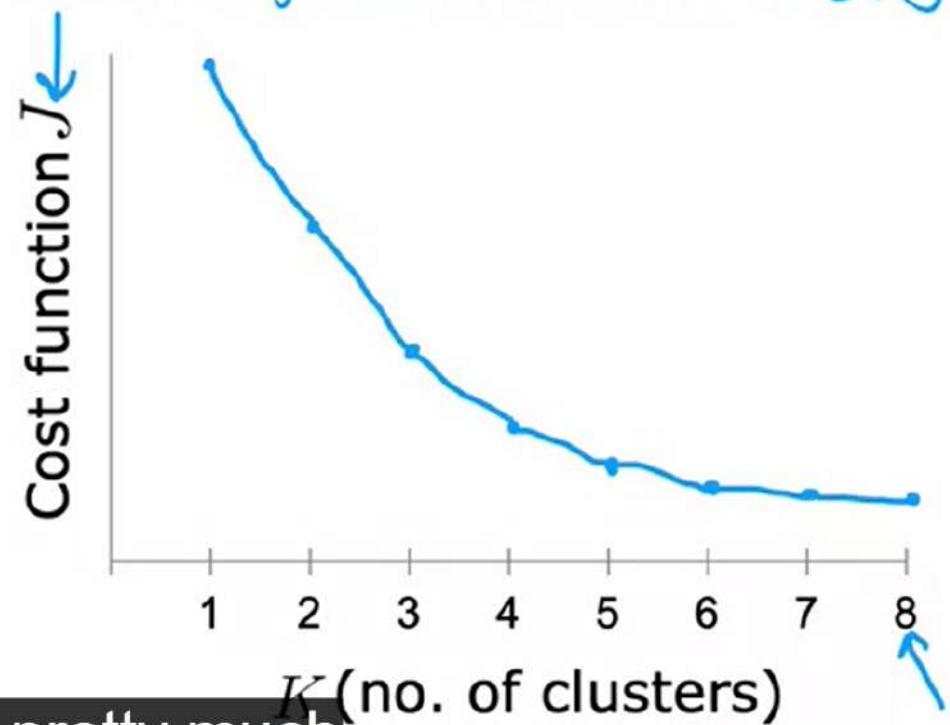
Pick set of clusters that gave lowest cost  $J$   
I plugged in the number up here as 100.

## Choosing the value of K

Elbow method



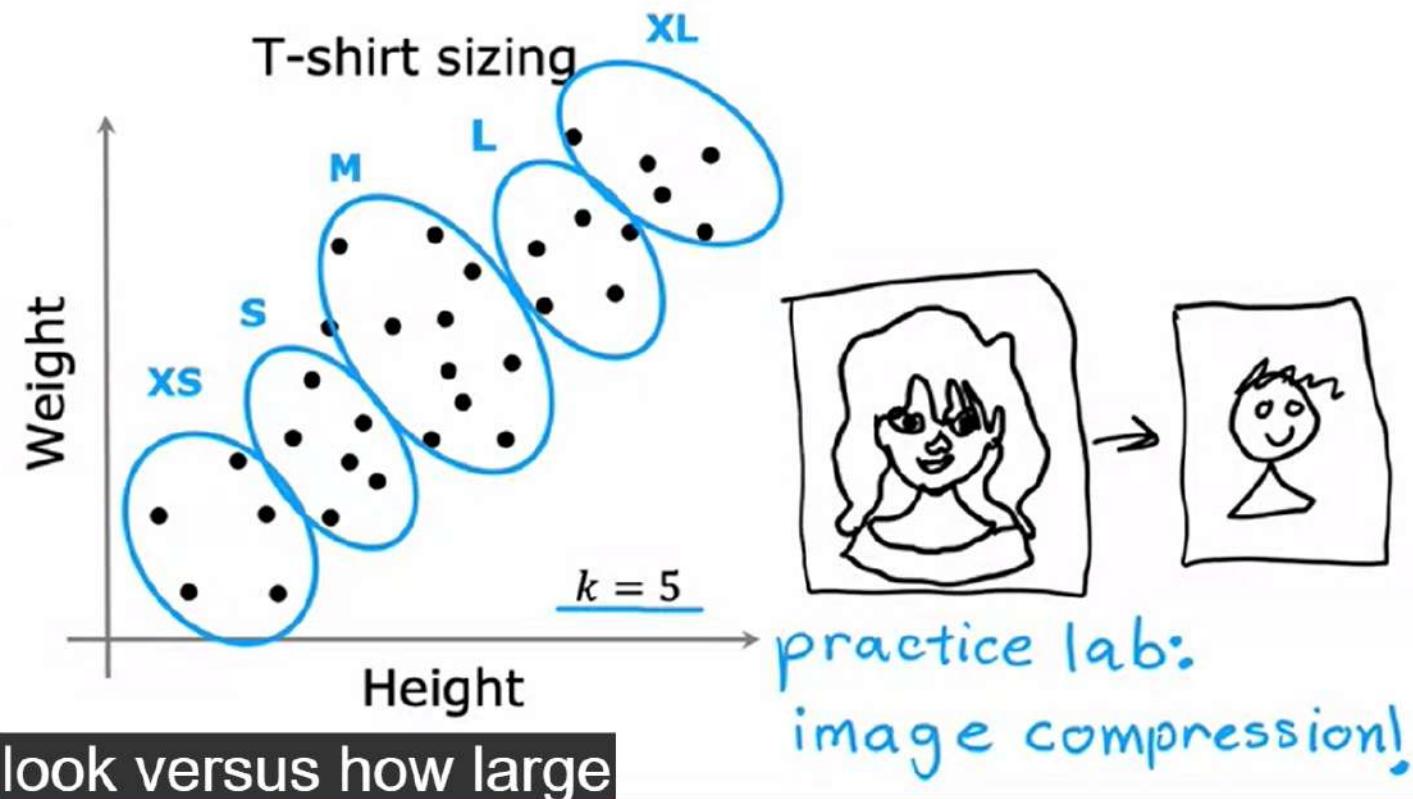
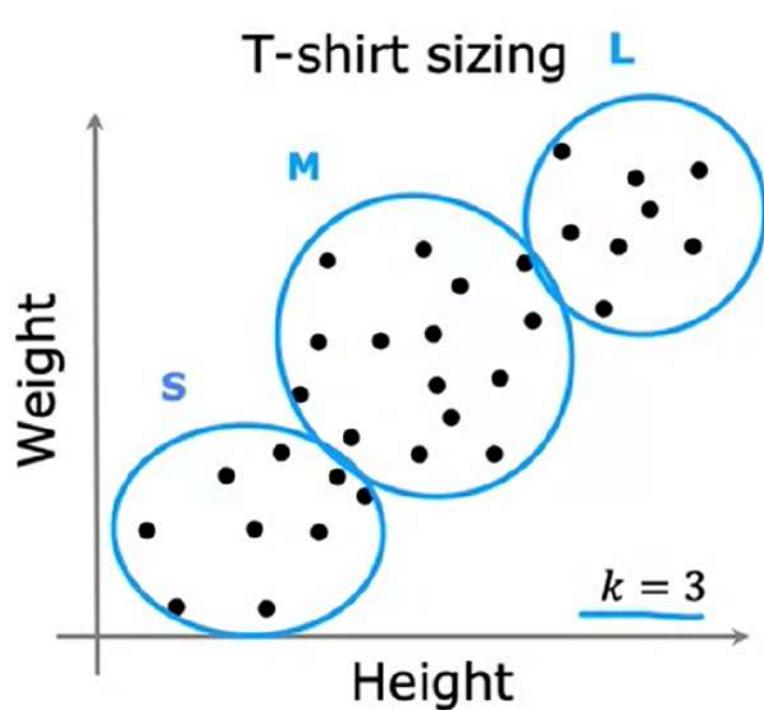
the right "K" is often ambiguous  
Don't choose  $K$  just to minimize cost  $J$



more clusters will pretty much

## Choosing the value of K

Often, you want to get clusters for some later (downstream) purpose.  
Evaluate K-means based on how well it performs on that later purpose.



## Anomaly detection example

Aircraft engine features:

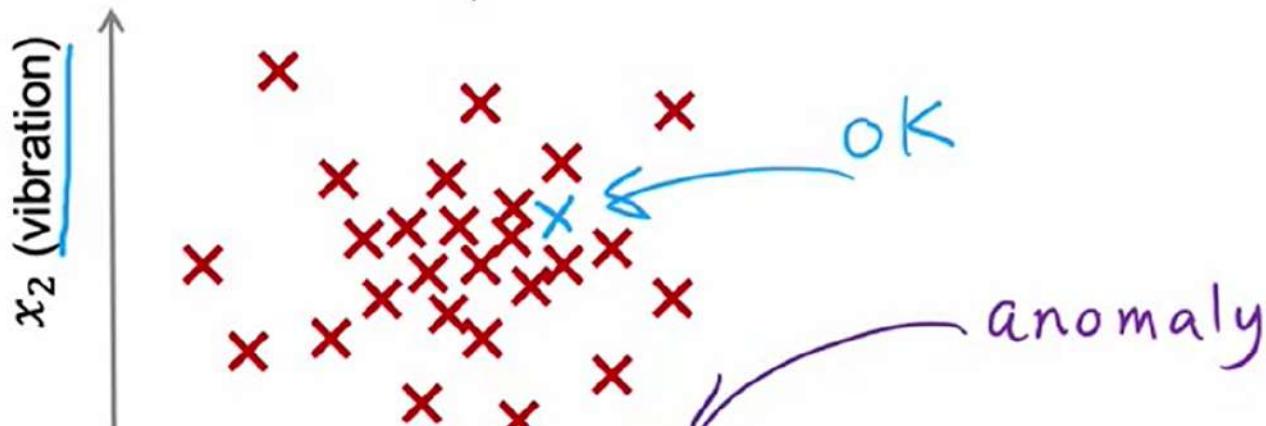
$x_1$  = heat generated

$x_2$  = vibration intensity

...

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

New engine:  $x_{test}$



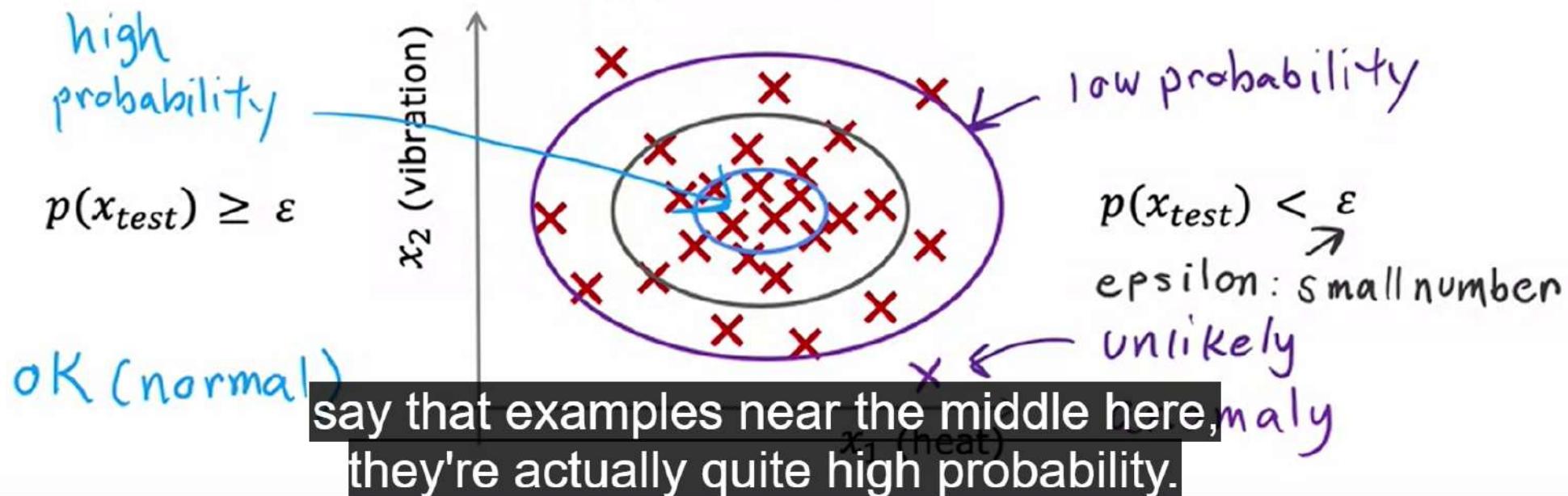
This doesn't look like the examples I've seen before, we better inspect this more

# Density estimation

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

probability of  $x$  being seen in dataset  
Model  $p(x)$

Is  $x_{test}$  anomalous?



## Anomaly detection example

Fraud detection:

- $x^{(i)}$  = features of user  $i$ 's activities
- Model  $p(x)$  from data.
- Identify unusual users by checking which have  $p(x) < \varepsilon$

how often log in?  
how many web pages visited?  
transactions?  
posts? typing speed?

perform additional checks to identify real fraud vs. false alarms

Manufacturing:

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

airplane engine

circuit board

smartphone

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.

it might be worth taking a look at that  
computer to see if something is wrong

# Gaussian (Normal) distribution

Say  $x$  is a number.

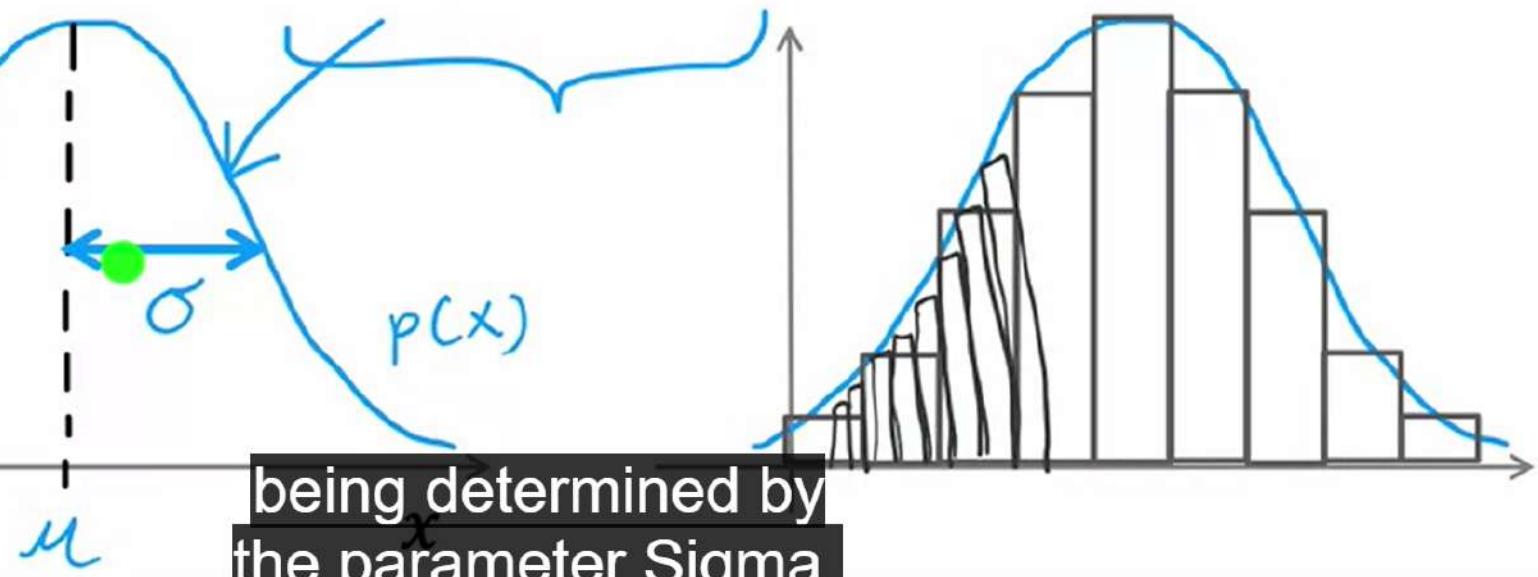
Probability of  $x$  is determined by a Gaussian with mean  $\mu$ , variance  $\sigma^2$ .

$\sigma$  standard deviation  
 $\sigma^2$  variance



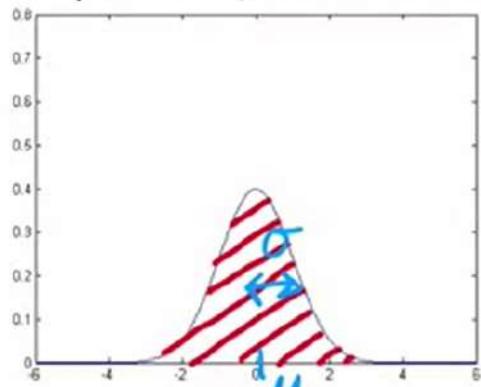
$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\pi = 3.14$$

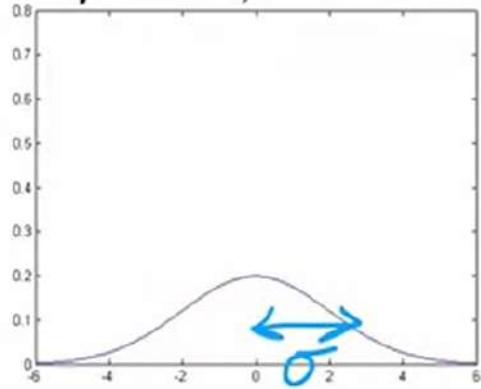


## Gaussian distribution example

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

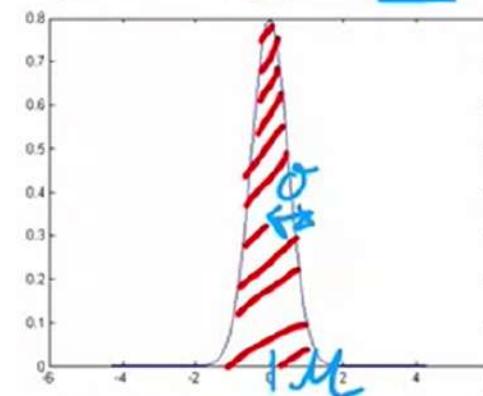


$$\mu = 0, \sigma = 2$$



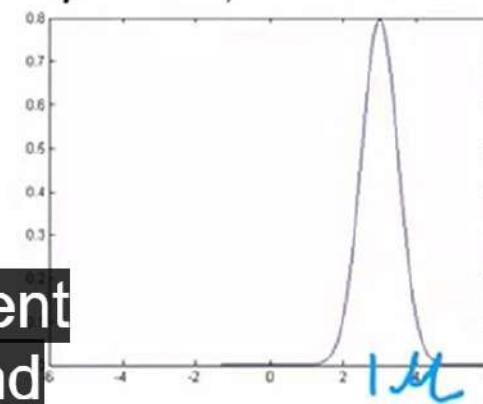
$$\sigma^2 = 4$$

$$\mu = 0, \underline{\sigma} = 0.5$$



$$\sigma^2 = 0.25$$

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

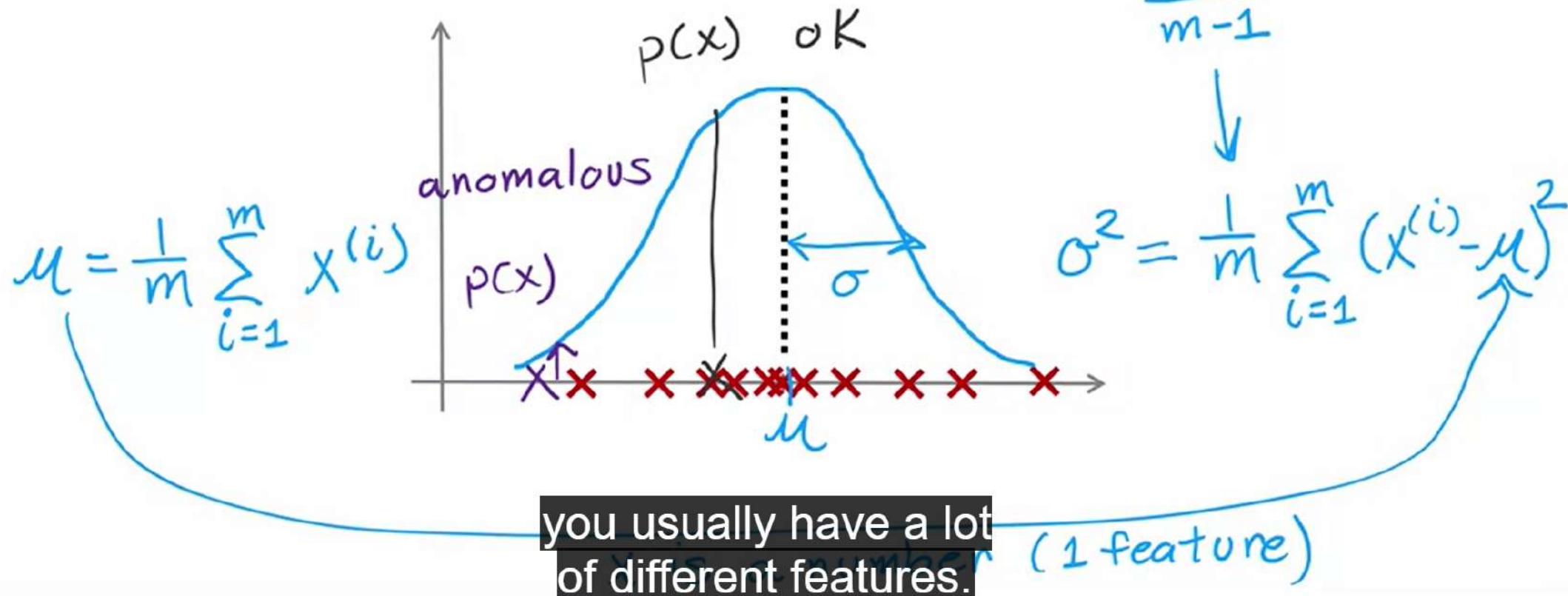


This is how different choices of Mu and

## Parameter estimation

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

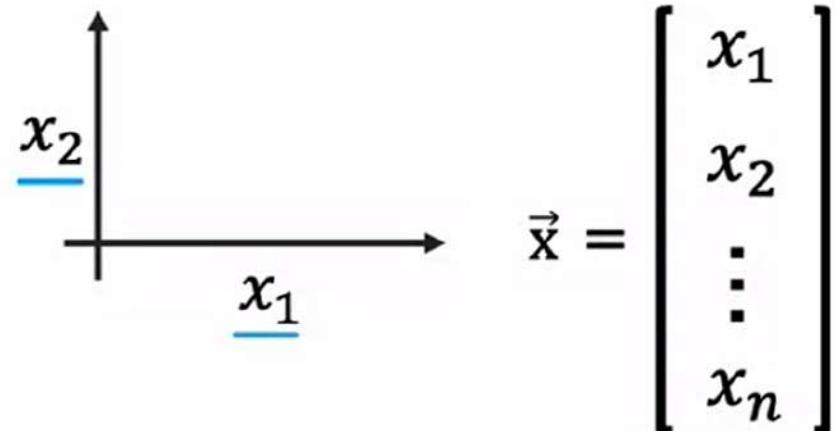
maximum likelihood  
for  $\mu, \sigma$



## Density estimation

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

Each example  $\vec{x}^{(i)}$  has  $n$  features



$$p(\vec{x}) = p(x_1; \mu_1, \sigma_1^2) * p(x_2; \mu_2, \sigma_2^2) * p(x_3; \mu_3, \sigma_3^2) * \dots * p(x_n; \mu_n, \sigma_n^2)$$

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

"add"      "multiply"

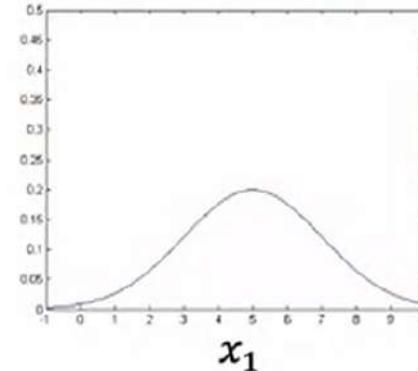
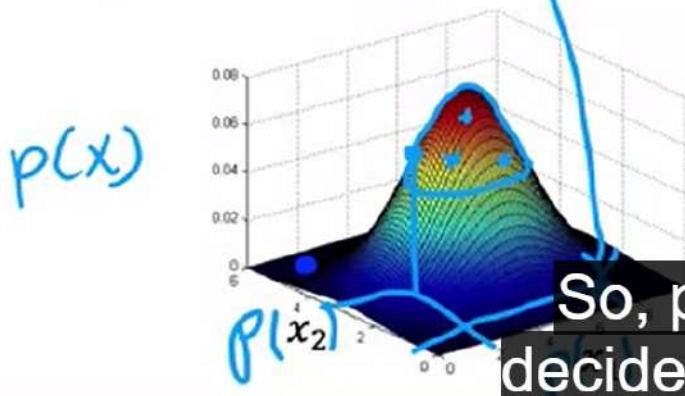
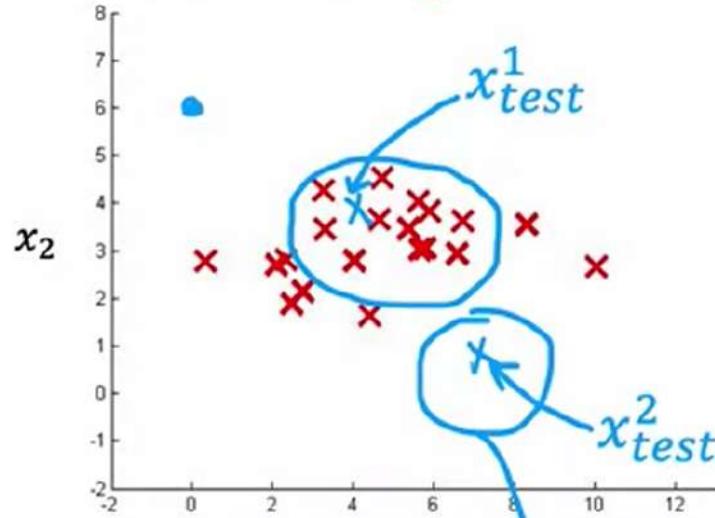
$$p(x_1 = \text{high temp}) = 1/10$$

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

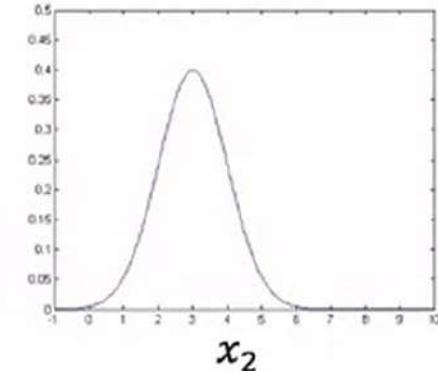
$$p(x_1, x_2) = p(x_1) * p(x_2)$$

symbol here corresponds to multiplying these terms over here for  $j = 1$  through  $n$ .  $\frac{1}{10} * \frac{1}{20} = \frac{1}{200}$

# Anomaly detection example



$$\mu_1 = 5, \sigma_1 = 2$$
$$p(x_1; \mu_1, \sigma_1^2)$$



$$\mu_2 = 3, \sigma_2 = 1$$
$$p(x_2; \mu_2, \sigma_2^2)$$

$$\varepsilon = 0.02$$

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

"ok" anomaly

So, pretty much as you might hope it decides that  $x$  test 1 looks pretty normal.

## Aircraft engines monitoring example

10000 good (normal) engines  
20 flawed engines (anomalous)

2 to 50

$y=0$

Training set: 6000 good engines

train algorithm on training set

CV: 2000 good engines ( $y = 0$ )  
use cross validation set

10 anomalous ( $y = 1$ )

Test: 2000 good engines ( $y = 0$ ),

tune  $\epsilon$  tune  $x_j$

10 anomalous ( $y = 1$ )

Alternative: No test set

Training set: 6000 good engines 2

CV: 4000 good engines ( $y = 0$ ) put all of that in the 1)  
cross validation set. 4j

## Anomaly detection vs. Supervised learning

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

Large number of negative ( $y = 0$ ) examples.

P(x)

y=1

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

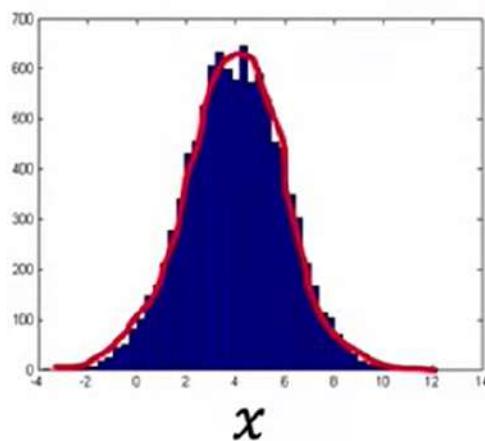
Large number of positive and negative examples.

**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.

some individuals are trying to commit financial fraud.

# Non-gaussian features



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

`plt.hist(x)`

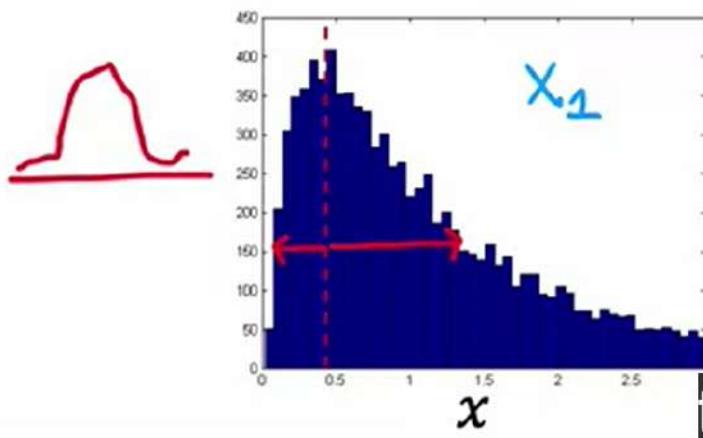
$$x_1 \leftarrow \log(x_1)$$

$$x_2 \leftarrow \log(x_2 + 1)$$

$$x_3 \leftarrow \sqrt{x_3} = x_3^{1/2}$$

$$x_4 \leftarrow x_4^{1/3}$$

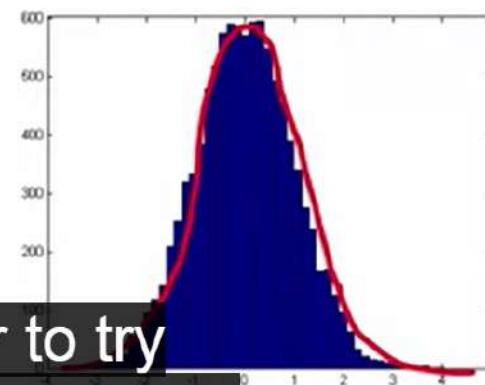
$$\log(x_2 + c)$$



`np.log(x)`



others, In order to try  
to make it more Gaussian.  $x$

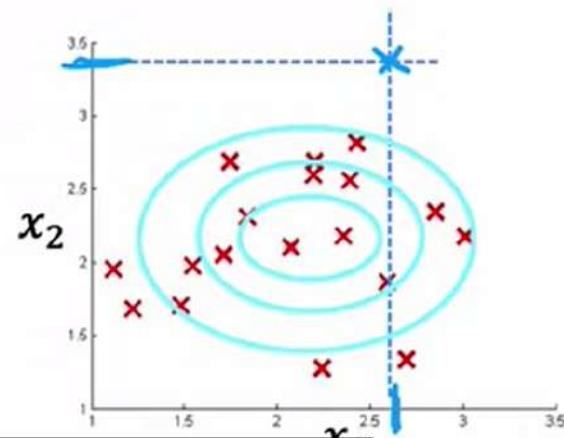
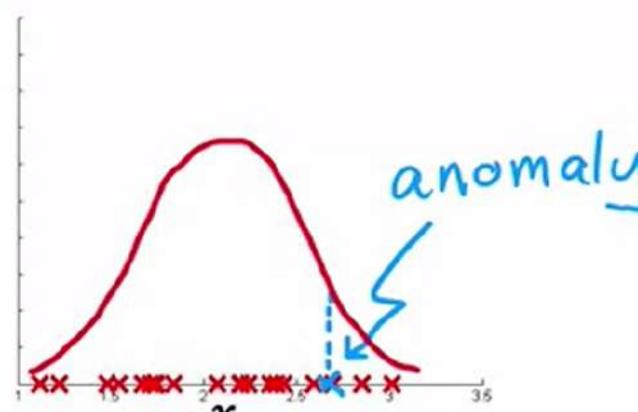


# Error analysis for anomaly detection

Want  $p(x) \geq \varepsilon$  large for normal examples  $x$ .  
 $p(x) < \varepsilon$  small for anomalous examples  $x$ .

Most common problem:

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



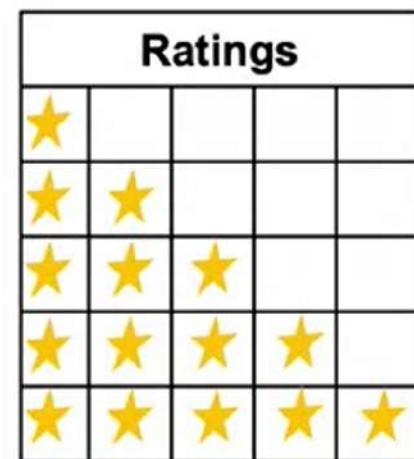
$x_1$  num transactions  
the creation of new features that would allow the algorithm to spot.

## Predicting movie ratings

User rates movies using one to five stars

*zero*

Movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)
Love at last	5	5	0	0
Romance forever	5	?	?	0
Cute puppies of love	?	4	0	?
Nonstop car chases	0	0	5	4
Swords vs. karate	0	0	5	?



Ratings

$n_u = \text{no. of users}$

$n_m = \text{no. of movies}$

$r(i,j) = 1$  if user  $j$  has rated movie  $i$

$$n_u = 4$$

$$r(1,1) = 1$$

$y^{(i,j)}$  = rating given by user  $j$  to movie  $i$

$$n_m = 5$$

$$r(3,1) = 0 \quad r^{(3,2)} = 4$$

(defined only if  $r(i,j)=1$ )

$r(i,j)=1$  if user  $j$  has rated movie  $i$  and

## What if we have features of the movies?

$$n_u = 4$$

$$n_m = 5$$

$$n = 2$$

Movie	Alice(1)	Bob(2)	Carol(3)	Dave(4)	$x_1$ (romance)	$x_2$ (action)	
Love at last	5	5	0	0	0.9	0	
Romance forever	5	?	?	0	1.0	0.01	$x^{(1)} = \begin{bmatrix} 0.9 \\ 0 \end{bmatrix}$
Cute puppies of love	?	4	0	?	0.99	0	
Nonstop car chases	0	0	5	4	0.1	1.0	$x^{(3)} = \begin{bmatrix} 0.99 \\ 0 \end{bmatrix}$
Swords vs. karate	0	0	5	?	0	0.9	

For user 1: Predict rating for movie  $i$  as:  $w^{(1)} \cdot x^{(i)} + b^{(1)}$  ← just linear regression

$$w^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix} \quad b^{(1)} = 0 \quad x^{(3)} = \begin{bmatrix} 0.99 \\ 0 \end{bmatrix}$$

$$w^{(1)} \cdot x^{(3)} + b^{(1)} = 4.95$$

→ For user  $j$ : Predict user  $j$ 's rating for movie  $i$  as  $w^{(j)} \cdot x^{(i)} + b^{(j)}$   
movie  $i$  as  $w(j) \cdot X(i) + b(j)$ .

# Cost function

Notation:

- $r(i,j) = 1$  if user  $j$  has rated movie  $i$  (0 otherwise)
- $y^{(i,j)}$  = rating given by user  $j$  on movie  $i$  (if defined)
- $w^{(j)}, b^{(j)}$  = parameters for user  $j$
- $x^{(i)}$  = feature vector for movie  $i$

For user  $j$  and movie  $i$ , predict rating:  $w^{(j)} \cdot x^{(i)} + b^{(j)}$

- $m^{(j)}$  = no. of movies rated by user  $j$

To learn  $w^{(j)}, b^{(j)}$

$$\min_{w^{(j)}, b^{(j)}} J(w^{(j)}, b^{(j)}) = \frac{1}{2m^{(j)}} \sum_{r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^n (w_k^{(j)})^2$$

number of features

you should end up with  
the same value of w and b.

# Cost function

To learn parameters  $\underline{w^{(j)}, b^{(j)}}$  for user  $j$  :

$$J(w^{(j)}, b^{(j)}) = \frac{1}{2} \sum_{i:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (w_k^{(j)})^2$$

To learn parameters  $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots w^{(n_u)}, b^{(n_u)}$  for all users :

$$J\begin{pmatrix} w^{(1)}, & \dots, & w^{(n_u)} \\ b^{(1)}, & \dots, & b^{(n_u)} \end{pmatrix} = \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

where that plays a role similar to  
the output  $f(x)$  of linear regression.

# Problem motivation

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$x_1$ (romance)	$x_2$ (action)
Love at last	5	5	0	0	?	?
Romance forever	5	?	?	0	?	?
Cute puppies of love	?	4	0	?	?	?
Nonstop car chases	0	0	5	4	?	?
Swords vs. karate	0	0	5	?	?	?

$$w^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}, w^{(2)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}, w^{(3)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}, w^{(4)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$$

$b^{(1)} = 0, b^{(2)} = 0, b^{(3)} = 0, b^{(4)} = 0$

using  $w^{(j)} \cdot x^{(i)} + b^{(j)}$

**That's what makes it possible to try to**

$$\rightarrow x^{(1)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

# Cost function

Given  $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(n_u)}, b^{(n_u)}$

to learn  $x^{(i)}$ :

$$J(x^{(i)}) = \frac{1}{2} \sum_{j:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2$$

→ To learn  $x^{(1)}, x^{(2)}, \dots, x^{(n_m)}$ :

$$J(x^{(1)}, x^{(2)}, \dots, x^{(n_m)}) = \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

of  $x_1$  through  $x^{n_m}$

# Collaborative filtering

Cost function to learn  $w^{(1)}, b^{(1)}, \dots, w^{(n_u)}, b^{(n_u)}$ :

$$\min_{w^{(1)}, b^{(1)}, \dots, w^{(n_u)}, b^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

Cost function to learn  $x^{(1)}, \dots, x^{(n_m)}$ :

$$\min_{x^{(1)}, \dots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

Put them together:

$$\frac{1}{2} \sum_{(i,j):r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

where we do have a rating

# Collaborative filtering

Cost function to learn  $w^{(1)}, b^{(1)}, \dots, w^{(n_u)}, b^{(n_u)}$ :

$$\min_{w^{(1)}, b^{(1)}, \dots, w^{(n_u)}, b^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

	$j=1$	$j=2$	$j=3$
	Alice	Bob	Carol
Movie1	5	5	?
Movie2	?	2	3

Cost function to learn  $x^{(1)}, \dots, x^{(n_m)}$ :

$$\min_{x^{(1)}, \dots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

Put them together:

$$\min_{\substack{w^{(1)}, \dots, w^{(n_u)} \\ b^{(1)}, \dots, b^{(n_u)} \\ x^{(1)}, \dots, x^{(n_m)}}} J(w, b, x) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

How do you minimize this cost function

# Gradient Descent

collaborative filtering

Linear regression (course 1)

repeat {

$$\cancel{w_i = w_i - \alpha \frac{\partial}{\partial w_i} J(w, b)}$$

$$\cancel{b = b - \alpha \frac{\partial}{\partial b} J(w, b)}$$

$$w_i^{(j)} = w_i^{(j)} - \alpha \frac{\partial}{\partial w_i^{(j)}} J(w, b, x)$$

$$b^{(j)} = b^{(j)} - \alpha \frac{\partial}{\partial b^{(j)}} J(w, b, x)$$

$$x_k^{(i)} = x_k^{(i)} - \alpha \frac{\partial}{\partial x_k^{(i)}} J(w, b, x)$$

}

parameters  $w, b, x$

$x$  is also a parameter  
w and b, as well as x.

## Example applications

- 1. Did user  $j$  purchase an item after being shown? 1, 0, ?
- 2. Did user  $j$  fav/like an item? 1, 0, ?
- 3. Did user  $j$  spend at least 30sec with an item? 1, 0, ?
- 4. Did user  $j$  click on an item? 1, 0, ?

Meaning of ratings:

- 1 - engaged after being shown item
  - 0 - did not engage after being shown item
  - ? - item not yet shown
- A labor of one means that the user engaged after being shown an item And

# From regression to binary classification

- Previously:
- Predict  $y^{(i,j)}$  as  $\underbrace{\mathbf{w}^{(j)} \cdot \mathbf{x}^{(i)} + \mathbf{b}^{(j)}}$
- For binary labels:  
Predict that the probability of  $y^{(i,j)} = 1$   
is given by  $\underbrace{g(\mathbf{w}^{(j)} \cdot \mathbf{x}^{(i)} + \mathbf{b}^{(j)})}$

where  $g(z) = \underbrace{\frac{1}{1+e^{-z}}}$

turn it into something that would  
be a lot like a logistic regression

# Cost function for binary application

Previous cost function:

$$\frac{1}{2} \sum_{(i,j):r(i,j)=1} \underbrace{(w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2}_{f(x)} + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (w_k^{(j)})^2$$

Loss for binary labels  $y^{(i,j)}$ :  $f_{(w,b,x)}(x) = g(w^{(j)} \cdot x^{(i)} + b^{(j)})$

$$L(f_{(w,b,x)}(x), y^{(i,j)}) = -y^{(i,j)} \log(f_{(w,b,x)}(x)) - (1 - y^{(i,j)}) \log(1 - f_{(w,b,x)}(x))$$

Loss for  
single  
example

$$J(w, b, x) = \sum_{(i,j):r(i,j)=1} L(f_{(w,b,x)}(x), y^{(i,j)})$$

we're going to use that loss function.

## Users who have not rated any movies

Movie	Alice(1)	Bob (2)	Carol (3)	Dave (4)	Eve (5)	
Love at last	5	5	0	0	?	○
Romance forever	5	?	?	0	?	○
Cute puppies of love	?	4	0	?	?	○
Nonstop car chases	0	0	5	4	?	○
Swords vs. karate	0	0	5	?	?	○

$$\begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}$$



→  $\min_{\substack{w^{(1)}, \dots, w^{(n_u)} \\ b^{(1)}, \dots, b^{(n_u)} \\ x^{(1)}, \dots, x^{(n_m)}}} \frac{1}{2} \sum_{(i,j): r(i,j)=1} (w^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^{n_u} \sum_{j=1}^n (w_k^{(j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$

Just to write out all the ratings including the question marks in a more

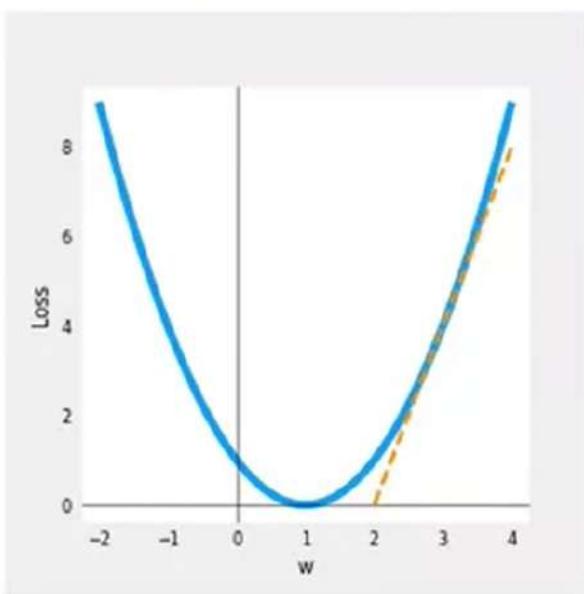
$w = \begin{bmatrix} \vdots \\ 0 \end{bmatrix} \quad b = 0 \quad w^{(s)} \cdot x^{(i)} + b^{(s)}$

$$J = (wx - 1)^2$$

Gradient descent algorithm  
Repeat until convergence

$$w = w - \alpha \frac{d}{dw} J(w, b)$$

Fix  $b = 0$  for this example



# Custom Training Loop

```
w = tf.Variable(3.0)
x = 1.0
y = 1.0 # target value
alpha = 0.01

iterations = 30
for iter in range(iterations):
    # Use TensorFlow's Gradient tape to record the steps
    # used to compute the cost J, to enable auto differentiation.
    with tf.GradientTape() as tape:
        fwb = w*x
        costJ = (fwb - y)**2

    # Use the gradient tape to calculate the gradients
    # of the cost with respect to the parameter w.
    [dJdw] = tape.gradient(costJ, [w])

    # Run one step of gradient descent by updating
    # the value of w to reduce the cost.
    w.assign_add(-alpha * dJdw)
```

Tf.variables are the parameters we want to optimize

**tf variables require special function to filtering algorithm correctly.**

# Implementation in TensorFlow

```
# Instantiate an optimizer.  
optimizer = keras.optimizers.Adam(learning_rate=1e-1) ←  
  
iterations = 200 ←  
for iter in range(iterations):  
    # Use TensorFlow's GradientTape  
    # to record the operations used to compute the cost  
    with tf.GradientTape() as tape: ←  
  
        # Compute the cost (forward_pass is included in cost)  
        cost_value = cofiCostFuncV(X, W, b, Ynorm, R, } ←  
            num_users, num_movies, lambda) ←  
            nu      nm ←  
        # Use the gradient tape to automatically retrieve  
        # the gradients of the trainable variables with respect to  
        # the loss  
        grads = tape.gradient(cost_value, [X,W,b]) ←  
  
        # Run one step of gradient descent by updating  
        # the value of the variables to minimize the loss.  
        optimizer.apply_gradients(zip(grads, [X,W,b]))
```

And then by asking it to give  
you grads equals tape.gradient,

Dataset credit: Harper and Konstan. 2015. The MovieLens Datasets: History and Context.

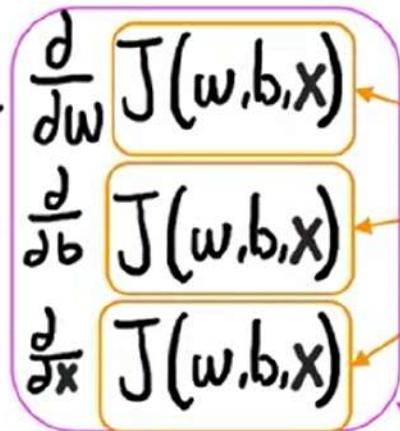
# Implementation in TensorFlow

Gradient descent algorithm

$$w = w - \alpha$$

$$b = b - \alpha$$

X



```
# Instantiate an optimizer.  
optimizer = keras.optimizers.Adam(learning_rate=1e-1)  
  
iterations = 200  
for iter in range(iterations):  
    # Use TensorFlow's GradientTape  
    # to record the operations used to compute the cost  
    with tf.GradientTape() as tape:  
  
        # Compute the cost (forward_pass is included in cost)  
        cost_value = cofiCostFuncV(X, W, b, Ynorm, R,  
            num_users, num_movies, lambda)  
            nu nm  
  
        # Use the gradient tape to automatically retrieve  
        # the gradients of the trainable variables with respect to  
        # the loss  
        grads = tape.gradient(cost_value, [X, W, b])  
  
        # Run one step of gradient descent by updating  
        # the value of the variables to minimize the loss.  
        optimizer.apply_gradients(zip(grads, [X, W, b]))
```

Dataset credit: Harper and Konstan. 2015. The MovieLens Datasets: History and Context.

And you repeat until convergence.

# Finding related items

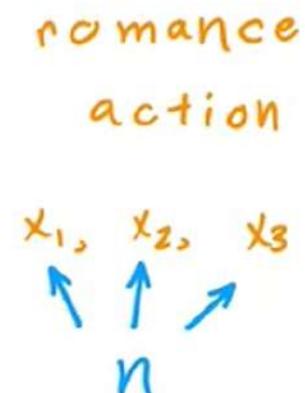
The features  $x^{(i)}$  of item  $i$  are quite hard to interpret.

To find other items related to it,

find item  $k$  with  $x^{(k)}$  similar to  $x^{(i)}$

i.e. with smallest  
distance

$$\sum_{l=1}^n (x_l^{(k)} - x_l^{(i)})^2$$
$$\|x^{(k)} - x^{(i)}\|^2$$



an even more powerful  
recommended system as well.

## Collaborative filtering vs Content-based filtering

- Collaborative filtering:

Recommend items to you based on ratings of users who gave similar ratings as you

- Content-based filtering:

Recommend items to you based on features of user and item to find good match

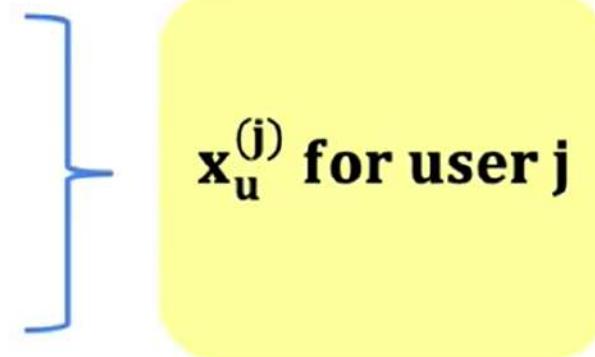
$r(i,j) = 1$  if user  $j$  has rated item  $i$

$y^{(i,j)}$  rating given by user  $j$  on item  $i$  (if defined)  
 $j$  to denote the rating

## Examples of user and item features

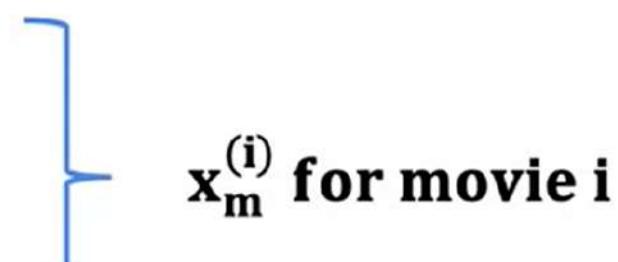
### User features:

- • Age
- • Gender (1 hot)
- • Country (1 hot, 200)
- • Movies watched (1000)
- • Average rating per genre
- ...



### Movie features:

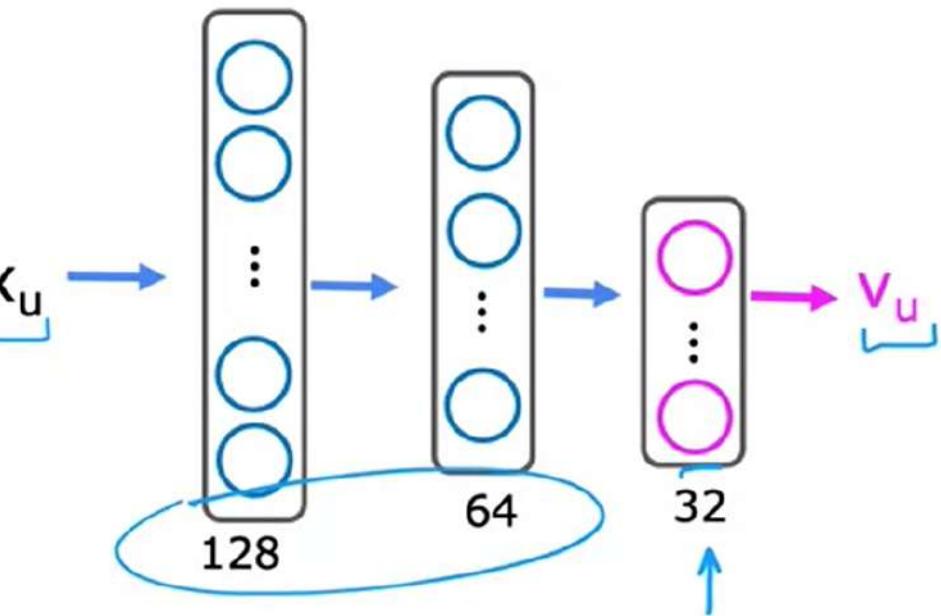
- • Year
- • Genre/Genres
- • Reviews
- • Average rating
- ...



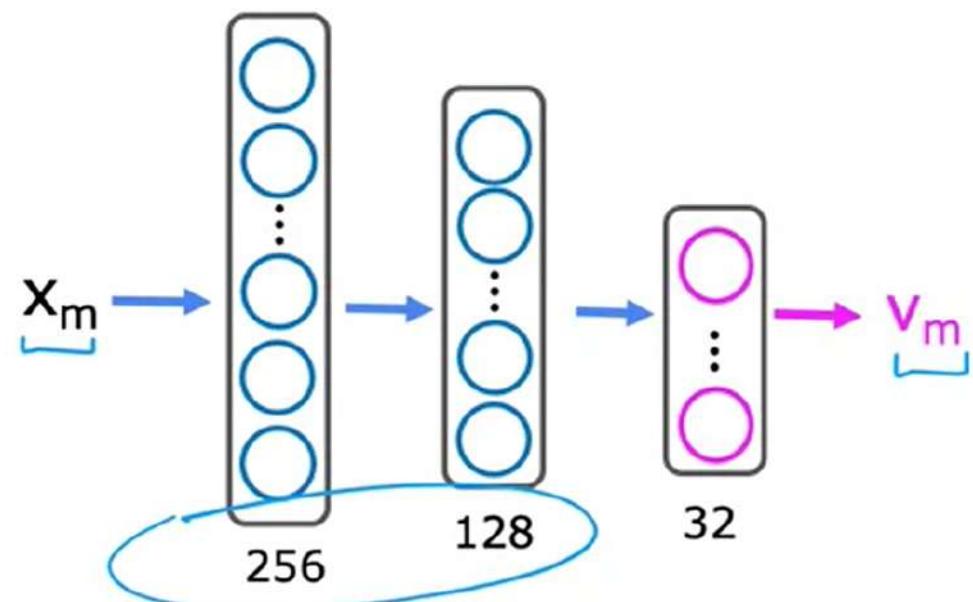
This feature again  
depends on the ratings

## Neural network architecture

$x_u \rightarrow v_u$  User network



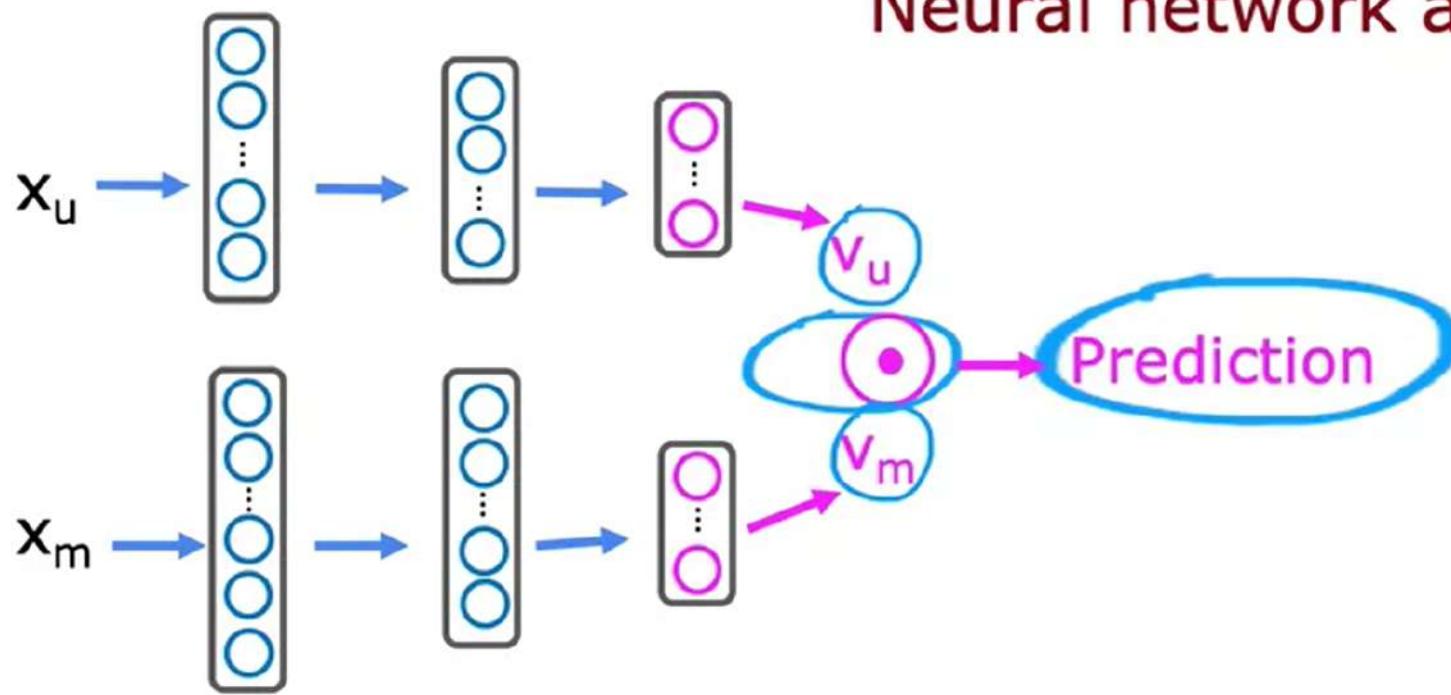
$x_m \rightarrow v_m$  Movie network



**Prediction :**  $v_u \cdot v_m$

All the output layer needs to

## Neural network architecture



Cost function 
$$J = \sum_{(i,j):r(i,j)=1} (v_u^{(j)} \cdot v_m^{(i)} - y^{(i,j)})^2 + \text{NN regularization term}$$
  
v\_m^i minus y^ij squared.

## Learned user and item vectors:

- $v_u^{(j)}$  is a vector of length 32 that describes user j with features  $x_u^{(j)}$
- $v_m^{(i)}$  is a vector of length 32 that describes movie i with features  $x_m^{(i)}$

To find movies similar to movie i:

$$\|v_m^{(k)} - v_m^{(i)}\|^2 \text{ small}$$

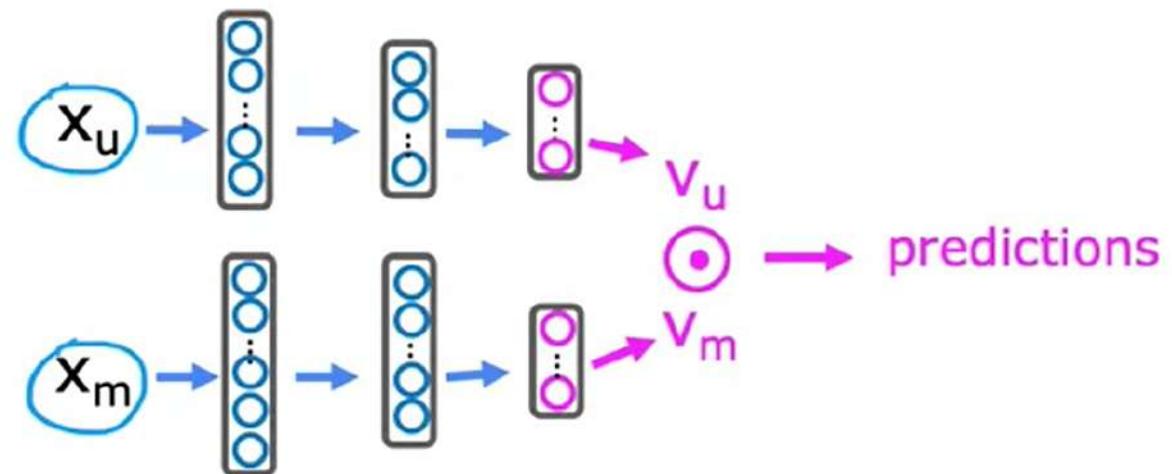
$$\|x^{(k)} - x^{(i)}\|^2$$

Note: This can be pre-computed ahead of time  
later when we talk about scaling

## How to efficiently find recommendation from a large set of items?

- • Movies
- • Ads
- • Songs
- • Products

1000+  
1m+  
10m+  
10m+



becomes computational e infeasible.

## Two steps: Retrieval & Ranking

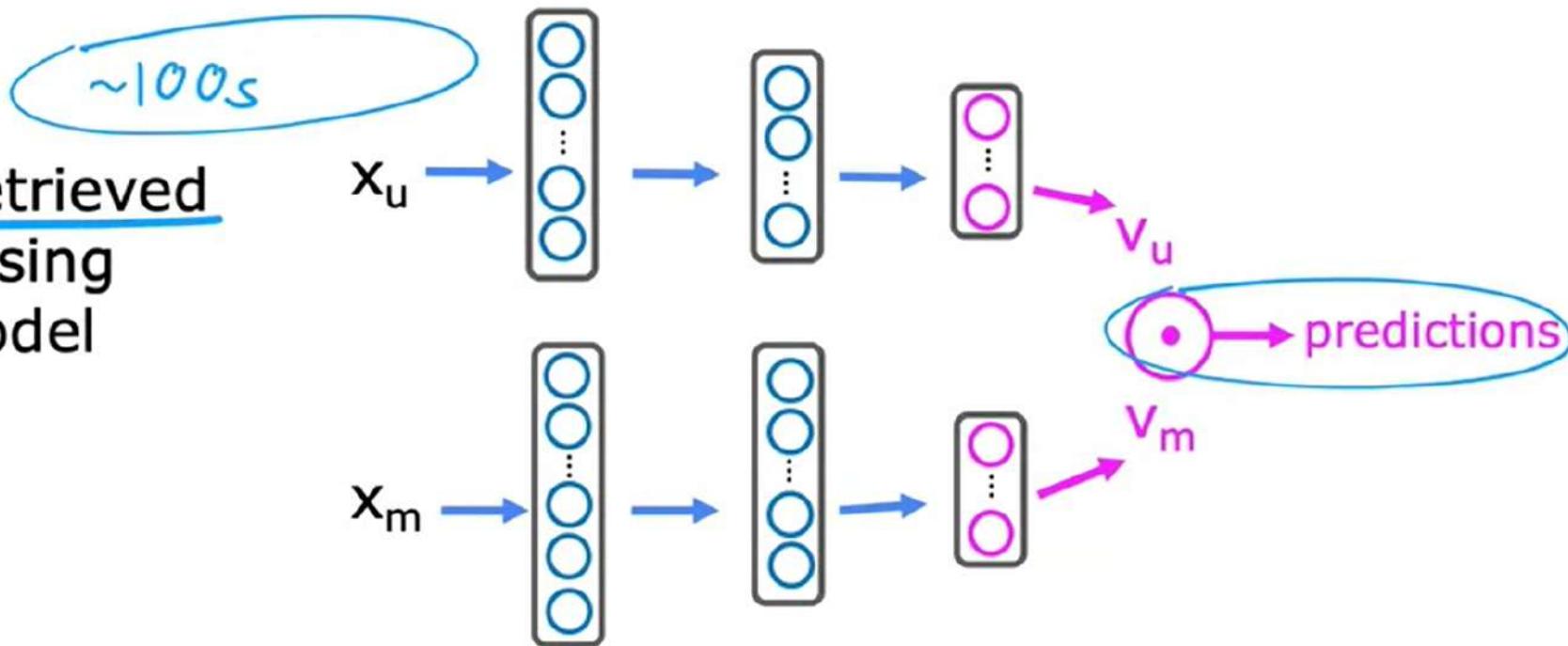
### Retrieval:

- • Generate large list of **plausible item** candidates  $\sim 100s$ 
  - e.g.
    - 1) For each of the last 10 movies watched by the user, find 10 most similar movies
$$\| v_m^{(k)} - v_m^{(i)} \|^2$$
    - 2) For most viewed 3 genres, find the top 10 movies
    - 3) Top 20 movies in the country
- • Combine **retrieved items** into list, removing duplicates and items already watched/purchased that the user has already purchased and

## Two steps: Retrieval & ranking

Ranking:

- Take list retrieved and rank using learned model



- Display ranked items to user

what you think the user will give.

# What is the goal of the recommender system?

Recommend:

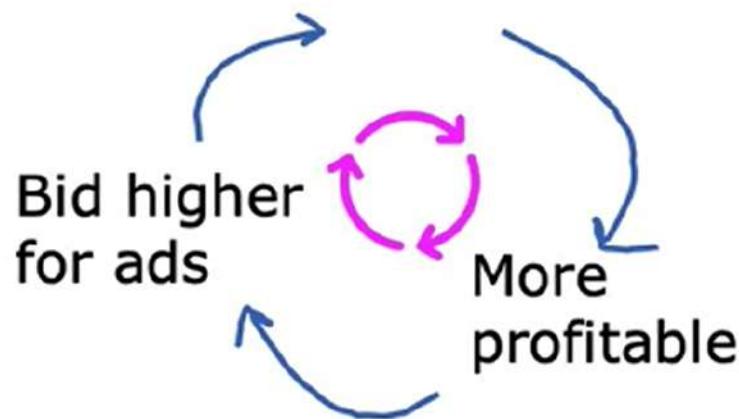
- • Movies most likely to be rated 5 stars by user
- • Products most likely to be purchased
- • Ads most likely to be clicked on *+high bid*
- • Products generating the largest profit
- • Video leading to maximum watch time

Specifically, websites that are

# Ethical considerations with recommender systems

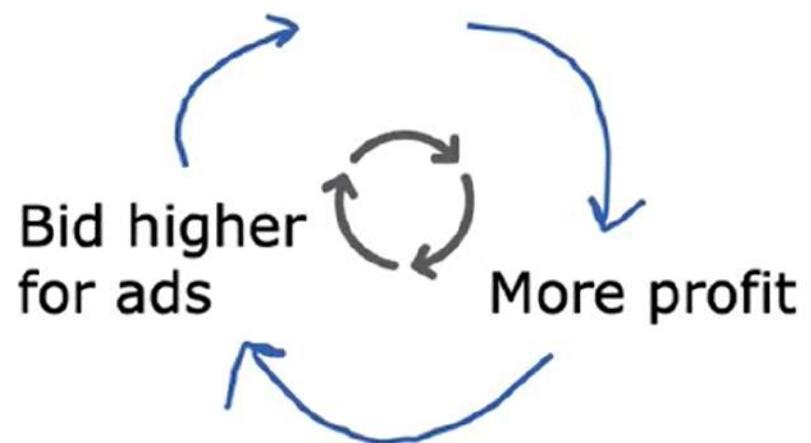
## Travel industry

Good travel experience  
to more users



## Payday loans

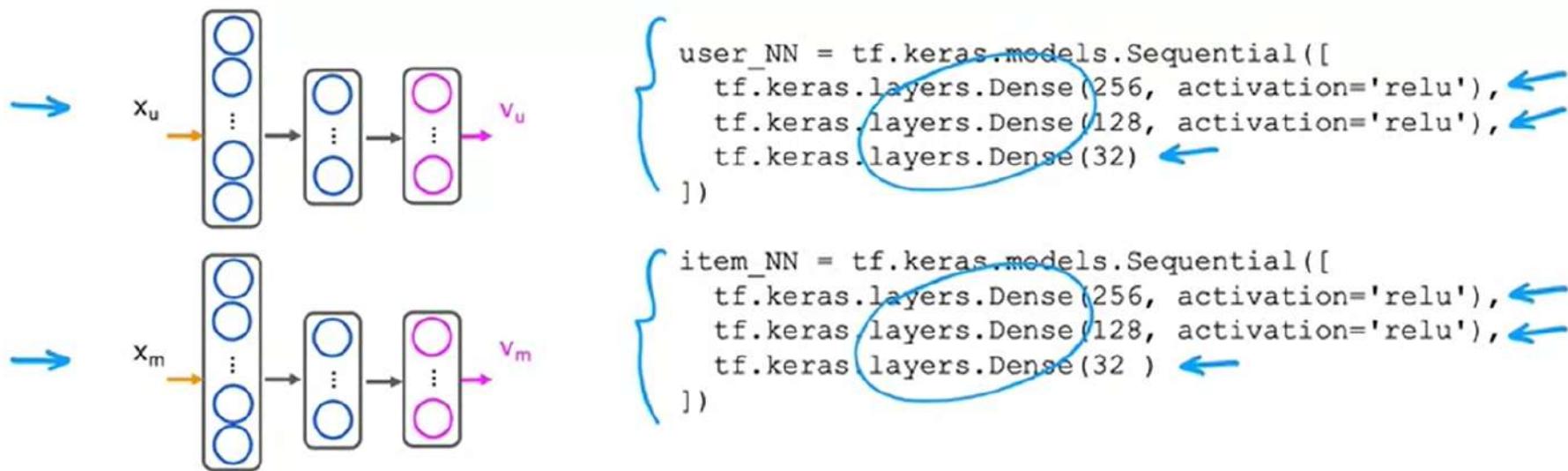
Squeeze customers  
more



to get sent more traffic.

## Other problematic cases:

- • Maximizing user engagement (e.g. watch time) has led to large social media/video sharing sites to amplify conspiracy theories and hate/toxicity
- Amelioration : Filter out problematic content such as hate speech, fraud, scams and violent content
- • Can a ranking system maximize your profit rather than users' welfare be presented in a transparent way?
- Amelioration : Be transparent with users
  - transparent with
  - users about why we're



```

# create the user input and point to the base network
input_user = tf.keras.layers.Input(shape=(num_user_features))
vu = user_NN(input_user)
vu = tf.linalg.l2_normalize(vu, axis=1)

# create the item input and point to the base network
input_item = tf.keras.layers.Input(shape=(num_item_features))
vm = item_NN(input_item)
vm = tf.linalg.l2_normalize(vm, axis=1)

# measure the similarity of the two vector outputs
output = tf.keras.layers.Dot(axes=1)([vu, vm])

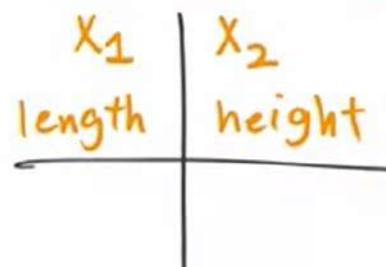
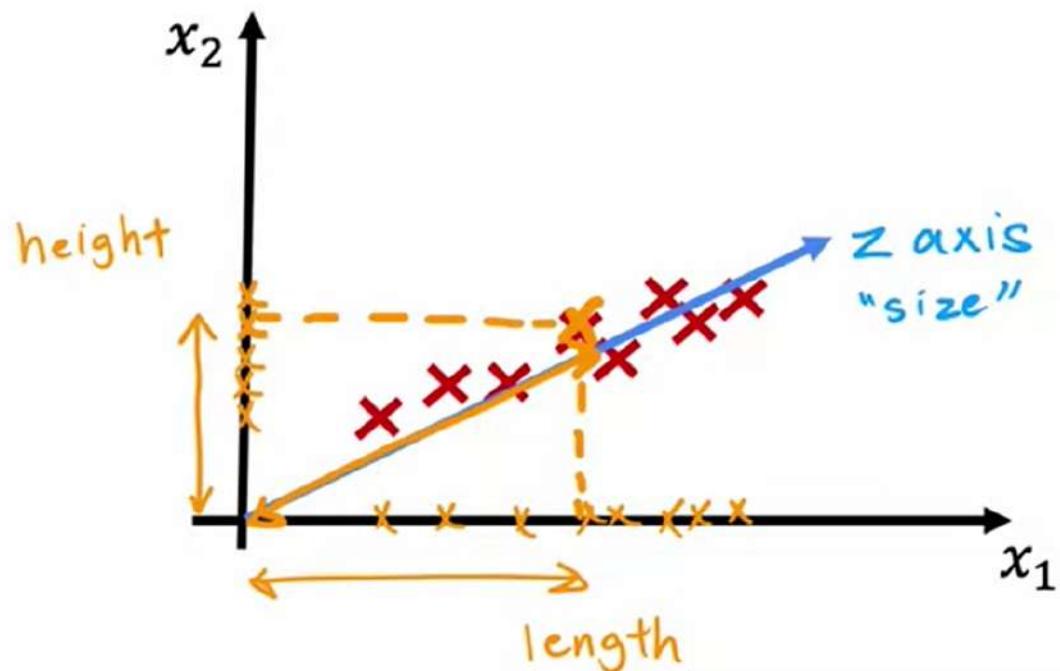
# specify the inputs and output of the model
model = Model([input_user, input_item], output)

# Specify the cost function
cost_fn = tf.keras.losses.MeanSquaredError()

```

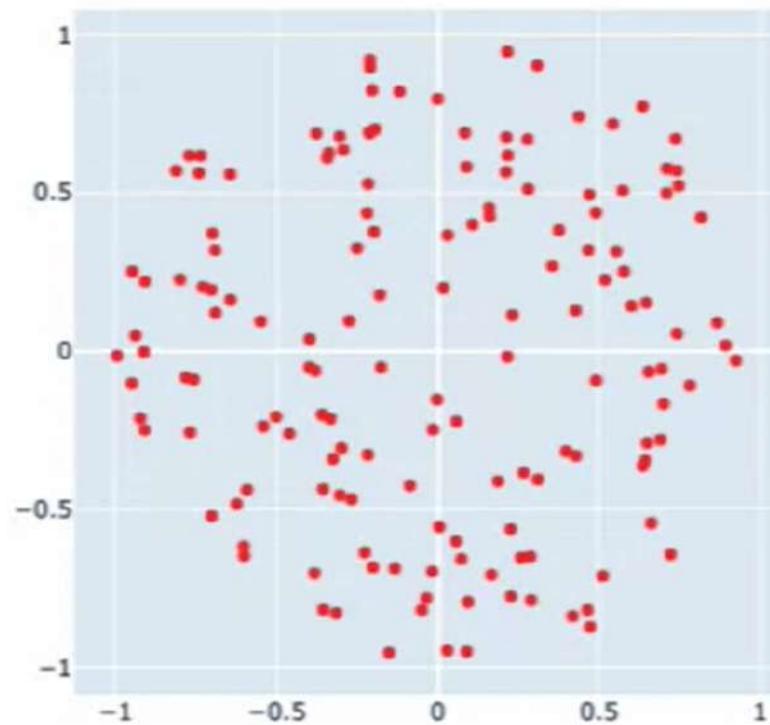
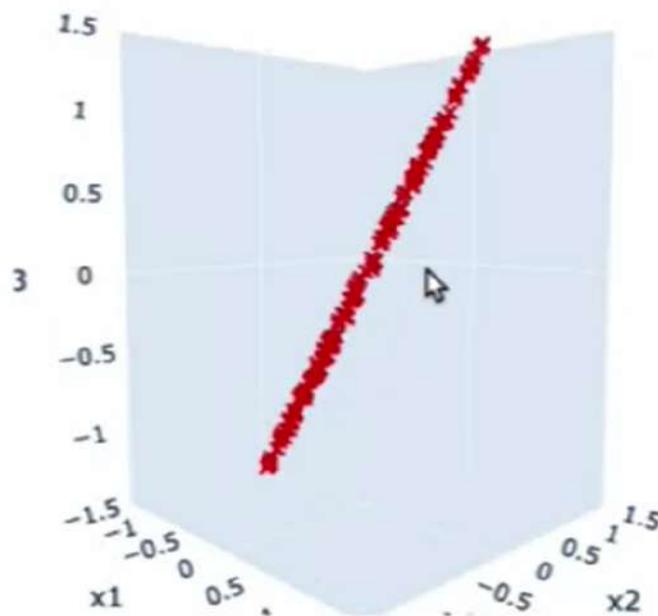
# Size

PCA: find new axis and coordinates



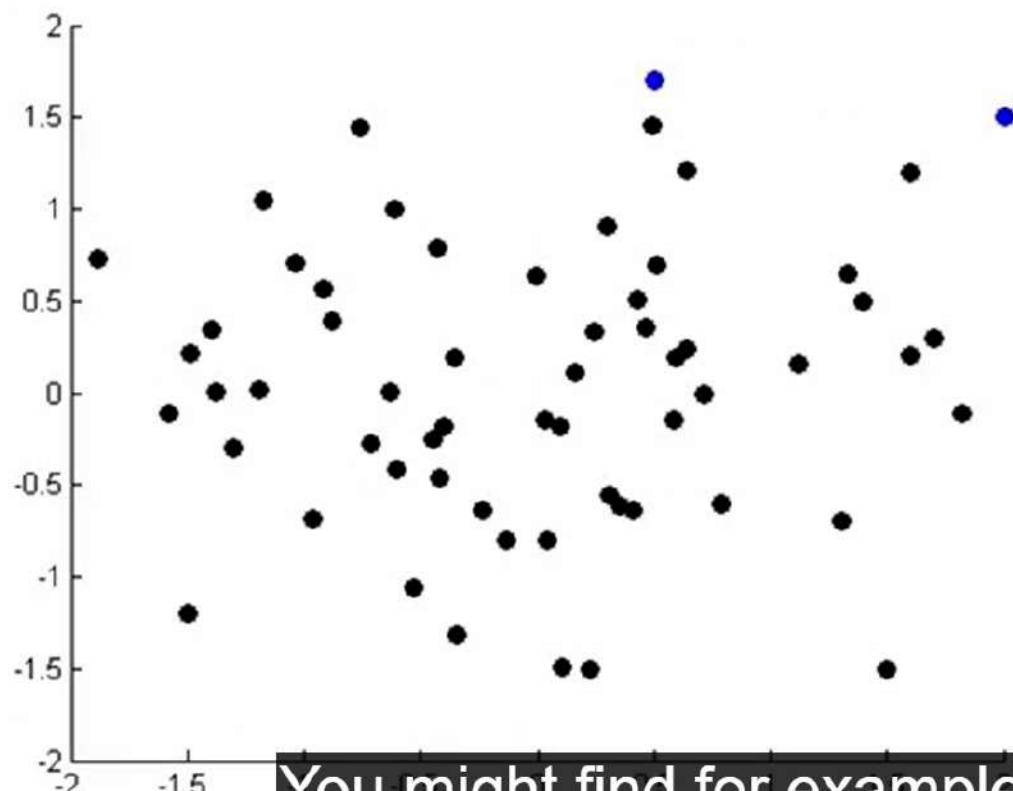
But the idea of PCA is to  
find one or more new axes,

# From 3D to 2D



on a two-dimensional pancake,

# Data visualization



You might find for example  
that  $Z_1$  loosely corresponds

# Anomaly detection algorithm

1. Choose  $n$  features  $x_i$  that you think might be indicative of anomalous examples.
2. Fit parameters  $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

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

Vectorized formula

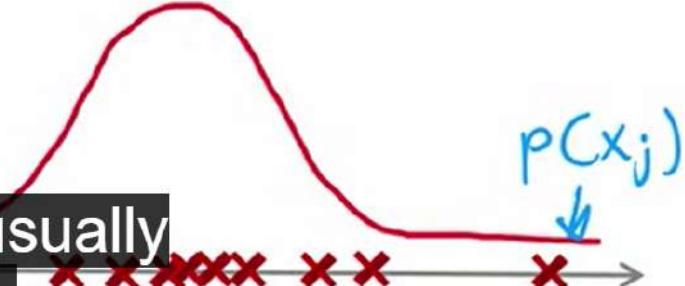
$$\vec{\mu} = \frac{1}{m} \sum_{i=1}^m \vec{x}^{(i)}$$
$$\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

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 if  $p(x) < \varepsilon$

has any features that are unusually large or unusually small.



50 features

Country	$x_1$ GDP (trillions of US\$)	$x_2$ Per capita GDP (thousands of int'l. \$)	$x_3$ Human Development Index	$x_4$ 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	—

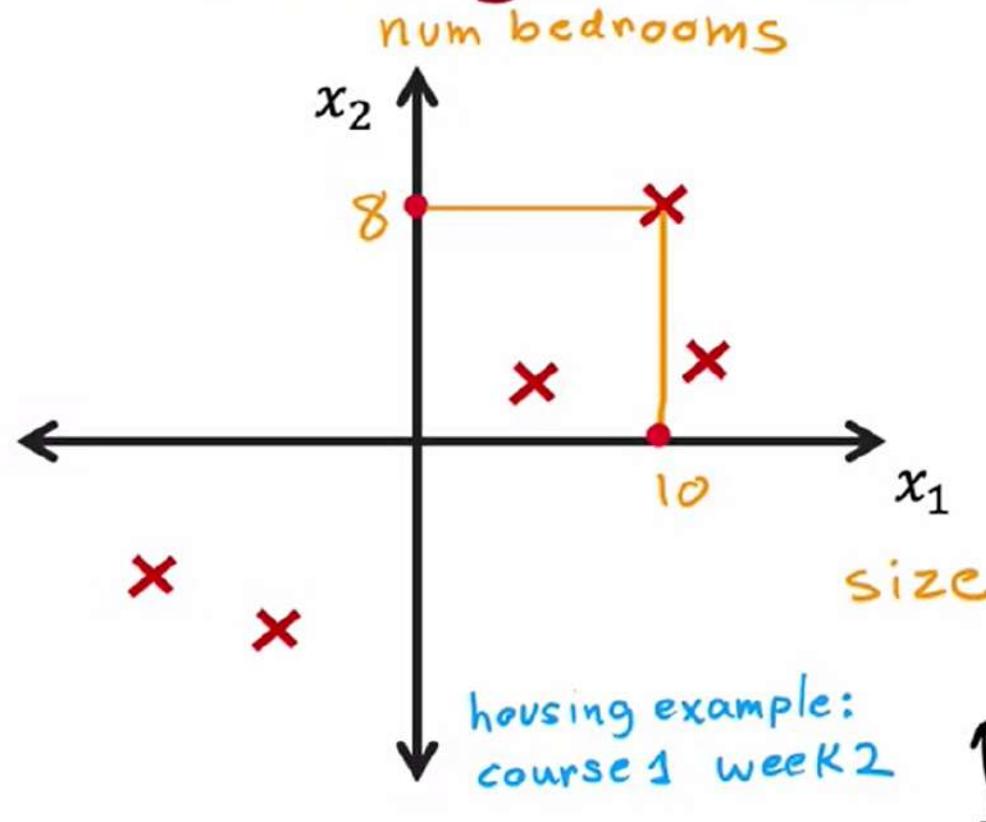
Country	$z_1$	$z_2$
Canada	1.6	1.2
China	1.7	0.3
India	1.6	0.2
Russia	1.	Z_1 and Z_2 and you can then

# PCA algorithm

coordinates

$$x_1 = 10 \quad x_2 = 8$$

Can we choose  
a different axis?



Preprocess features

Normalized to  
have zero mean

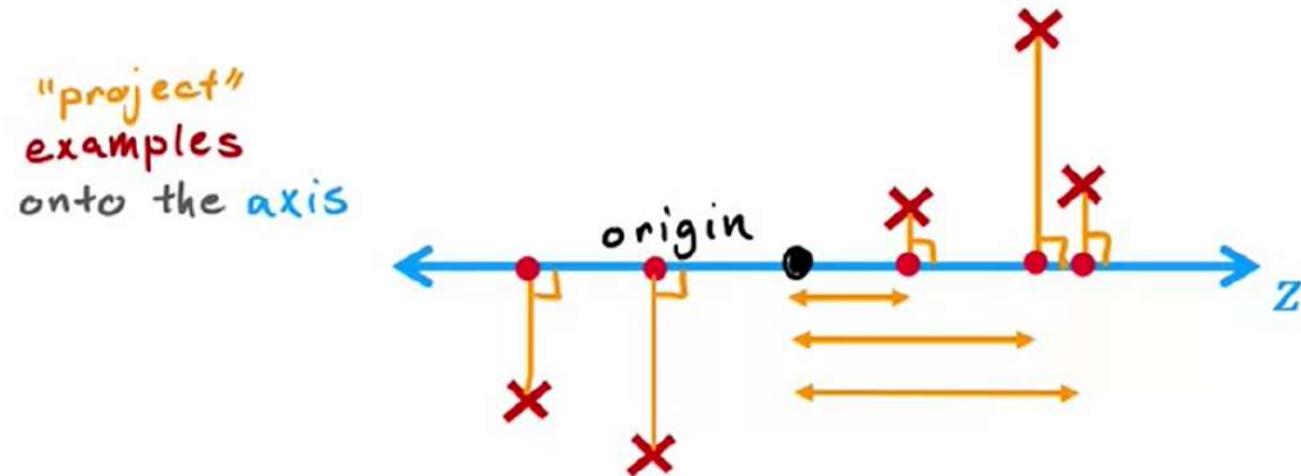


feature scaling



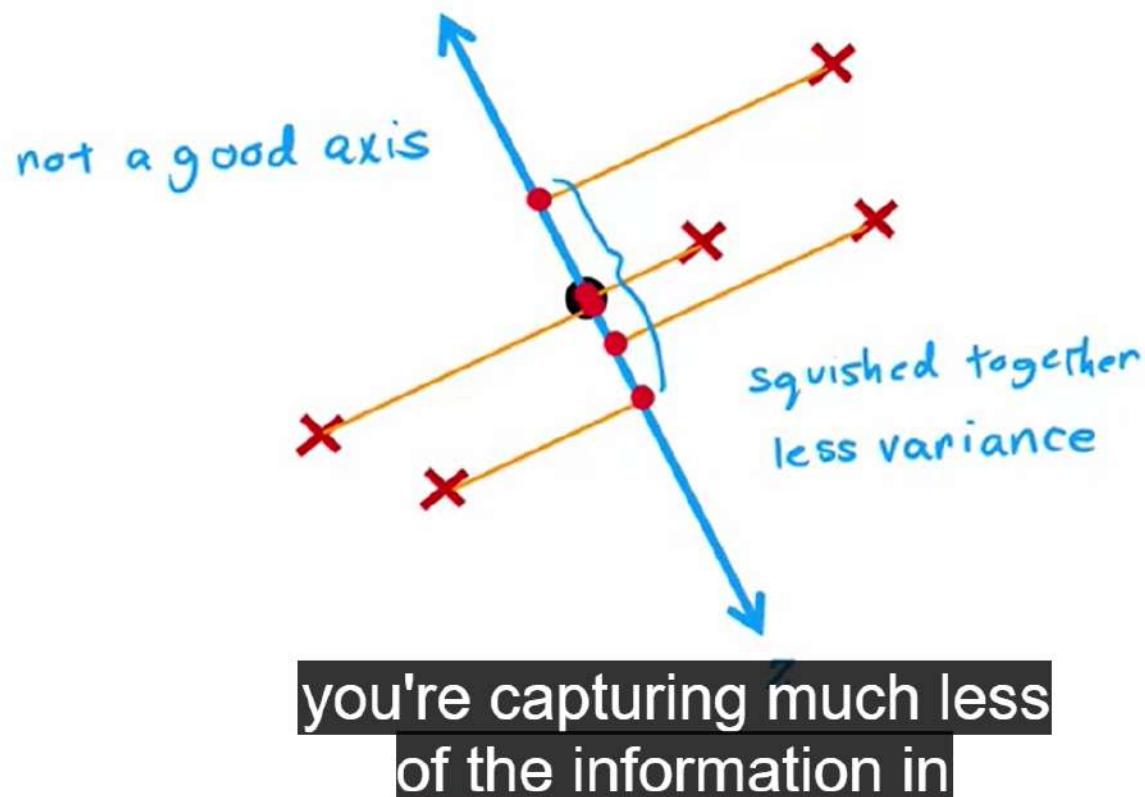
To examine what PCA does,

# Choose an axis

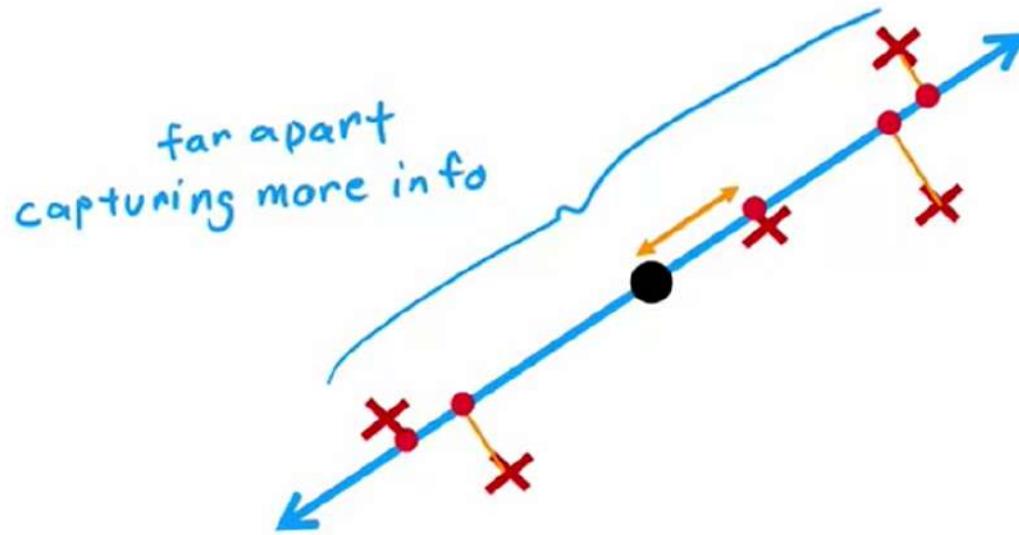


This choice isn't  
too bad because

# Choose an axis

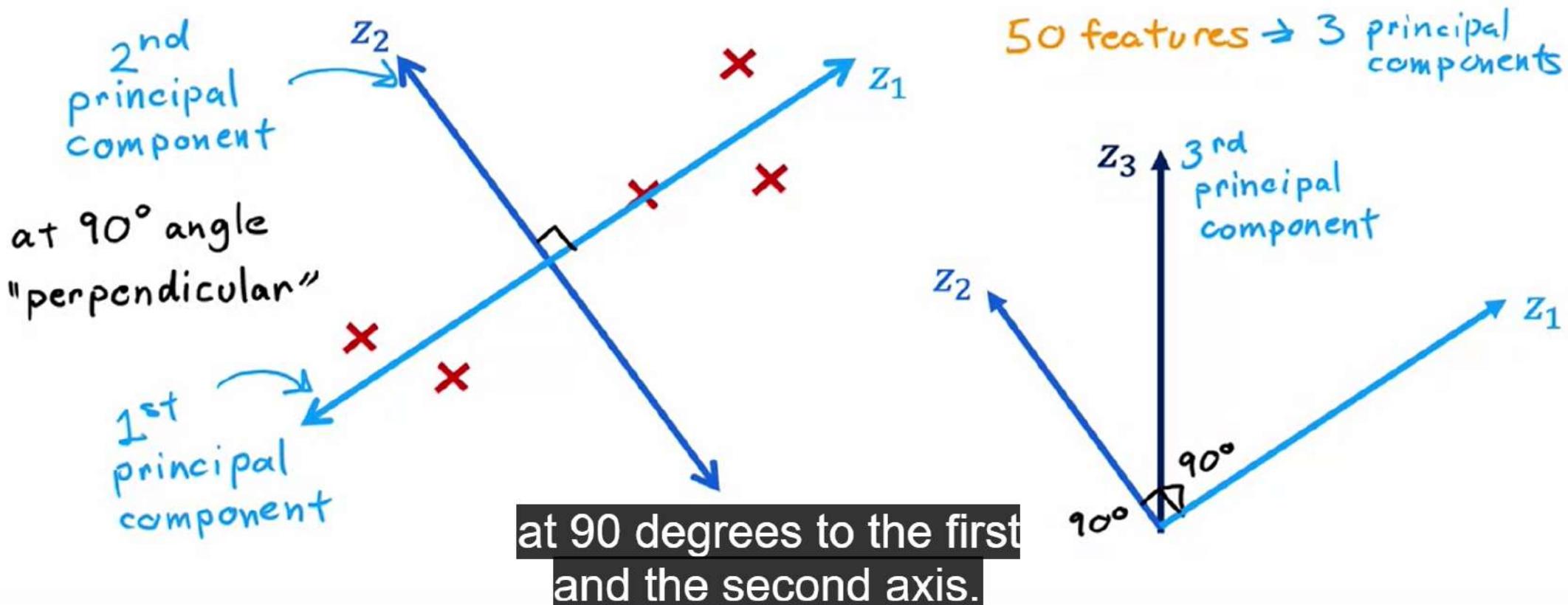


# Choose an axis

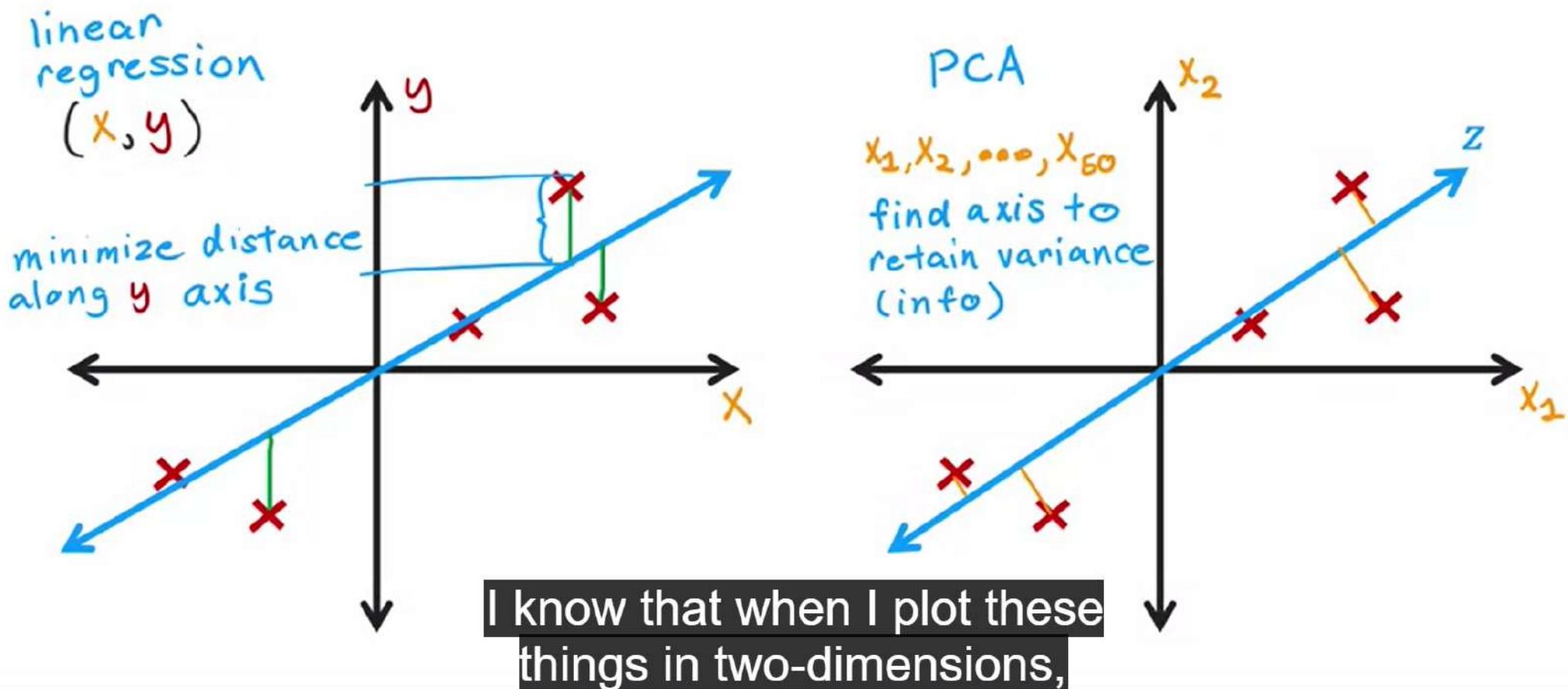


this axis is called the  
principal component.

# More principal components

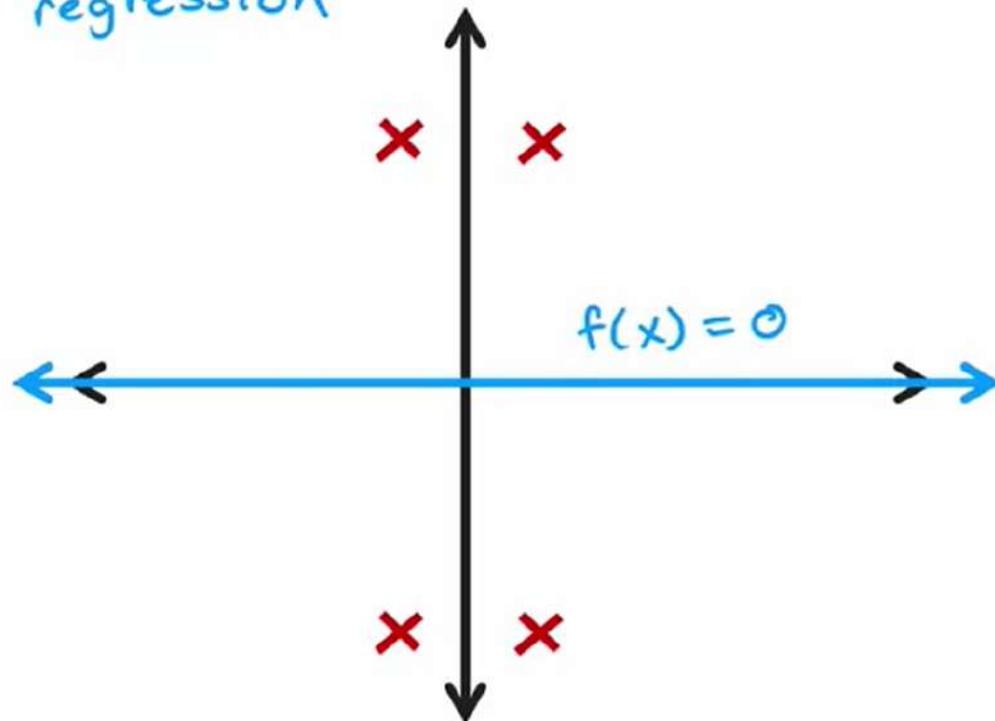


# PCA is not linear regression

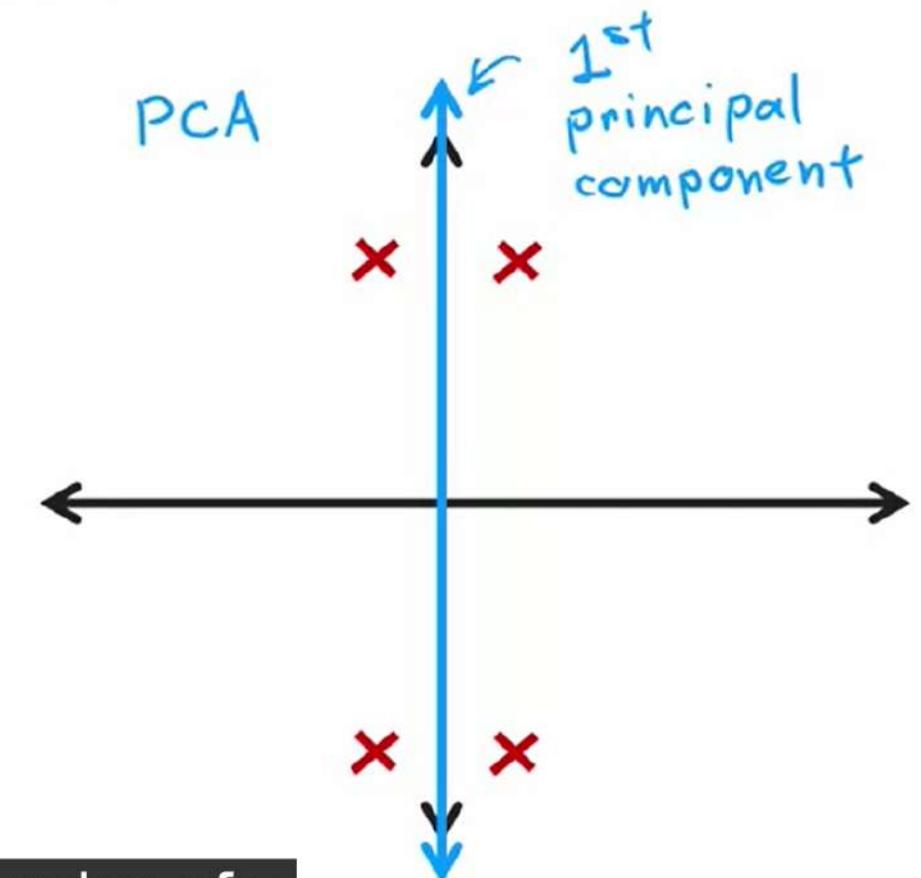


# PCA is not linear regression

linear  
regression



PCA



to predict the value of  $y$ ,

# PCA in scikit-learn

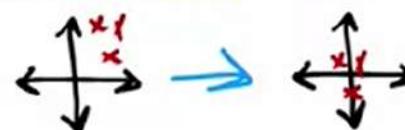


Optional pre-processing: Perform feature scaling

*for visualization*

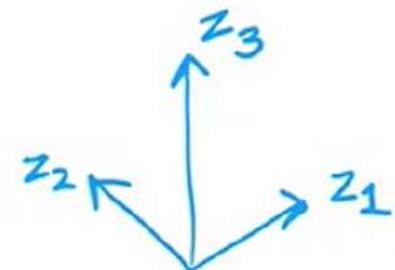
1. "fit" the data to obtain 2 (or 3) new axes (principal components)

*fit includes mean normalization*



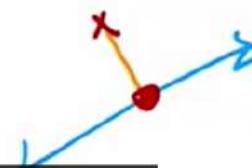
2. Optionally examine how much variance is explained by each principal component.

*info explained\_variance\_ratio*

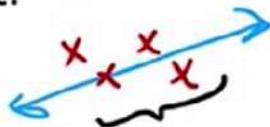


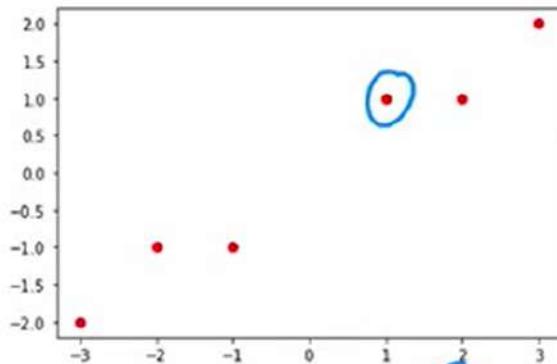
3. Transform (project) the data onto the new axes

*transform*



**to visualize your data.**



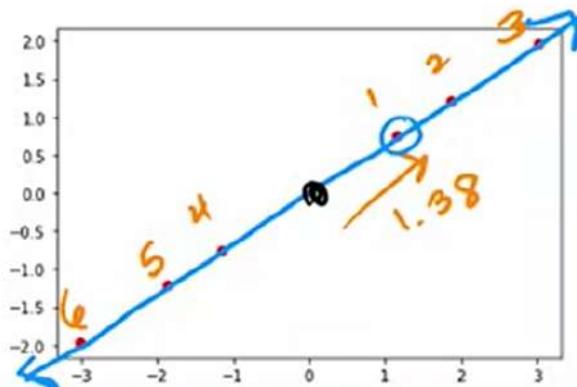


# Example

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

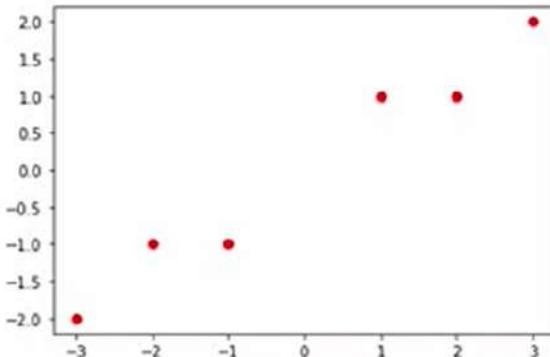
2D

```
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 ]])
```

This data is  
two-dimensional data,



# Example

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

2D

```
pca_2 = PCA(n_components=2)
```

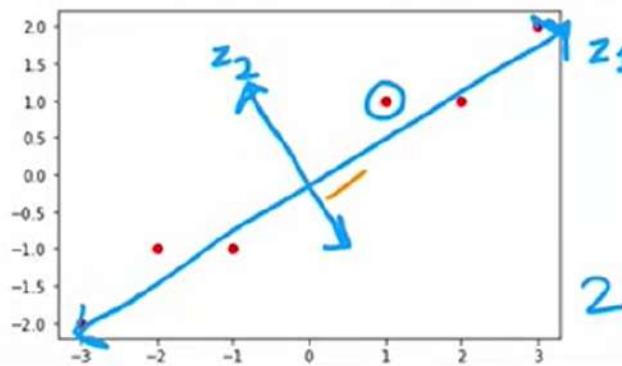
```
pca_2.fit(X)
```

```
pca_2.explained_variance_ratio_
```

$z_1$      $z_2$   
0.992    0.008

```
X_trans_2 = pca.transform(X)
```

```
X_reduced_2 = pca.inverse_transform(X_trans_2)
```



2D  
[1] has a distance of  
1.38 on the  $z_1$  axis,

$z_1$                    $z_2$   
array([  
 → [ 1.38340578, 0.2935787 ],  
 [ 2.22189802, -0.25133484 ],  
 [ 3.6053038 , 0.04224385 ],  
 [-1.38340578, -0.2935787 ],  
 [-2.22189802, 0.25133484 ],  
 [-3.6053038 , -0.04224385 ]])

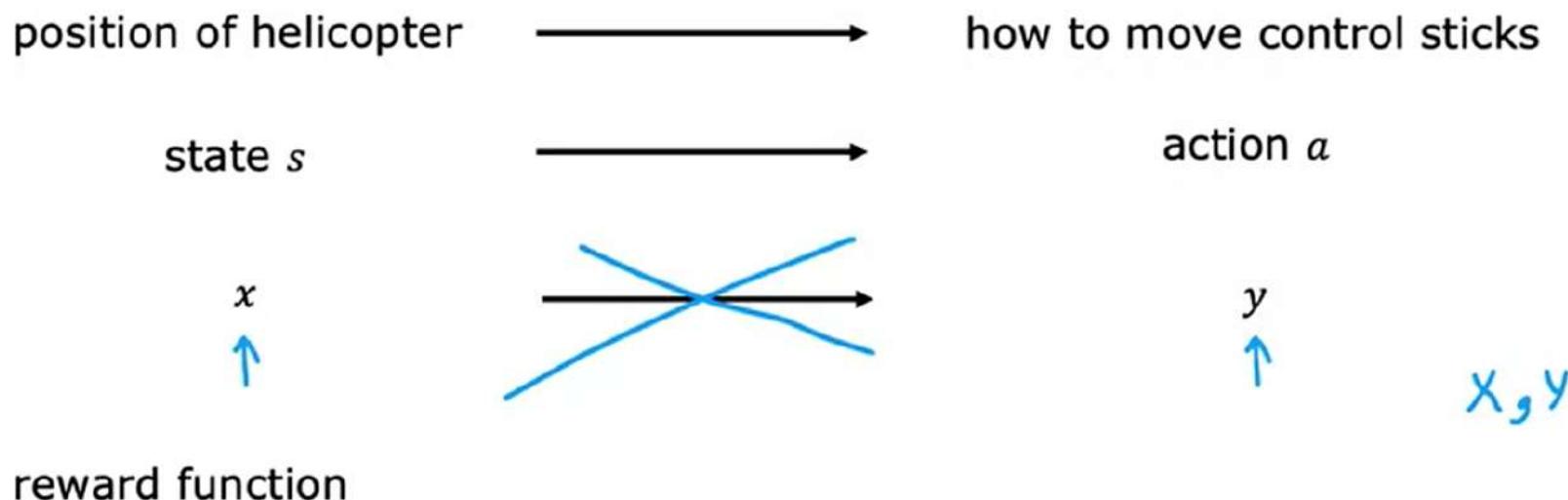


# Reinforcement Learning Introduction

---

What is Reinforcement  
specialization, but  
I'm looking forward to this week,

# Reinforcement Learning



positive reward : helicopter flying well +1

negative reward : helicopter flying poorly -1000

time flying well and  
hopefully to never crash.

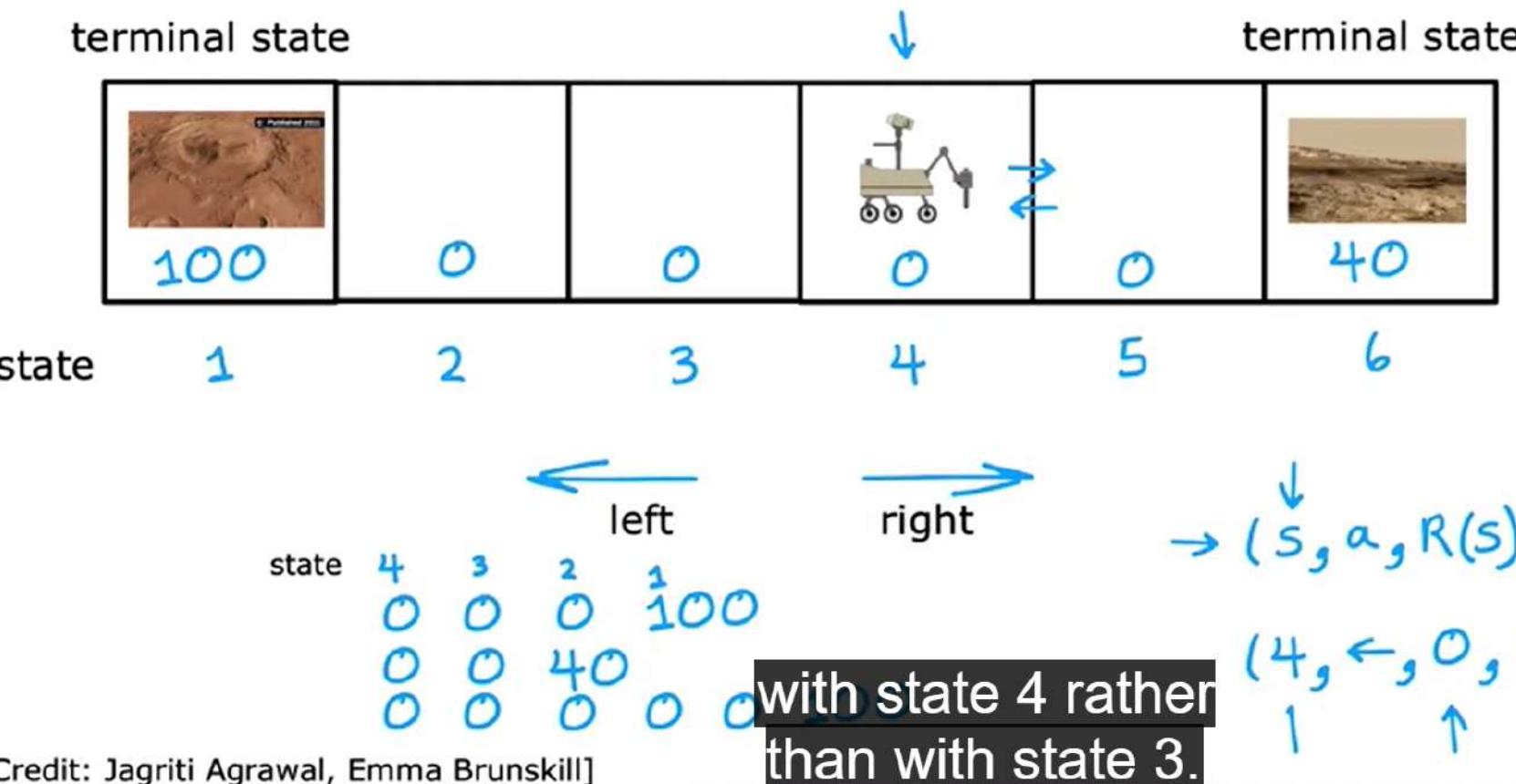
# Applications

- • Controlling robots
- • Factory optimization
- • Financial (stock) trading
- Playing games (including video games)

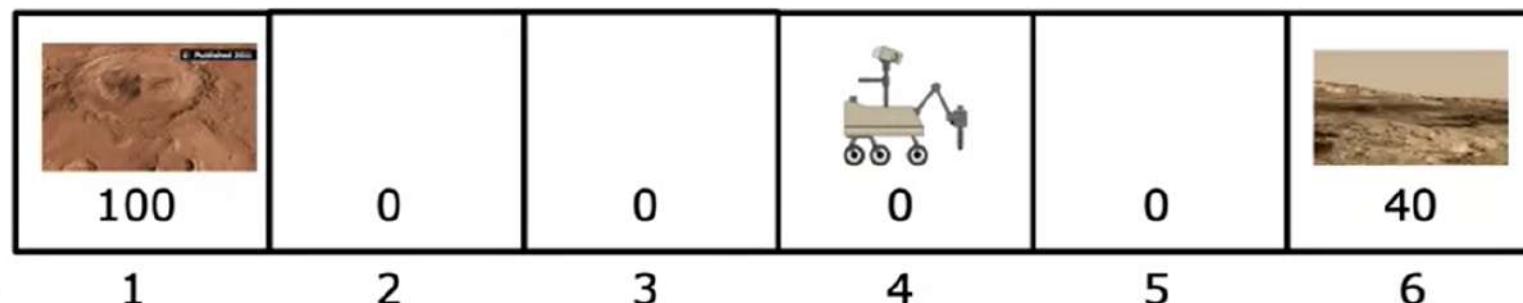


Finally, there have also been many  
applications of reinforcement

# Mars Rover Example



# Return



$$\text{Return} = 0 + (0.9)0 + (0.9)^20 + (0.9)^3100 = 0.729 \times 100 = 72.9$$

$$\text{Return} = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots \quad (\text{until terminal state})$$

Discount Factor  $\gamma = 0.9 \quad 0.99 \quad 0.999$

$$\gamma = 0.5$$

$$\text{Return} = 0 + (0.5)0 + (0.5)^20 + (0.5)^3100 = 12.5$$

and this turns out to  
be a return of 12.5.

# Example of Return

100	50	25	12.5	6.25	40
100	0	0	0	0	40
1	2	3	4	5	6

← return  
← reward

$$\gamma = 0.5$$

The return depends on the actions you take.

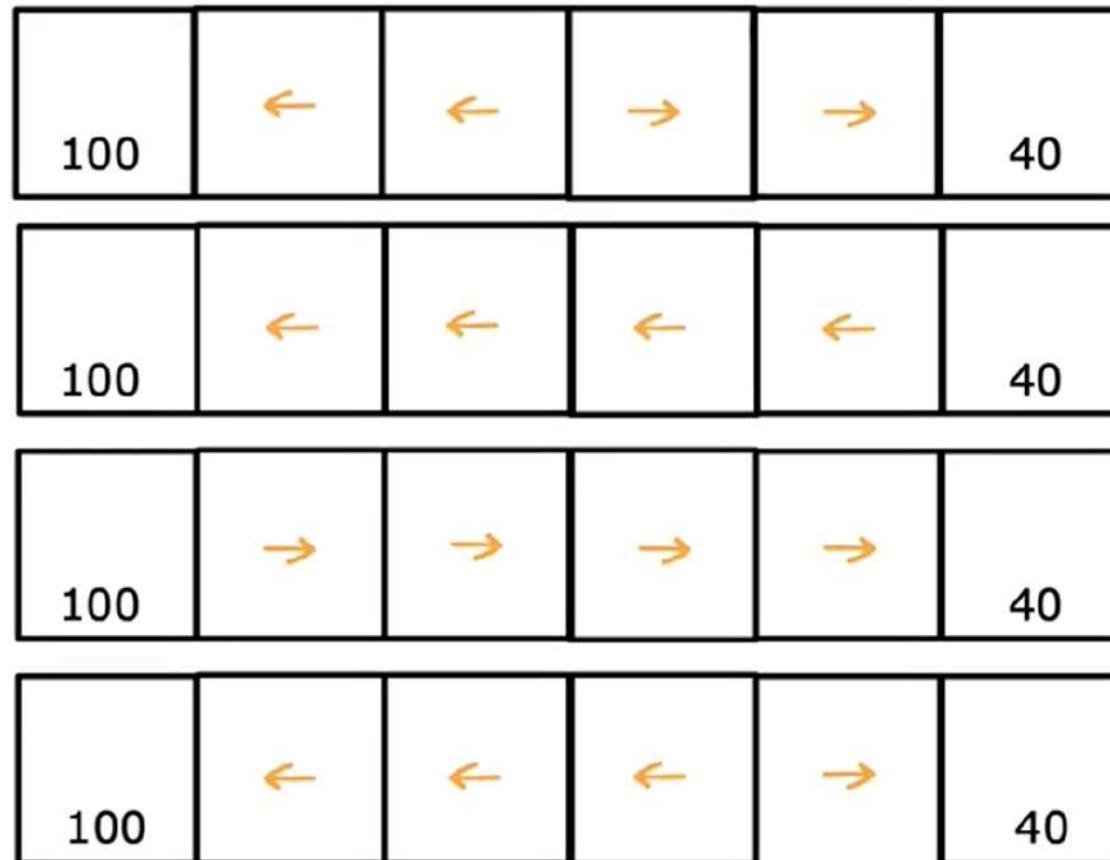
100	2.5	5	10	20	40
100	0	0	0	0	40
1	2	3	4	5	6

$$0 + (0.5)0 + (0.5)^2 40 = 10$$

and if you were to  
always go to the right.

# Policy

state  $s$     policy  $\pi$     action  $a$



$$\begin{aligned}\pi(s) &= a \\ \pi(2) &= \leftarrow \\ \pi(3) &= \leftarrow \\ \pi(4) &= \leftarrow \\ \pi(5) &= \rightarrow\end{aligned}$$

A policy is a function  $\pi(s) = a$ , mapping from states to actions, that tells you what action  $a$  to take in a given state.

Mars rover



Helicopter



Chess



→ states

6 states

→ actions



→ rewards

100, 0, 40

→ discount factor  $\gamma$ 

0.5

→ return

 $R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$ → policy  $\pi$ 

position of helicopter

how to move  
control stick

+1, -1000

0.99

 $R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$ Find  $\pi(s) = a$ 

actually has a name.

pieces on board

possible move

+1, 0, -1

0.995

 $R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$ Find  $\pi(s) = a$ 