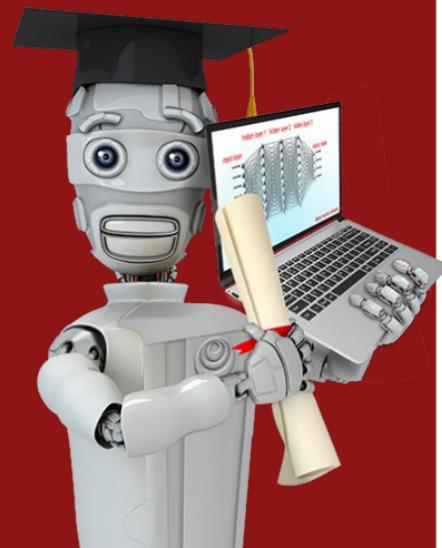


# Recommender Systems

---



## Recommender System

### Making recommendations

# Predicting movie ratings

User rates movies using one to five stars

~~zero~~

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

Ratings				
★				
★	★			
★	★	★		
★	★	★	★	
★	★	★	★	★

$n_u$  = no. of users

$n_m$  = no. of movies

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

$$n_u = 4$$

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

$$n_m = 5$$

$$r(3,1) = 0$$

$$y^{(3,2)} = 4$$

$y^{(i,j)}$  = rating given by user  $j$  to movie  $i$  (defined only if  $r(i,j)=1$ )

## Collaborative Filtering

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  $i$ .

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

For user  $j$ : Predict user  $j$ 's rating for movie  $i$  as

This can be rewritten as  $w^{(j)} \cdot X^{(i)}$  if there's  $X_0^{(i)} = 1$ .

$w^{(j)} \cdot X^{(i)} + b^{(j)}$  just linear regression

For example, if  $X_1$  is romance,  $X_2$  is action.  $X^{(1)} = \begin{pmatrix} 0.9 \\ 0 \end{pmatrix}$ ,  $X^{(3)} = \begin{pmatrix} 0.99 \\ 0 \end{pmatrix}$

and for user 1:  $w^{(1)} = \begin{pmatrix} 5 \\ 0 \end{pmatrix}$ ,  $b^{(1)} = 0$ ,

then predict rating for movie 3 as:  $w^{(1)} \cdot X^{(3)} + b^{(1)} = 4.95$

The aim: to learn  $w^{(j)}, b^{(j)}$

Cost function: for single user  $j$ , given  $X^{(i)}$

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

Only sum over the movies  $i$  that user  $j$  has actually rated.

$\cancel{m^{(j)}}$ :  $m^{(j)}$  is eliminated because it's a constant.

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

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

# Problem motivation

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 → ?
Cute puppies of love	?	4	0	?	0.99 → ?	-0 → ?
Nonstop car chases	0	0	5	4	-0.1 → ?	-1.0 → ?
Swords vs. karate	0	0	5	?	-0 → ?	-0.9 → ?

The cost function

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

$x^{(i)}$  is unknown

is built on the assumption that the feature vector for movie  $i$ ,  $x^{(i)}$ , is known beforehand.

But, in reality,  $x^{(i)}$  is unknown.

# Collaborative Filtering Algorithm

## Problem Motivation

Since  $X^{(i)}$  is unknown, so we need to learn  $X^{(i)}$ . But how?

Let's assume  $w^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$ ,  $w^{(2)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$ ,  $w^{(3)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$ ,  $w^{(4)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$ ,  $b^{(1)} = b^{(2)} = b^{(3)} = b^{(4)} = 0$

then  $\begin{aligned} w^{(1)} \cdot X^{(1)} &\approx 5 \\ w^{(2)} \cdot X^{(1)} &\approx 5 \\ w^{(3)} \cdot X^{(1)} &\approx 0 \\ w^{(4)} \cdot X^{(1)} &\approx 0 \end{aligned}$  From these four equations, we can derive  $X^{(1)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

Notice that this works only because we have parameters for four users. That's what allows us to try to guess  $X^{(1)}$ . This is why in a typical linear regression application if we had just a single user, we don't actually have enough information to figure out what should be the features,  $X_1$  and  $X_2$ , which is why in the linear regression contexts that we saw in Course 1, we can't come up with features  $X_1$  and  $X_2$  from scratch.

But in collaborative filtering, we have ratings from multiple users of the same item with the same movie. That's what makes it possible to try to guess the features.

## Cost function:

Given  $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(n_u)}, b^{(n_u)}$  for all users, we want to learn  $X^{(i)}$ :

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

To learn  $X^{(1)}, X^{(2)}, \dots, X^{(n_m)}$

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

Now we are ready to put two cost functions together. Thus, the term collaborative filtering.

Cost function to learn  $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(n_u)}, b^{(n_u)}$ :

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

Cost function to learn  $x^{(1)}, x^{(2)}, \dots, x^{(n_m)}$ :

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

Put them together:

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

Gradient Descent:

Step 1: Initialize  $x^{(1)}, x^{(2)}, \dots, x^{(n_m)}, w^{(1)}, w^{(2)}, \dots, w^{(n_u)}, b^{(1)}, b^{(2)}, \dots, b^{(n_u)}$  to small random values.

Step 2:

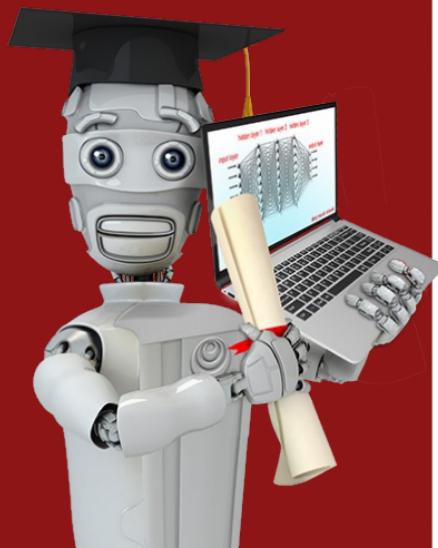
Repeat {

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

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

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

}



# Collaborative Filtering

Binary labels:  
favs,  
likes and clicks

# Binary labels

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

1 : hit like

0 : not hit like

? : not watch the movie

The binary labels can represent many features.  
It depends on the questions.

# 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

# From regression to binary classification

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

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

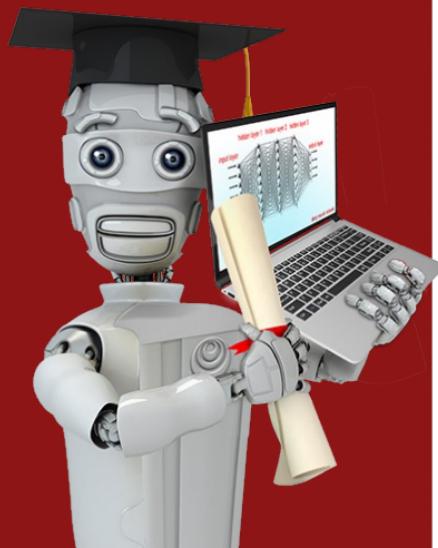
$$L(f_{(w,b,x)}(x), y^{(i,j)}) = -y^{(i,j)} \log(f_{(w,b,x)}(x)) - (1 - y^{(i,j)}) \log(1 - f_{(w,b,x)}(x))$$

Loss for single example

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

$\leftarrow$  cost for all examples

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



# Recommender Systems implementation

## Mean normalization

In the first course, Andrew showed that feature normalization can help the gradient descent run faster. It turns out the same idea can also help recommender systems perform better.

# 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	?	?	○

[ 5 5 0 0 ?  
5 ? ? 0 ?  
? 4 0 ? ?  
0 0 5 4 ?  
0 0 5 0 ? ]

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

For the new user Eve, who hasn't had any rating, its parameters should be  $w^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ ,  $b^{(s)} = 0$ . Thus, the prediction  $w^{(s)} \cdot x^{(1)} + b^{(s)} = 0$ , which is not making too much sense.

$$w^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad b^{(s)} = 0 \quad w^{(s)} \cdot x^{(1)} + b^{(s)}$$

# Mean Normalization

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

$\rightarrow 2.5 = \frac{10}{4}$   
 $\rightarrow 2.5 = \frac{5}{2}$   
 $\rightarrow 2 = \frac{4}{2}$   
 $\rightarrow 2.25 = \frac{9}{4}$   
 $\rightarrow 1.25 = \frac{5}{4}$

$$\mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix}$$

$$\gamma - \mu = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

For user  $j$ , on movie  $i$  predict:

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

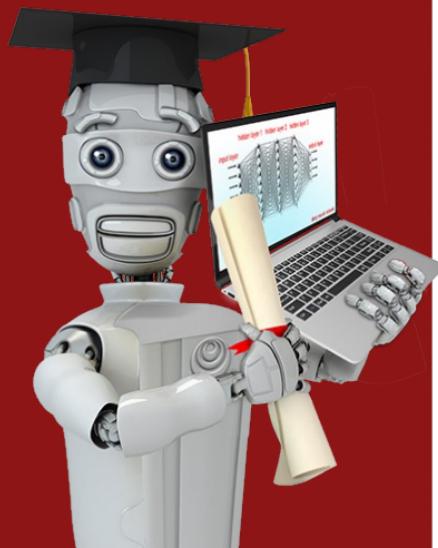
User 5 (Eve):

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

$$w^{(5)} \cdot x^{(1)} + b^{(5)} + \mu_1 = 2.5$$

$\underbrace{\phantom{w^{(5)} \cdot x^{(1)}}}_{0}$

Now, the prediction is 2.5 instead of 0.



# Recommender Systems implementational detail

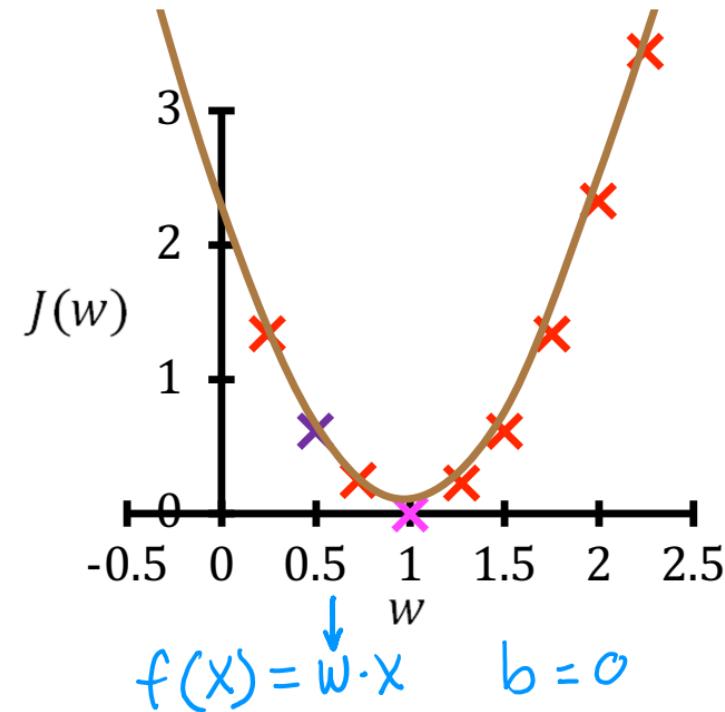
## TensorFlow implementation

# Derivatives in ML

Gradient descent algorithm

Repeat until convergence

$$\begin{aligned} w &= w - \alpha \frac{d}{dw} J(w, b) && \text{Learning rate} \\ b &= b - \alpha \frac{d}{db} J(w, b) \leftarrow b = 0 && \text{Derivative} \end{aligned}$$



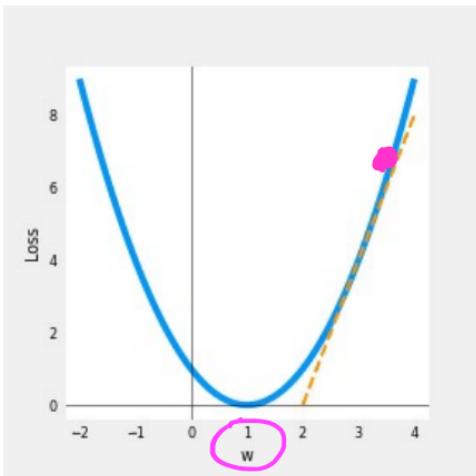
$$J = (wx - 1)^2$$

$f(x)$        $y$

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)
```

$$\frac{\partial}{\partial w} J(w)$$

tf.variables require special function to modify

# Implementation in TensorFlow

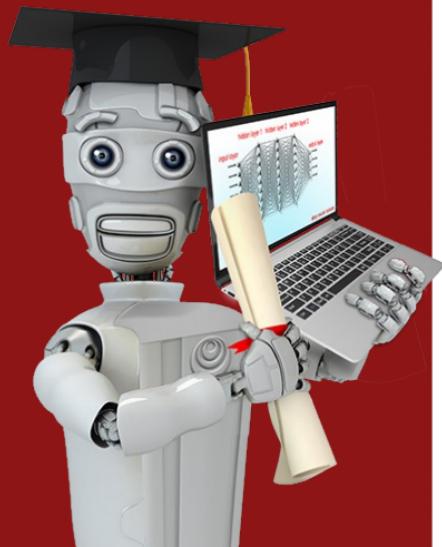
Gradient descent algorithm

Repeat until convergence

$$\begin{aligned} w &= w - \alpha \frac{\partial}{\partial w} J(w, b, X) \\ b &= b - \alpha \frac{\partial}{\partial b} J(w, b, X) \\ X &= X - \alpha \frac{\partial}{\partial X} J(w, b, X) \end{aligned}$$

```
# 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, n_u, n_m)  
        num_users, num_movies, lambda  
        y(i,j) r(i,j)  
        # 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



After training the model,  $X^{(1)}, X^{(2)}, \dots, X^{(n_m)}$  are known. But, we don't know what the meaning is for these features. E.g., we don't know if  $x_1$  is romance or  $x_2$  is action, etc. But we can still find similar items by calculating distances among features.

## Collaborative Filtering

# Finding related items

# Finding related items

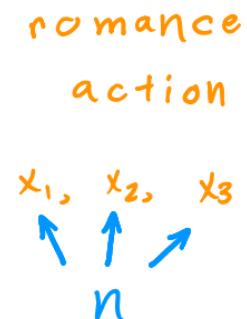
The features  $x^{(i)}$  of item  $i$  are quite hard to interpret.

To find other items related to it,

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

i.e. with smallest  
distance

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



# Limitations of Collaborative Filtering

## → Cold start problem. How to

- • rank new items that few users have rated?
- • show something reasonable to new users who have rated few items?

## → Use side information about items or users:

Collaborative filtering cannot use side information about items or users.

- • Item: Genre, movie stars, studio, ....
- • User: Demographics (age, gender, location), expressed preferences, ...

}

## Content-based Filtering

Collaborative filtering v.s. Content-based filtering

Collaborative filtering :

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

Content-based filtering:

Recommend items to you based on features of users and features of items to find good match.

$r(i, j) = 1$  if user  $j$  has rated movie  $i$  (0 otherwise)

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

Examples of user and item features:

User features:

- Age
- Gender
- Country
- Movies watched
- Average rating per genre
- ...

$X_u^{(j)}$  for user  $j$

vector size could be different.

Movie features:

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

$X_m^{(i)}$  for movie  $i$

Content-based filtering: Learning to match.

Previously, in collaborative filtering, Predict rating of user j on movie i as

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

But now, in content-based filtering,

$$w^{(j)} \cdot x^{(i)} + b^{(j)} \rightarrow v_u^{(j)} \cdot v_m^{(i)}$$

dot product

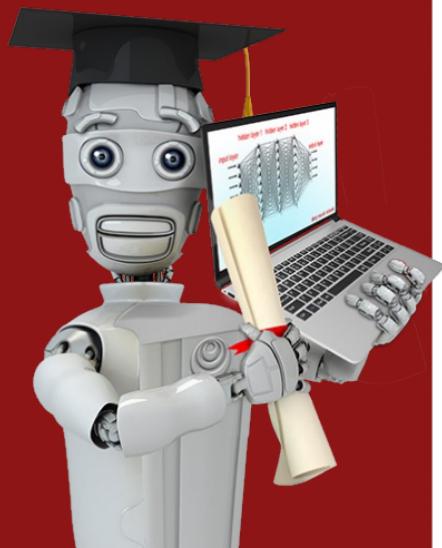
The diagram illustrates the calculation of the dot product. It shows three components:  $w^{(j)}$  (blue),  $x^{(i)}$  (green), and  $b^{(j)}$  (red). Arrows point from  $w^{(j)}$  to  $v_u^{(j)}$  and from  $x^{(i)}$  to  $v_m^{(i)}$ . A red 'X' is placed over  $b^{(j)}$ . An arrow points from the expression  $w^{(j)} \cdot x^{(i)} + b^{(j)}$  to the result  $v_u^{(j)} \cdot v_m^{(i)}$ , which is labeled as the 'dot product'.

$v_u^{(j)}$ : computed from  $x_u^{(i)}$

$v_m^{(i)}$ : computed from  $x_m^{(i)}$

Note that  $x_u^{(j)}$ ,  $x_m^{(i)}$  could be different size but  $v_u^{(j)}$  and  $v_m^{(i)}$  must be the same size.

The question is how to calculate  $v_u^{(j)}$  and  $v_m^{(i)}$ . Ans: Deep Learning

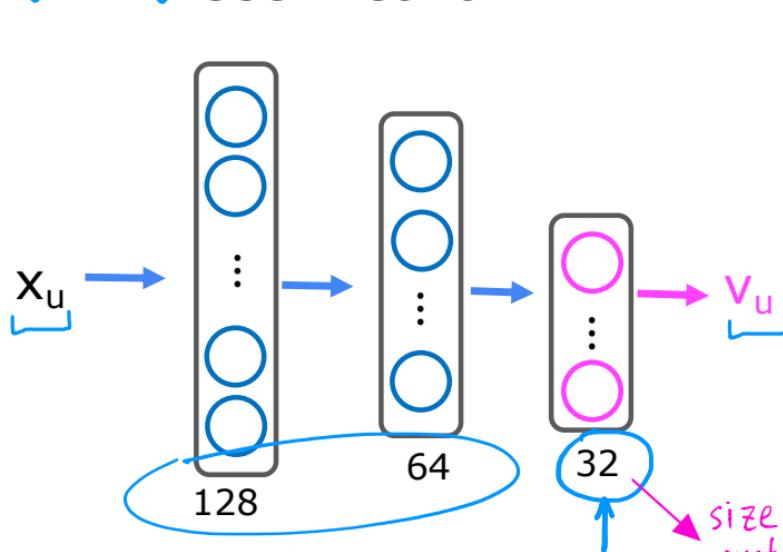


## Content-based Filtering

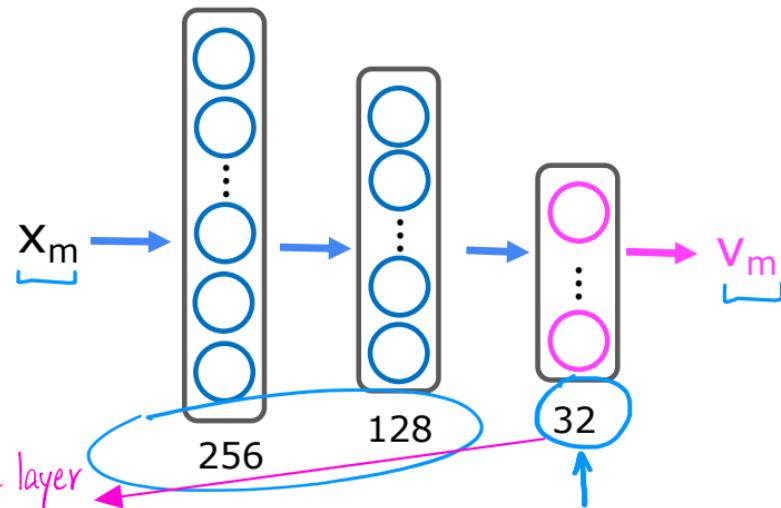
Deep learning for  
content-based filtering

# Neural network architecture

$x_u \rightarrow v_u$  User network



$x_m \rightarrow v_m$  Movie network



sigmoid function  $\mathbf{g}(\mathbf{v}_u^{(j)} \cdot \mathbf{v}_m^{(i)})$  to predict the probability that  $y^{(i,j)}$  is 1 (Only do it if output label is binary.)

# Neural network architecture

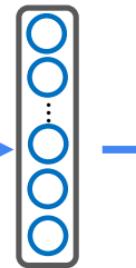
User network

$$x_u$$



Movie network

$$x_m$$



normalization

$$v_u$$

$$v_m$$

Train User network and Movie network altogether.

Prediction

Cost function

$$J =$$

$$\sum_{(i,j):r(i,j)=1} (v_u^{(j)} \cdot v_m^{(i)} - y^{(i,j)})^2 + \text{NN regularization term}$$

## 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  $x_m^{(i)}$

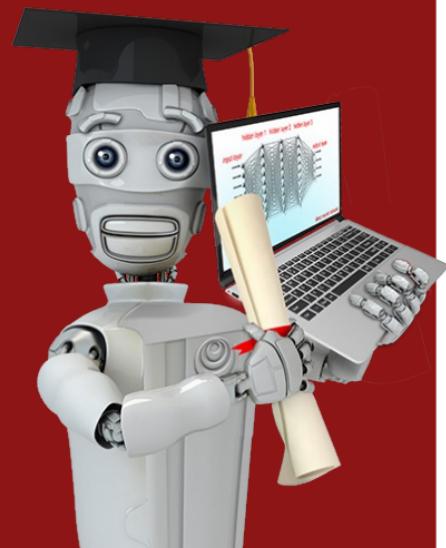
To find movies similar to movie i:

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

For each movie i, similar K movies can be figured out  
when features vectors  $v_m^{(i)}$  are learned.

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

- Note: This can be pre-computed ahead of time



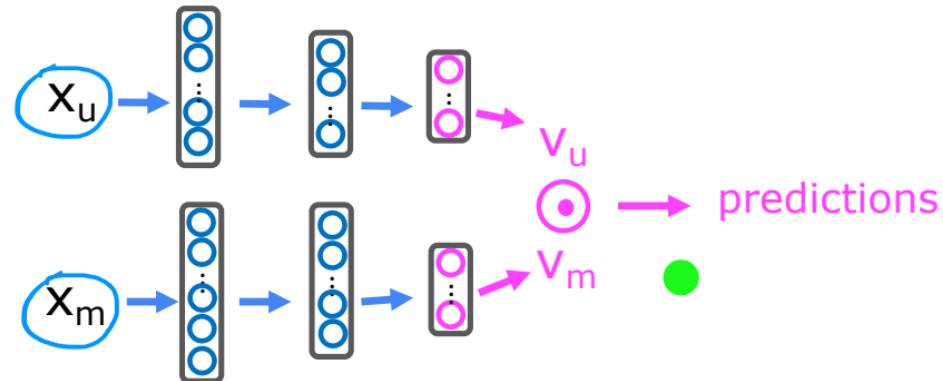
## Advanced implementation

Recommending from  
a large catalogue

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

The problem is that when a user Xu visits my website and I have to run thousands or millions of inferences to predict the best item, This is super inefficient.

- • Movies      1000+
- • Ads            1m+
- • Songs        10m+
- • Products    10m+



## Two steps: Retrieval & Ranking

Retrieval:

- • Generate large list of plausible item candidates  
e.g.
  - 1) For each of the last 10 movies watched by the user,  
find 10 most similar movies
  - 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

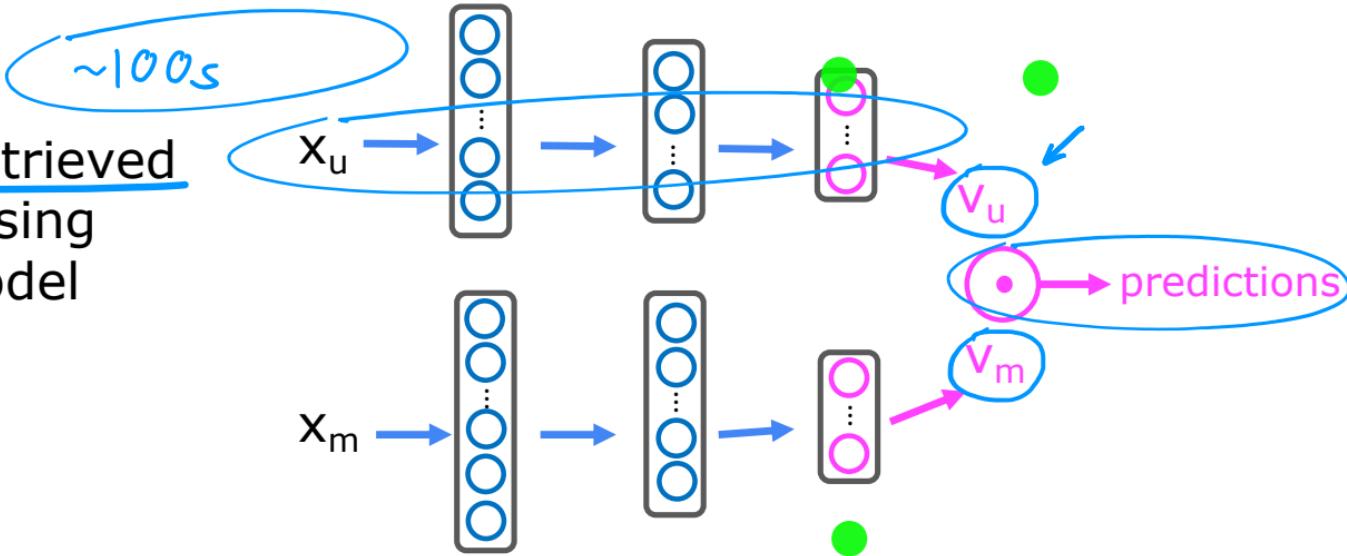
$$\| V_m^{(k)} - V_m^{(i)} \|^2$$

$10 \times 10$   
 $\sim 100s$

## Two steps: Retrieval & ranking

Ranking:

- Take list retrieved and rank using learned model

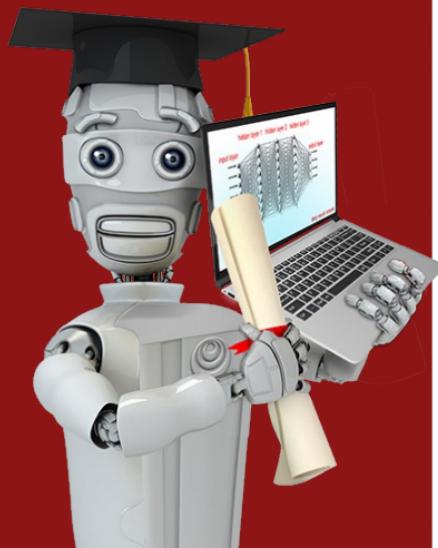


- Display ranked items to user

## Retrieval step

- Retrieving more items results in better performance, but slower recommendations.
- To analyse/optimize the trade-off, carry out offline experiments to see if retrieving additional items results in more relevant recommendations (i.e.,  $p(y^{(i,j)}) = 1$  of items displayed to user are higher).

100    500



## Advanced implementation

Ethical use of  
recommender systems

# 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

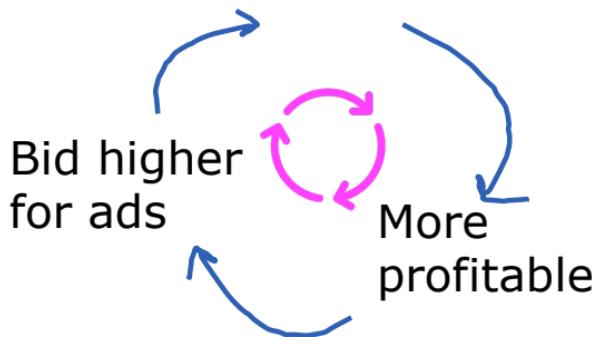


這就是我為什麼這麼喜歡 Andrew 的原因

# Ethical considerations with recommender systems

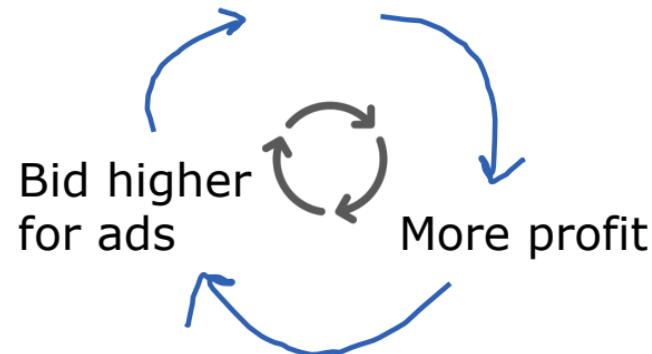
## Travel industry

Good travel experience  
to more users



## Payday loans

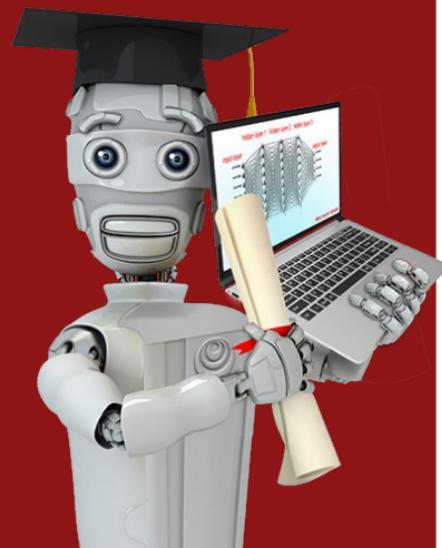
Squeeze customers more



Amelioration: Do not accept ads from exploitative businesses

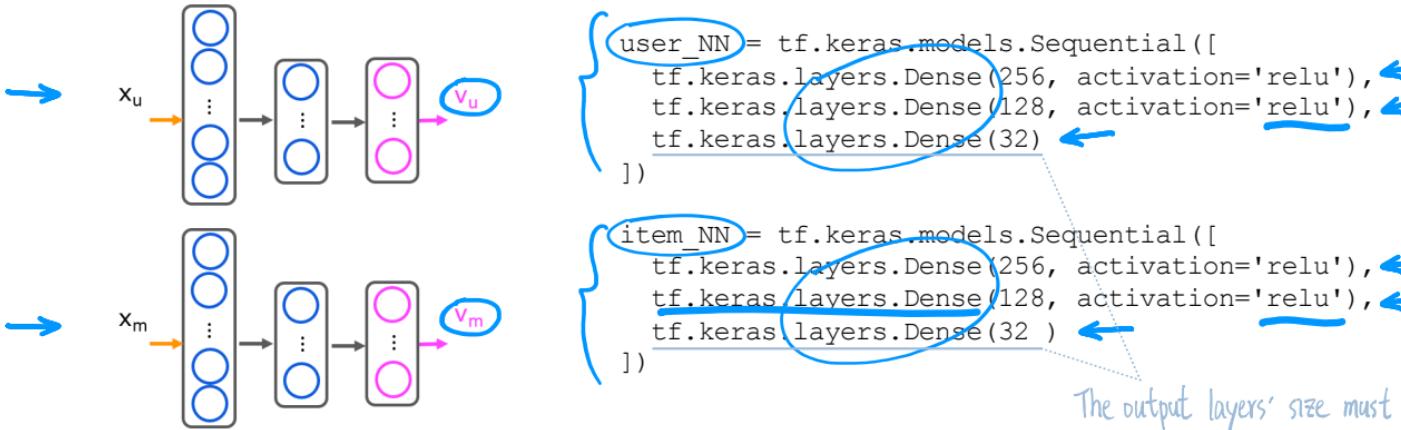
## Other problematic cases:

- • 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



# Content-based Filtering

## TensorFlow Implementation



```

# 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) # normalization

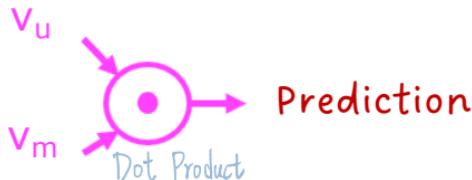
# create the item input and point to the base network
input_item = tf.keras.layers.Input(shape=(num_item_features))
vm = item_NN(input_item)
vm = tf.linalg.l2_normalize(vm, axis=1) # normalization

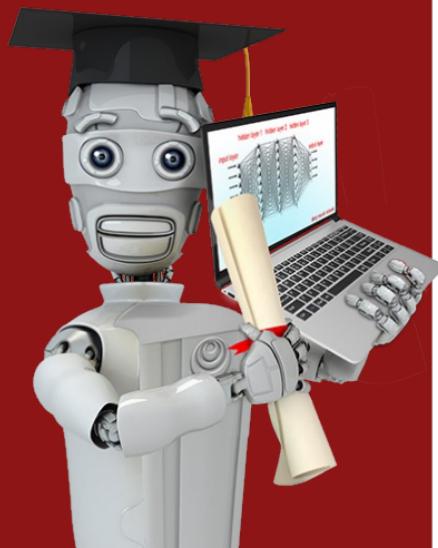
# measure the similarity of the two vector outputs
output = tf.keras.layers.Dot(axes=1)([vu, vm])

# 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()

```



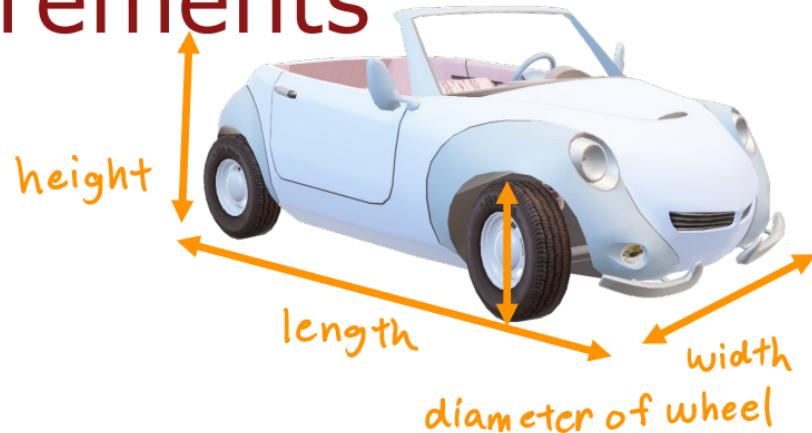


# Principal Component Analysis (Optional)

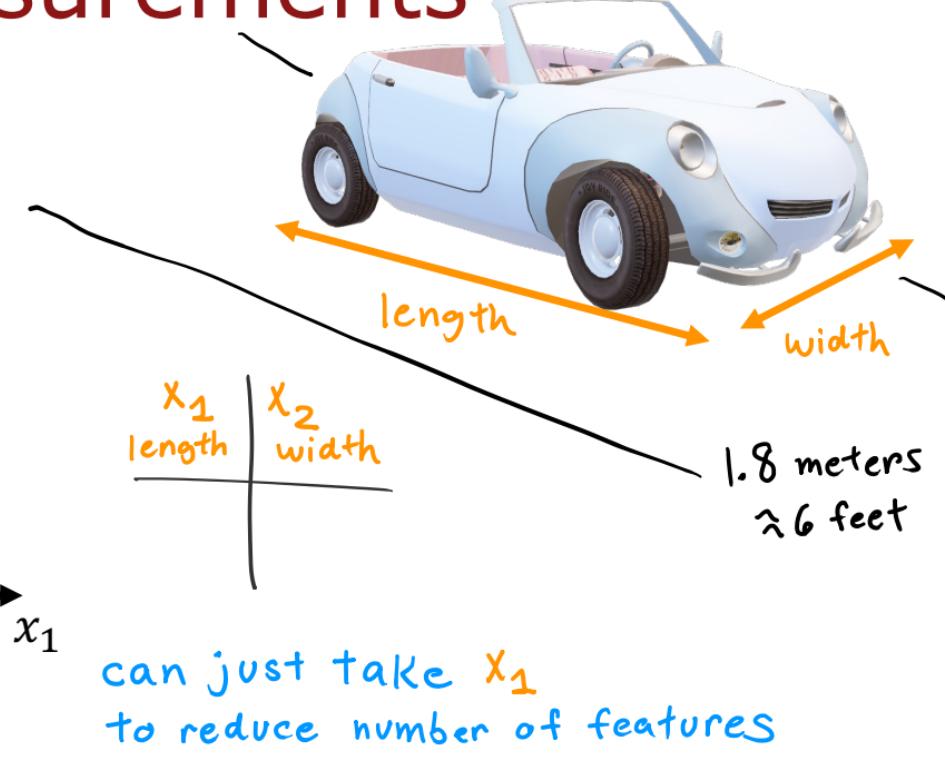
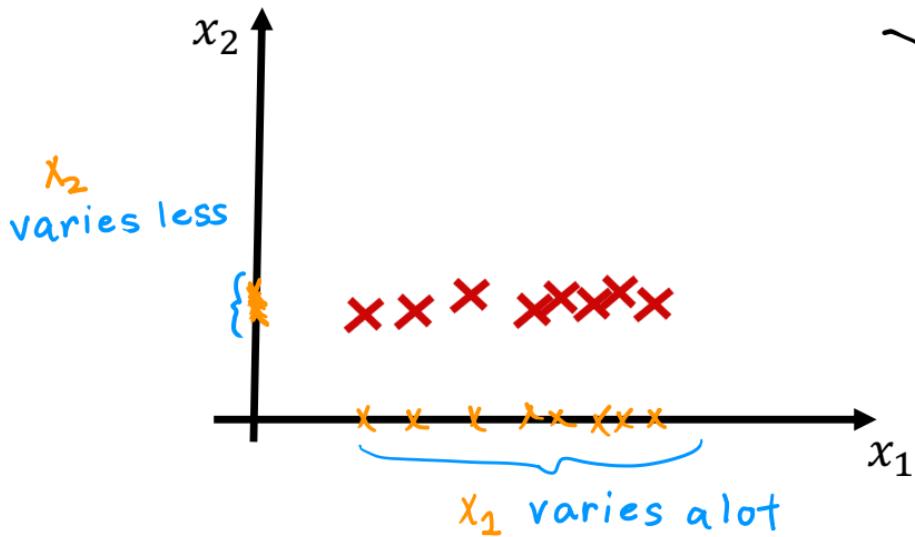
Reducing the number of features

# Car measurements

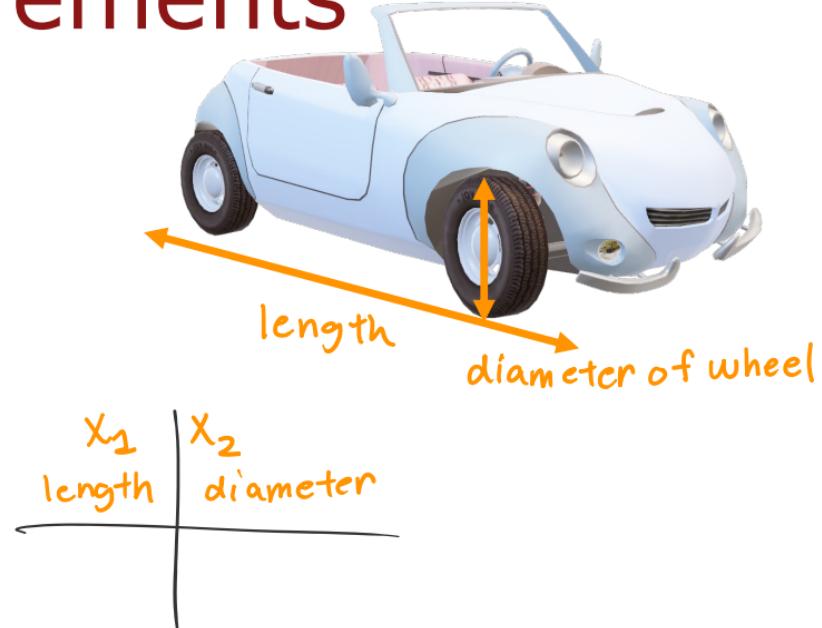
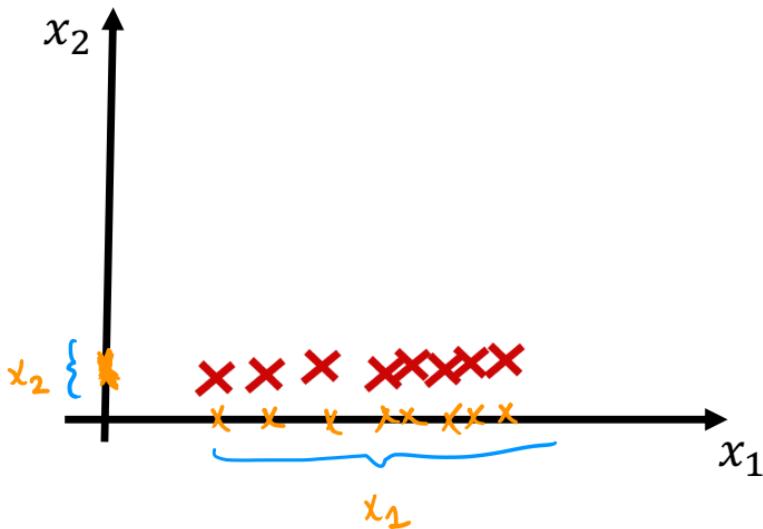
Given cars' features as shown in the right image



# Car measurements



# Car measurements



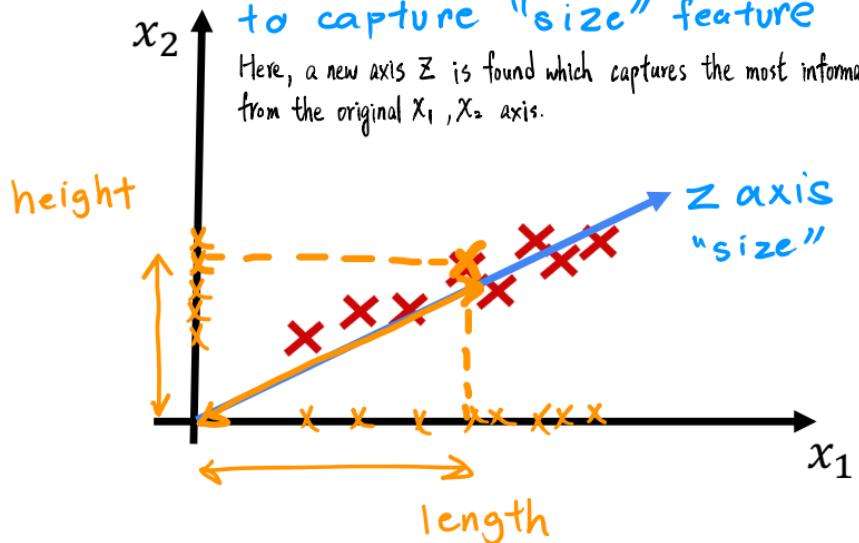
can just take  $x_1$   
to reduce number of features

# Size

PCA: find new axis and coordinates

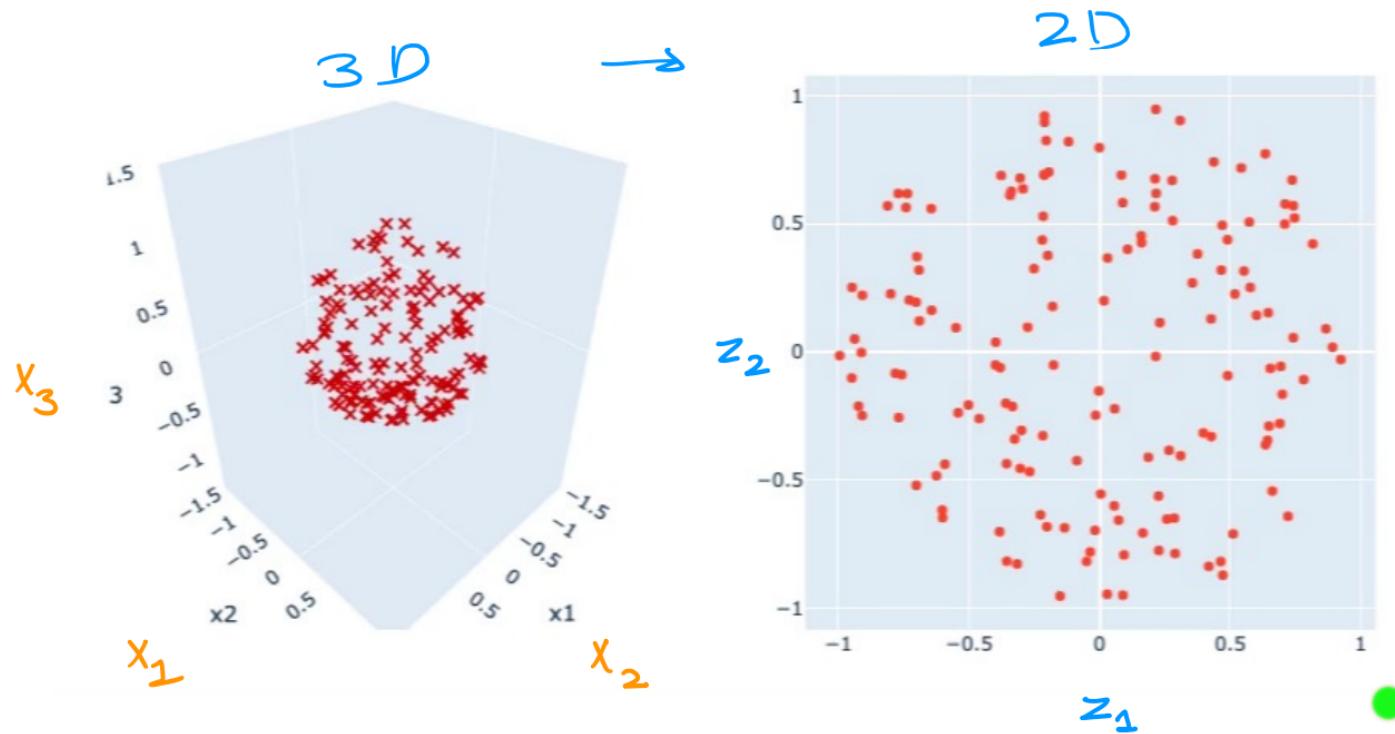
use fewer numbers  
to capture "size" feature

Here, a new axis  $Z$  is found which captures the most information from the original  $x_1, x_2$  axis.

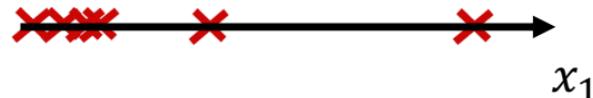


2 → 1 features  
many features → 2 or 3 features

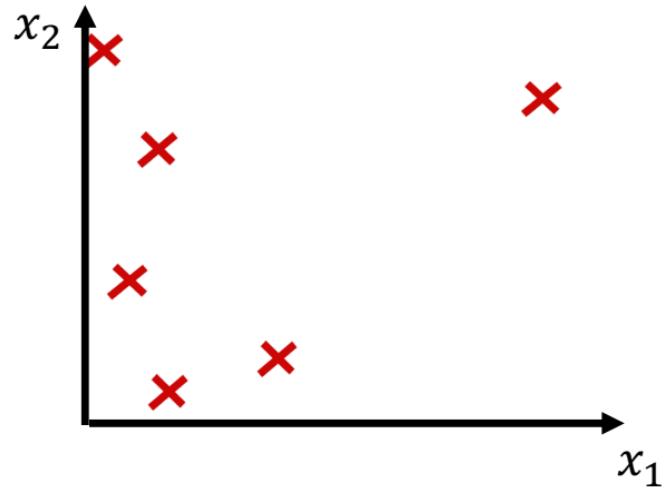
# From 3D to 2D



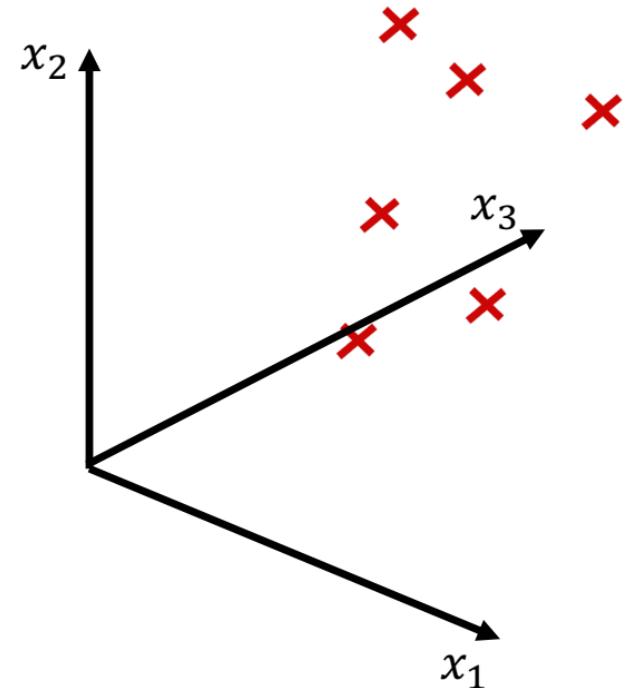
Country	GDP (trillions of US\$)
Canada	1.577
China	5.878
India	1.632
Russia	1.48
Singapore	0.223
USA	14.527
...	...



Country	GDP (trillions of US\$)	Per capita GDP (thousands of intl. \$)
Canada	1.577	39.17
China	5.878	7.54
India	1.632	3.41
Russia	1.48	19.84
Singapore	0.223	56.69
USA	14.527	46.86

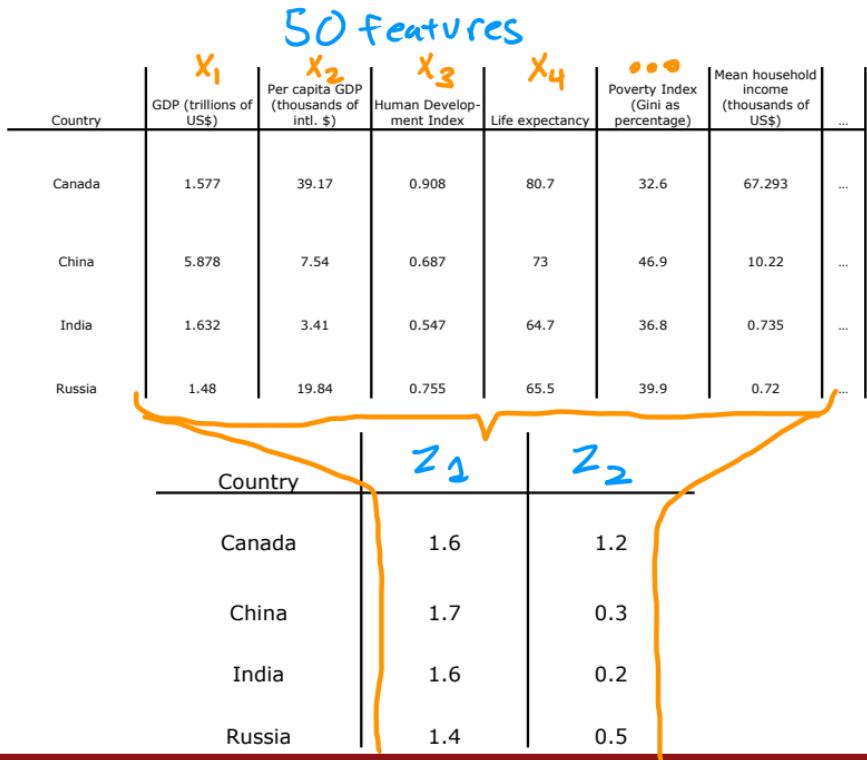


Country	GDP (trillions of US\$)	Per capita GDP (thousands of intl. \$)	Human Development Index
Canada	1.577	39.17	0.908
China	5.878	7.54	0.687
India	1.632	3.41	0.547
Russia	1.48	19.84	0.755
Singapore	0.223	56.69	0.866
USA	14.527	46.86	0.91
...	...	...	...

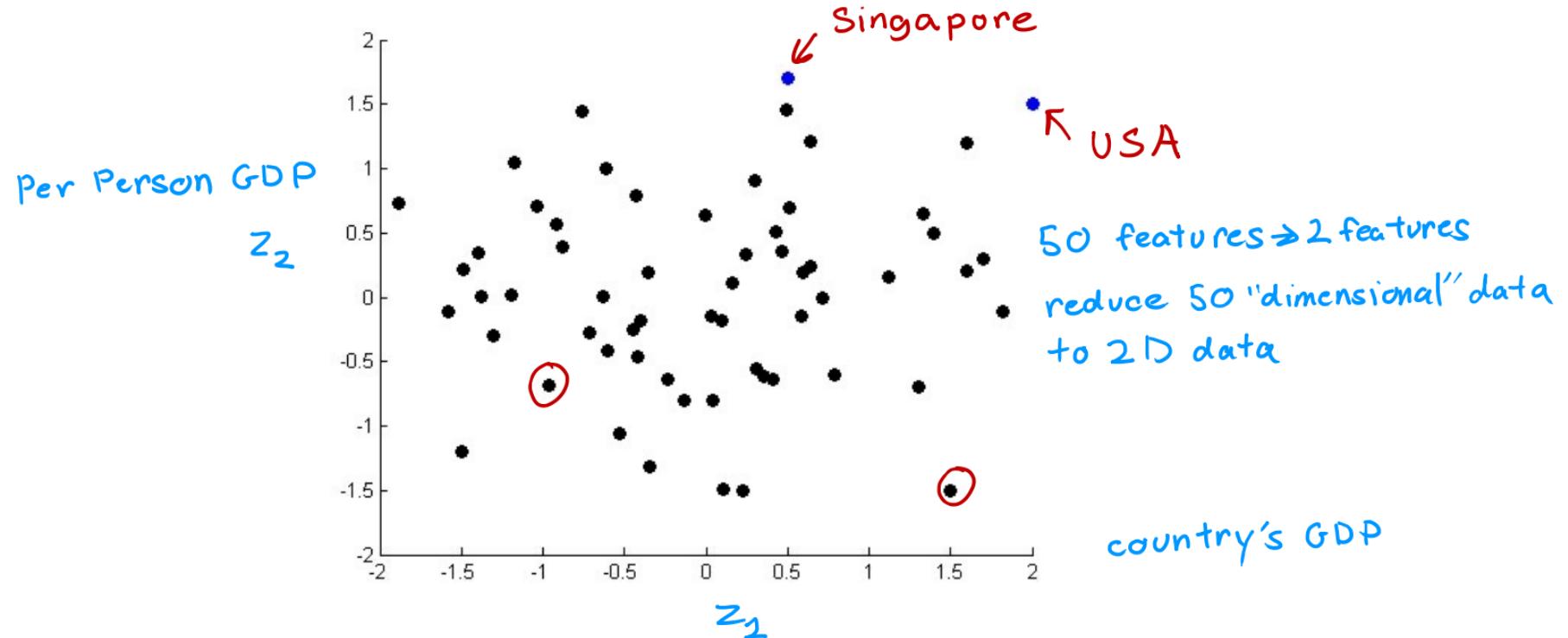


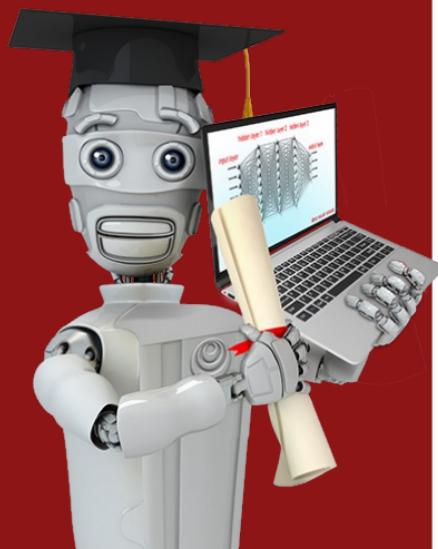
Country	GDP (trillions of US\$)	Per capita GDP (thousands of intl. \$)	Human Development Index	Life expectancy
Canada	1.577	39.17	0.908	80.7
China	5.878	7.54	0.687	73
India	1.632	3.41	0.547	64.7
Russia	1.48	19.84	0.755	65.5
Singapore	0.223	56.69	0.866	80
USA	14.527	46.86	0.91	78.3
...	...	...	...	...

what if 50 features?



# Data visualization





# Principal Component Analysis

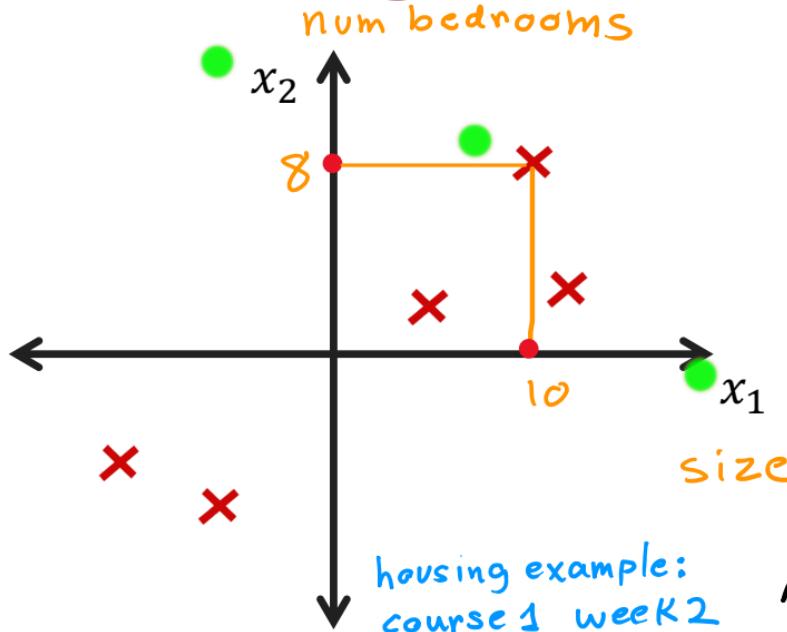
## PCA Algorithm

# PCA algorithm

coordinates

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

Can we choose  
a different axis?



Preprocess features

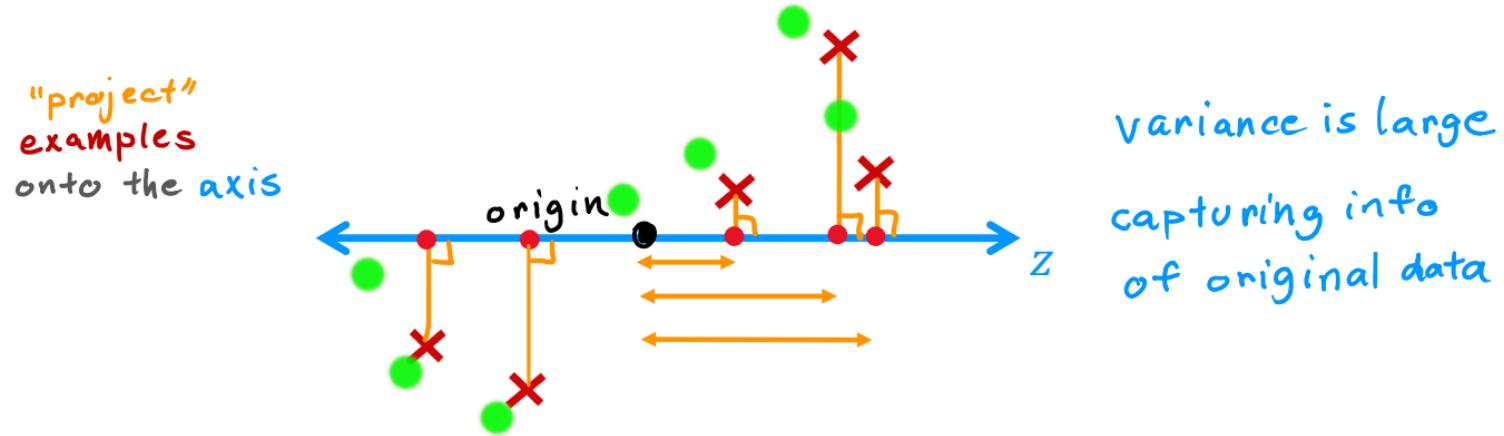
Normalized to  
have zero mean



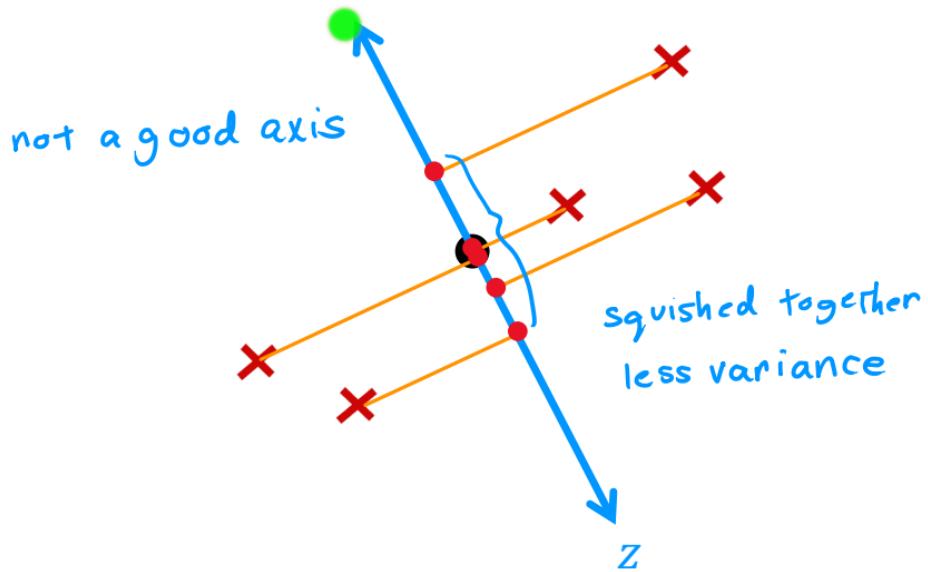
feature scaling



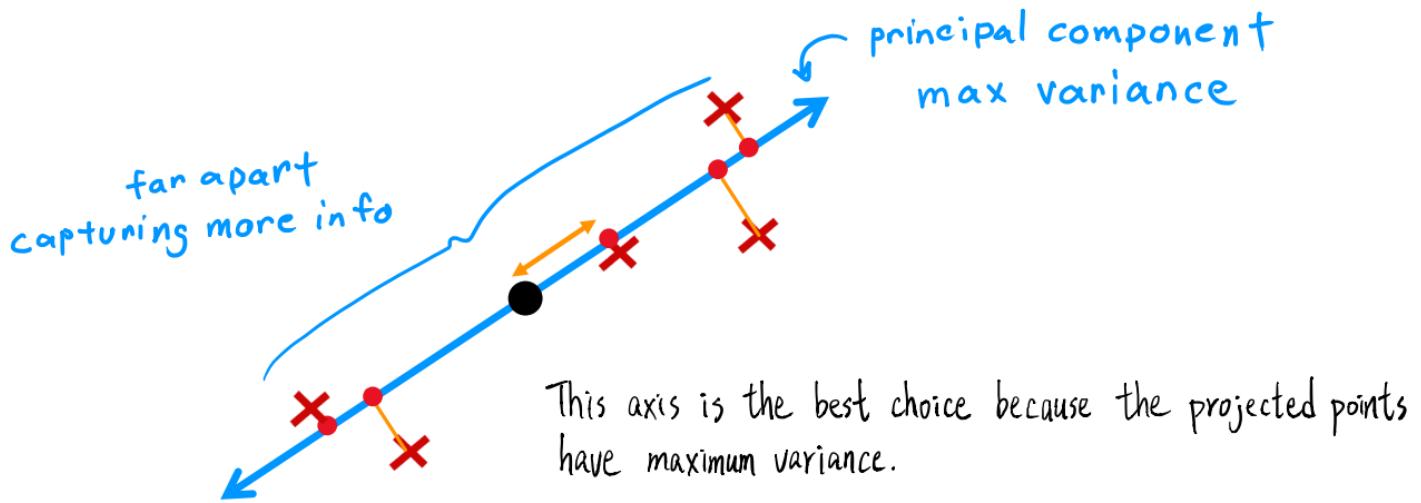
# Choose an axis



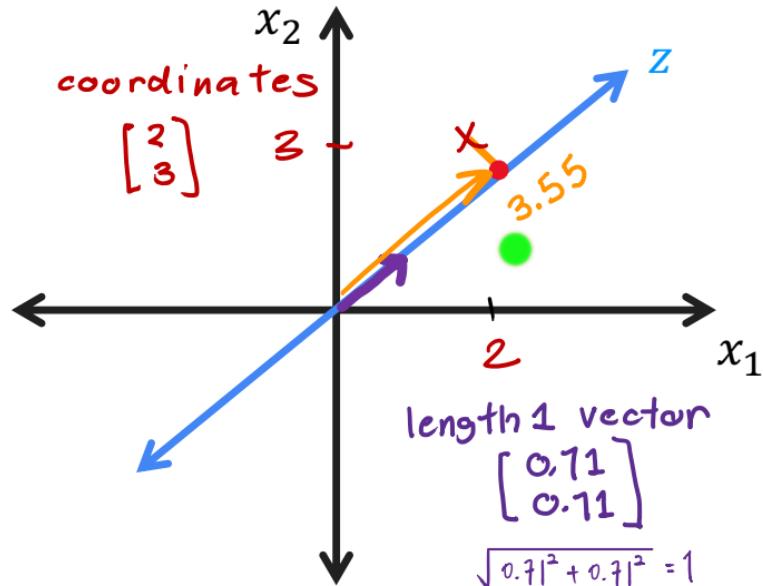
# Choose an axis



# Choose an axis



# Coordinate on the new axis

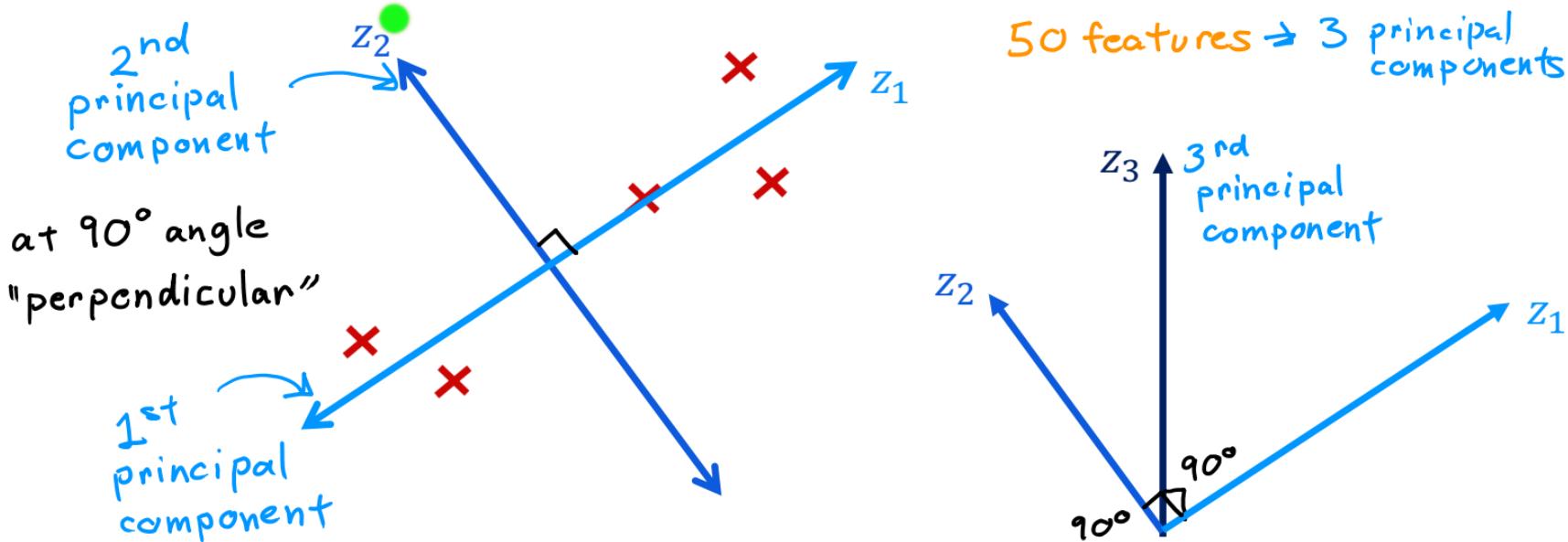


dot product

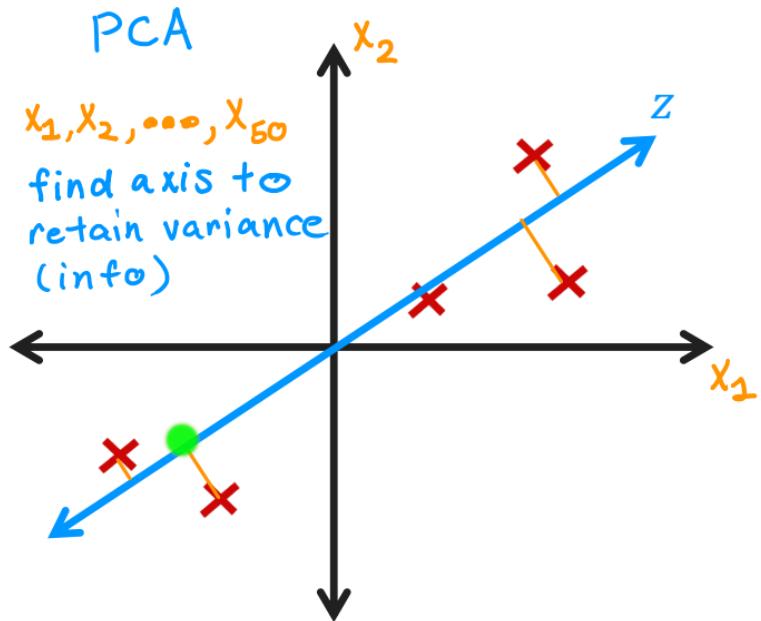
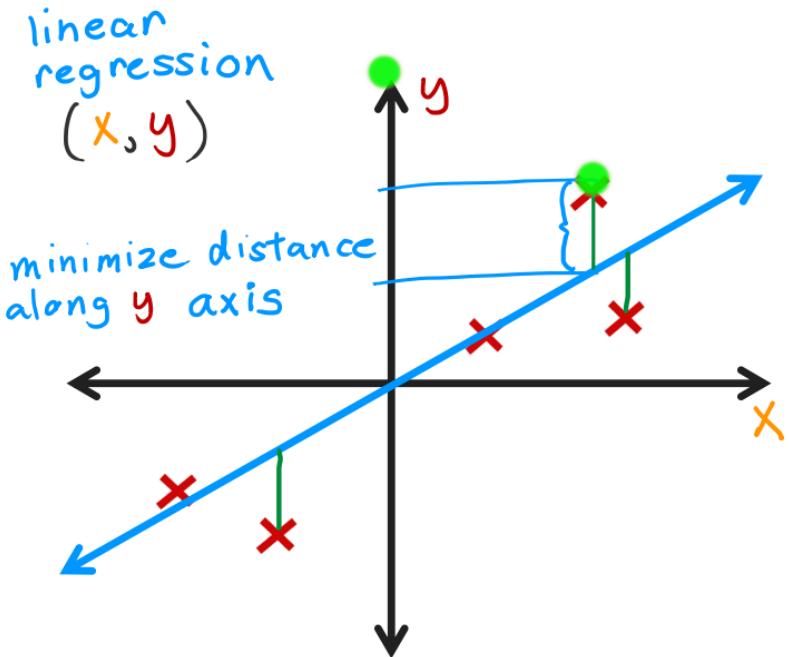
$$\begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix}$$

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

# More principal components

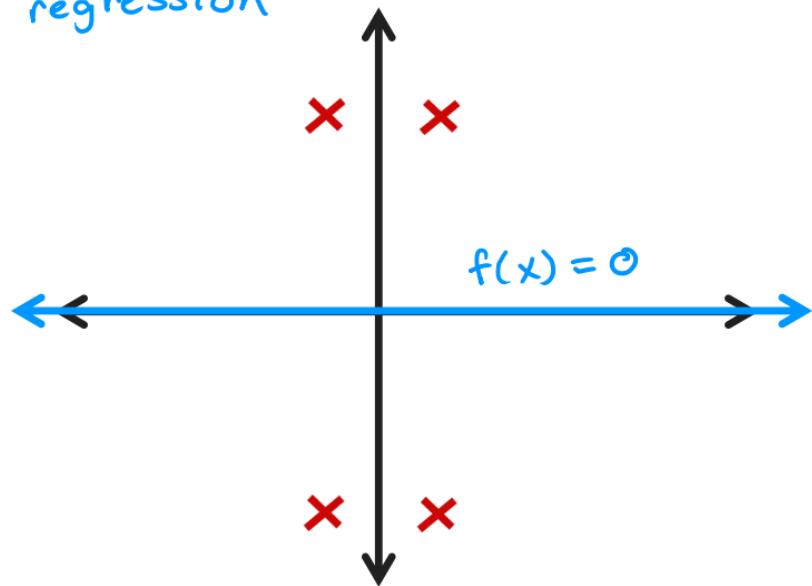


# PCA is not linear regression

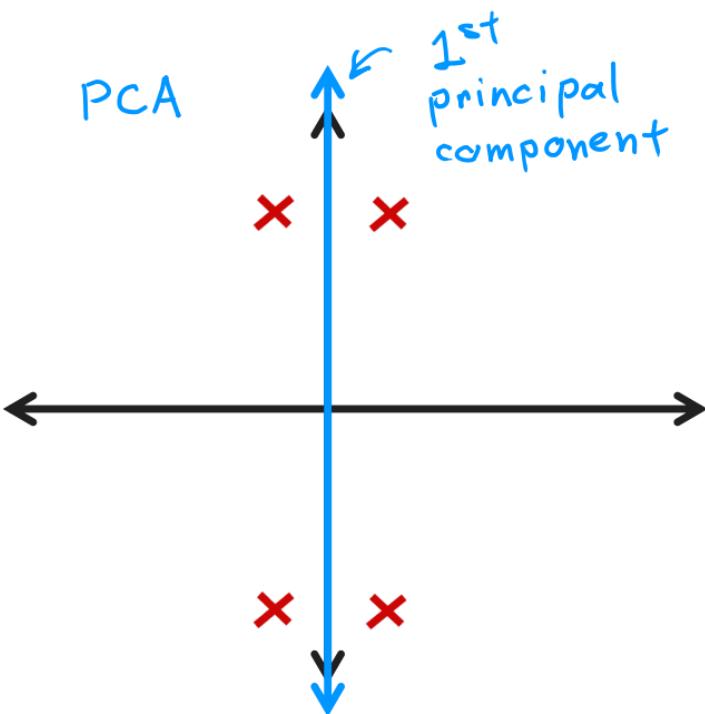


# PCA is not linear regression

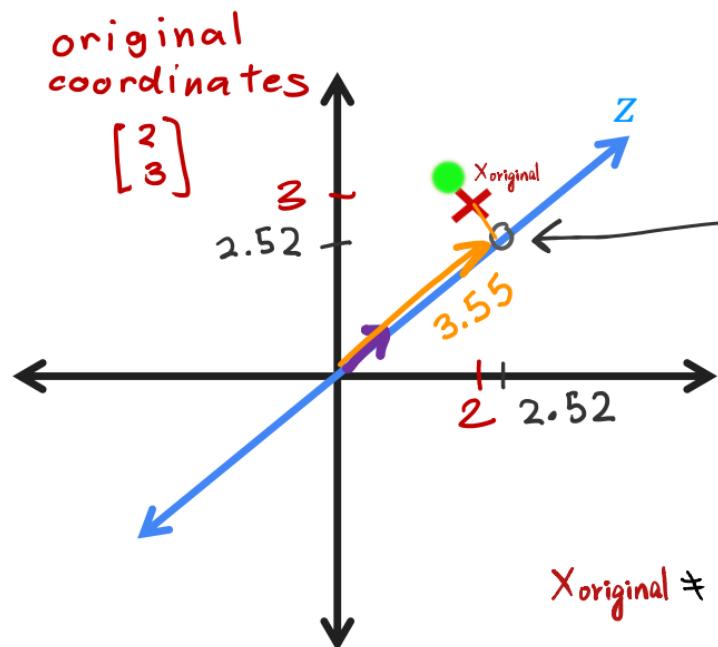
linear  
regression



PCA



# Approximation to the original data

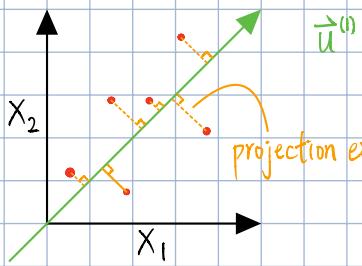


given  $z = 3.55$ ,  
find original  $(x_1, x_2)$  (approximately)  
"reconstruction"

$$3.55 \times \begin{bmatrix} 0.71 \\ 0.71 \end{bmatrix} = \begin{bmatrix} 2.52 \\ 2.52 \end{bmatrix}$$

$x_{\text{approximate}}$

# PCA Algorithm



Given training set  $X^{(1)}, X^{(2)}, \dots, X^{(m)}$ ,  $X^{(i)} \in \mathbb{R}^n$

The basic idea is to find a direction (a vector  $U^{(1)} \in \mathbb{R}^n$ ) onto which to project the data so as to minimize the **projection error**.

Then, we can find  $k$  principal components  $U^{(1)}, U^{(2)}, \dots, U^{(k)} \in \mathbb{R}^n$

Data Preprocessing:

Calculate mean:  $\mu_j = \frac{1}{m} \sum_{i=1}^m X_j^{(i)}$ ,  $j = 1, \dots, n$

Move toward mean: replace each  $X_j^{(i)}$  with  $X_j^{(i)} - \mu_j$

Rescale features: E.g.  $X_1$  = size of house,  $X_2$  = number of bedrooms. We can rescale  $X_1$  and  $X_2$  to have comparable range of values.

Compute "covariance matrix".  $\Sigma = \frac{1}{m} \sum_{i=1}^m X^{(i)} X^{(i)T}$ .  $\Sigma \in \mathbb{R}^{n \times n}$

Compute "eigen vectors" of  $\Sigma$ :  $u, s, v = \text{svd}(\Sigma)$

$$U = \begin{bmatrix} | & | & | \\ U^{(1)}, U^{(2)}, \dots, U^{(n)} \\ | & | & | \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad U_{\text{reduce}} = \begin{bmatrix} | & | & | \\ U^{(1)}, U^{(2)}, \dots, U^{(k)} \\ | & | & | \end{bmatrix} \in \mathbb{R}^{n \times k}, \quad Z^{(i)} = U_{\text{reduce}}^T X^{(i)} \in \mathbb{R}^{k \times 1}$$

Reconstruction:  $\hat{X}_{\text{approximate}}^{(i)} = U_{\text{reduce}} \cdot Z^{(i)}$ .  $\hat{X}_{\text{approximate}}^{(i)} \in \mathbb{R}^n$

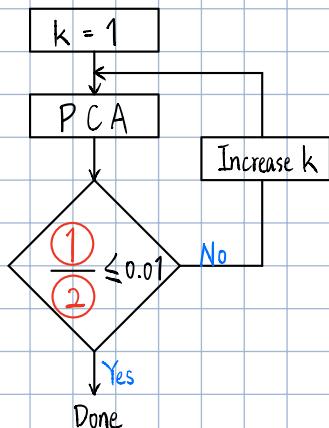
Now, the question is how to determine  $k$ ? I.e., how many principal components are enough to capture most information?

Calculate average squared projection error:  $\frac{1}{m} \sum_{i=1}^m \|X^{(i)} - \hat{X}_{\text{approximate}}^{(i)}\|^2$  —①

$$S = \begin{bmatrix} S_{11} & & 0 \\ S_{21} & \ddots & \\ 0 & & S_{nn} \end{bmatrix}$$

Calculate total variation in the data:  $\frac{1}{m} \sum_{i=1}^m X^{(i)} X^{(i)T}$  —②

Run:



$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01$$

$$\Rightarrow \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99 \quad \text{99% of variance is retained.}$$



Stanford  
ONLINE



PCA

PCA in Code

# 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*

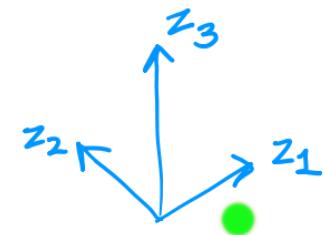


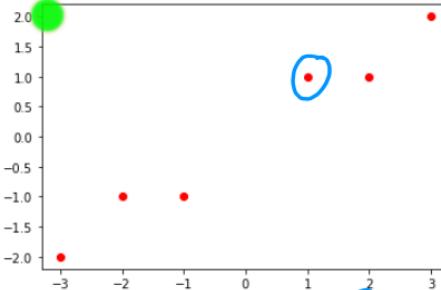
2. Optionally examine how much variance is explained by each principal component.  
*info*

`explained_variance_ratio` *E.g., if it's 0.992, it means the 1st principal component captures 99.2 % of information.*

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

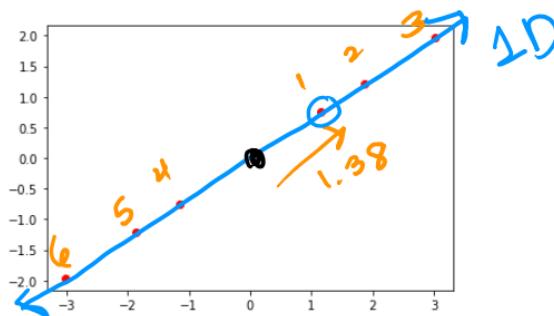
pca\_1 = PCA(n\_components=1) *k=1 Only 1 principle components*

pca\_1.fit(X)

pca\_1.explained\_variance\_ratio\_ *0.992*

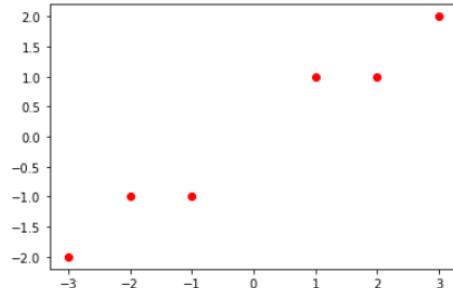
X\_trans\_1 = pca\_1.transform(X)

X\_reduced\_1 = pca.inverse\_transform(X\_trans\_1)



1D

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



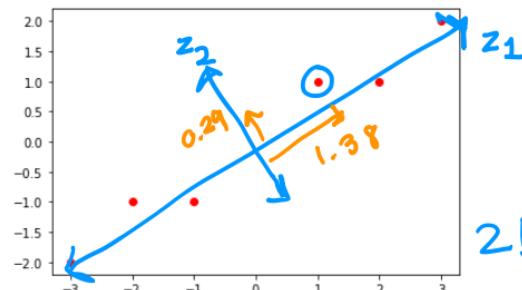
# Example

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

2D

```
pca_2 = PCA(n_components=2) # k=2 Only 2 principle components
pca_2.fit(X)
pca_2.explained_variance_ratio_
X_trans_2 = pca.transform(X)
X_reduced_2 = pca.inverse_transform(X_trans_2)
```

$z_1$      $z_2$   
0.992    0.008



2D

$z_1$                    $z_2$

```
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]])
```

$0.992 + 0.002 = 1$ , which means  $z_1$  and  $z_2$  together explain 100% of data, which makes sense because  $X$  is only 2-D.

# Applications of PCA

Note: Don't use PCA to prevent overfitting. Valuable information in data may be deprived due to PCA.

❖ Visualization *reduce to 2 or 3 features*

Less frequently used for:

- Data compression Andrew says it's not really often used nowadays because hardware improves a lot.  
(to reduce storage or transmission costs)  $50 \rightarrow 10$
- Speeding up training of a supervised learning model

$$n = 1000 \rightarrow 100$$

The idea is by reducing feature's dimension, the dataset becomes smaller, which may make older generation of algorithms, such as Support Vector Machine, run faster. But SVM is not as good as neural networks, and neural networks work well even with lots of features. Thus, it's not necessary to reduce dimension of features.