

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	?

$n_u = 4$ $r(1,1) = 1$ $y^{(1,1)} = 1$
 $n_m = 5$ $r(3,1) = 0$ $y^{(3,1)} = 0$

$n_u = \text{no. of users}$
 $n_m = \text{no. of movies}$
 $r(i,j) = 1$ if user j has rated movie i
 $y^{(i,j)}$ = rating given by user j to movie i (defined only if $r(i,j)=1$)

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



What if we have features of the movies?

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

$n_u = 4$
 $n_m = 5$
 $n = 2$

$x^{(1)} = \begin{bmatrix} 0.9 \\ 0 \end{bmatrix}$
 $x^{(3)} = \begin{bmatrix} 0.99 \\ 0 \end{bmatrix}$

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

$$w^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix} \quad b^{(1)} = 0 \quad x^{(1)} = \begin{bmatrix} 0.9 \\ 0 \end{bmatrix} \quad w^{(1)} \cdot x^{(1)} + b^{(1)} = 4.95$$

→ For user j : Predict user j 's rating for 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 $w^{(j)}, b^{(j)}$

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

↑ predicting user j's ratings for
other movies.

Cost function

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

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

↑ learning all the parameters for
all of the users.

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	?	$x^{(2)}$
Cute puppies of love	?	4	0	?	?	$x^{(3)}$
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 $w^{(j)} \cdot x^{(i)} + b^{(j)}$

$$\left. \begin{array}{l} w^{(1)} \cdot x^{(1)} \approx 5 \\ w^{(2)} \cdot x^{(1)} \approx 5 \\ w^{(3)} \cdot x^{(1)} \approx 0 \\ w^{(4)} \cdot x^{(1)} \approx 0 \end{array} \right\} \rightarrow x^{(1)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

4.46 / 13:55



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

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

that will give us our
collaborative filtering algorithm.

7:54 / 13:55



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:

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

as a function of w , b , and x ?

Gradient Descent

collaborative filtering

Linear regression (course 1)

repeat {

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

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

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

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

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

}

parameters w, b, x

\times is also a parameter

w and b, as well as x.

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	? ←

we can then decide how much we should recommend these items to them.



From regression to binary classification

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

1 that is of the user having engaged with or like the item using this model.



Cost function for binary application

Previous cost function:

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

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

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

cost for all examples



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	?	?

Matrix representation of user ratings:

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

An arrow points from the matrix to the cost function equation below.

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

sustained and more compact way.
 $w^{(s)} \cdot x^{(i)} + b^{(s)}$



Mean Normalization

$$\begin{array}{c}
 \xrightarrow{\quad} \begin{bmatrix} 5 & 5 & 0 & 0 & ? & 2.5 \\ 5 & ? & ? & 0 & ? & 2.5 \\ ? & 4 & 0 & ? & ? & 2 \\ 0 & 0 & 5 & 4 & ? & 2.25 \\ 0 & 0 & 5 & 0 & ? & 1.25 \end{bmatrix} \\
 \xrightarrow{\quad} \mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \\
 \xrightarrow{\quad} \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? & ? \\ 2.5 & ? & ? & -2.5 & ? & ? \\ ? & 2 & -2 & ? & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? & ? \end{bmatrix} \\
 \uparrow \uparrow \uparrow \uparrow \uparrow \quad \uparrow \uparrow \uparrow \quad \downarrow y^{(i,j)} \\
 \text{For user } j, \text{ on movie } i \text{ predict:} \\
 w^{(j)} \cdot x^{(i)} + b^{(j)} + \mu_i
 \end{array}$$

User 5 (Eve):

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

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

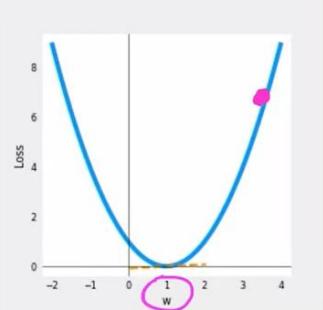
8.09 / 8.45

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

Gradient descent algorithm
Repeat until convergence

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

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

Tf.variables are the parameters we want to optimize

Auto Diff
Auto Grad

```

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

```

# 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)

    # 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



Finding related items

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

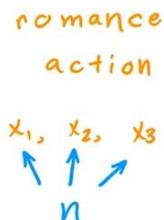
To find other items related to it,

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.



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:

- Item: Genre, movie stars, studio,
- User: Demographics (age, gender, location), expressed preferences, ...
you guess what might the user be interested in,

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)

content-based filtering is that

Examples of user and item features

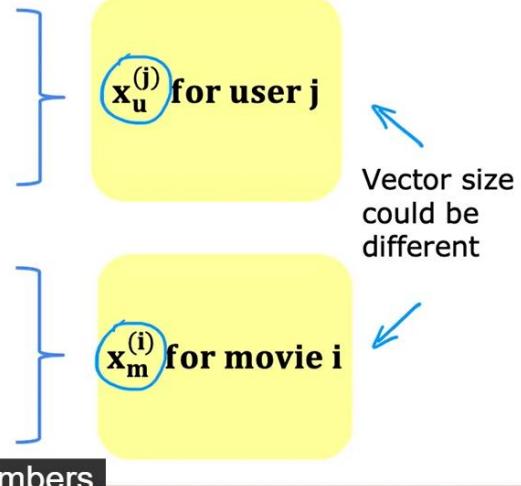
User features:

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

Movie features:

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

could be just 50 numbers.



Content-based filtering: Learning to match

Predict rating of user j on movie i as

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

computed from $x_u^{(j)}$

computed from $x_m^{(i)}$

user's preferences $v_u^{(j)} = \begin{bmatrix} 4.9 \\ 0.1 \\ \vdots \\ 3.0 \end{bmatrix}$ likes likes

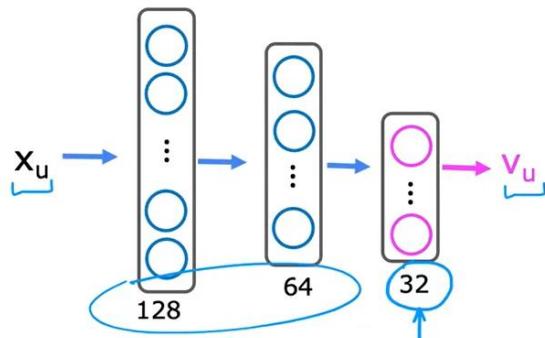
movie features $v_m^{(i)} = \begin{bmatrix} 4.5 \\ 0.2 \\ \vdots \\ 3.5 \end{bmatrix}$ romance action

maybe both of these are say 32 numbers.

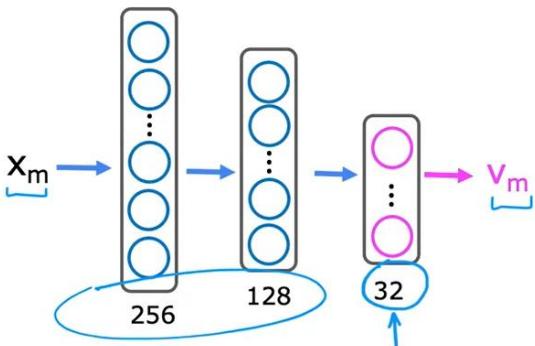


Neural network architecture

$x_u \rightarrow v_u$ User network



$x_m \rightarrow v_m$ Movie network



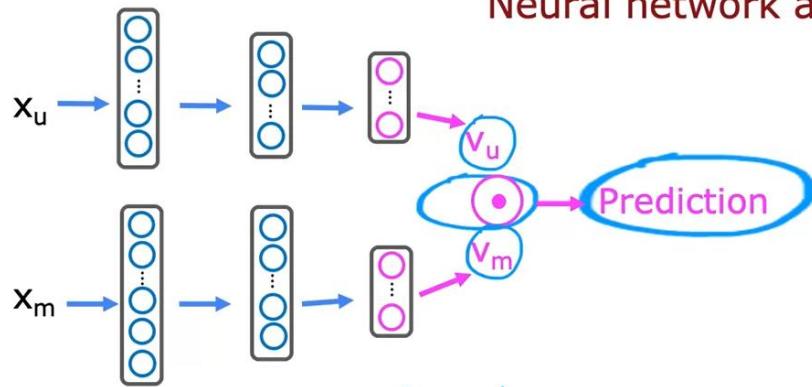
Prediction : $v_u^{(j)} \cdot v_m^{(i)}$

$g(v_u^{(j)} \cdot v_m^{(i)})$ to predict the probability that $y^{(i,j)}$ is 1

draw them together in



Neural network architecture



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

as well. Let's take a look.



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 :

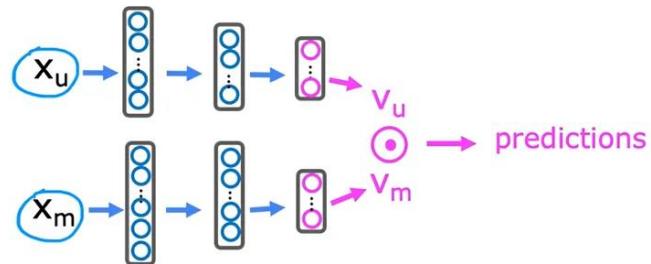
$$\|v_m^{(k)} - v_m^{(i)}\|^2 \text{ small}$$
$$\|x^{(k)} - x^{(i)}\|^2$$

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



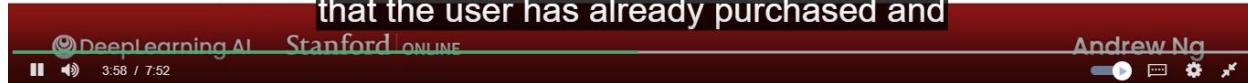
becomes computationally infeasible.



Two steps: Retrieval & Ranking

Retrieval:

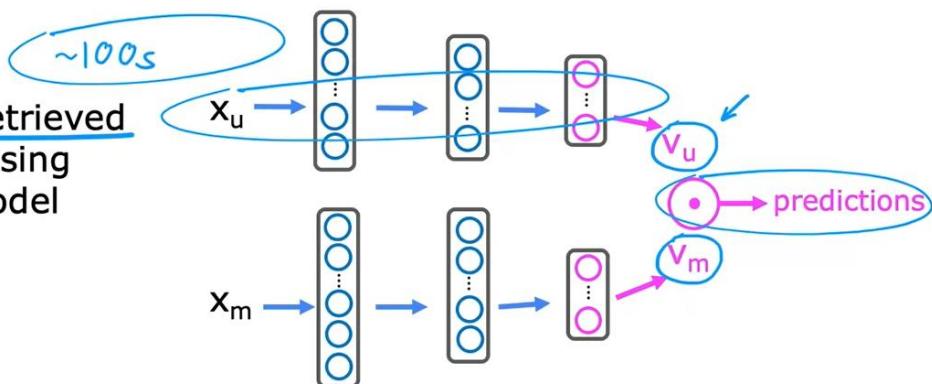
- • Generate large list of plausible item candidates $\sim 100s$
 - e.g.
 - 1) For each of the last 10 movies watched by the user, find 10 most similar movies
 - $$\| V_m^{(k)} - V_m^{(i)} \|^2$$
 - 2) For most viewed 3 genres, find the top 10 movies
 - 3) Top 20 movies in the country
- • Combine retrieved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



Two steps: Retrieval & ranking

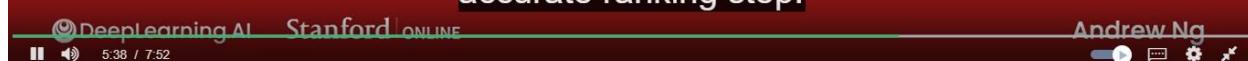
Ranking:

- Take list retrieved and rank using learned model



- Display ranked items to user

To feed into the more accurate ranking step.



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

not worth doing the more detailed
influence and inner product on.



What is the goal of the recommender system?

Recommend:

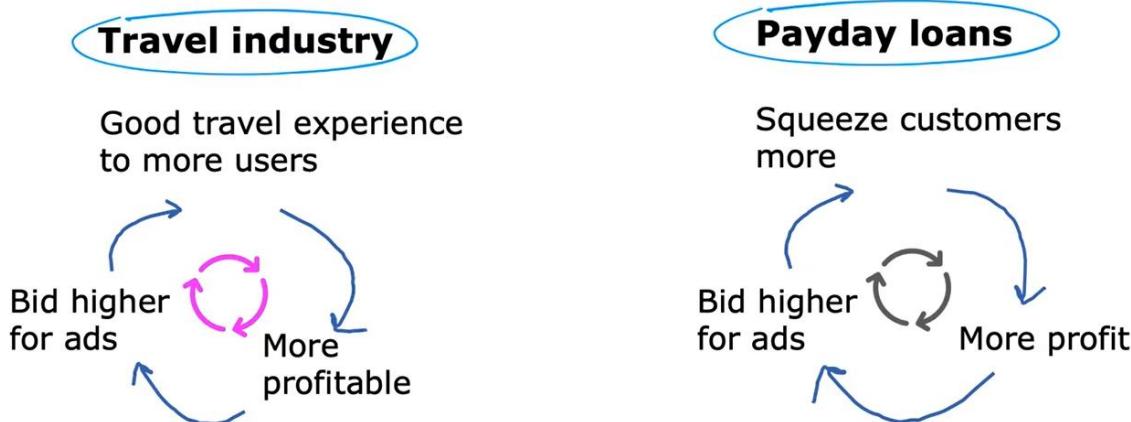
- • Movies most likely to be rated 5 stars by user
- • Products most likely to be purchased
- • Ads most likely to be clicked on *+high bid*
- • Products generating the largest profit
- • Video leading to maximum watch time



Or they could also be



Ethical considerations with recommender systems



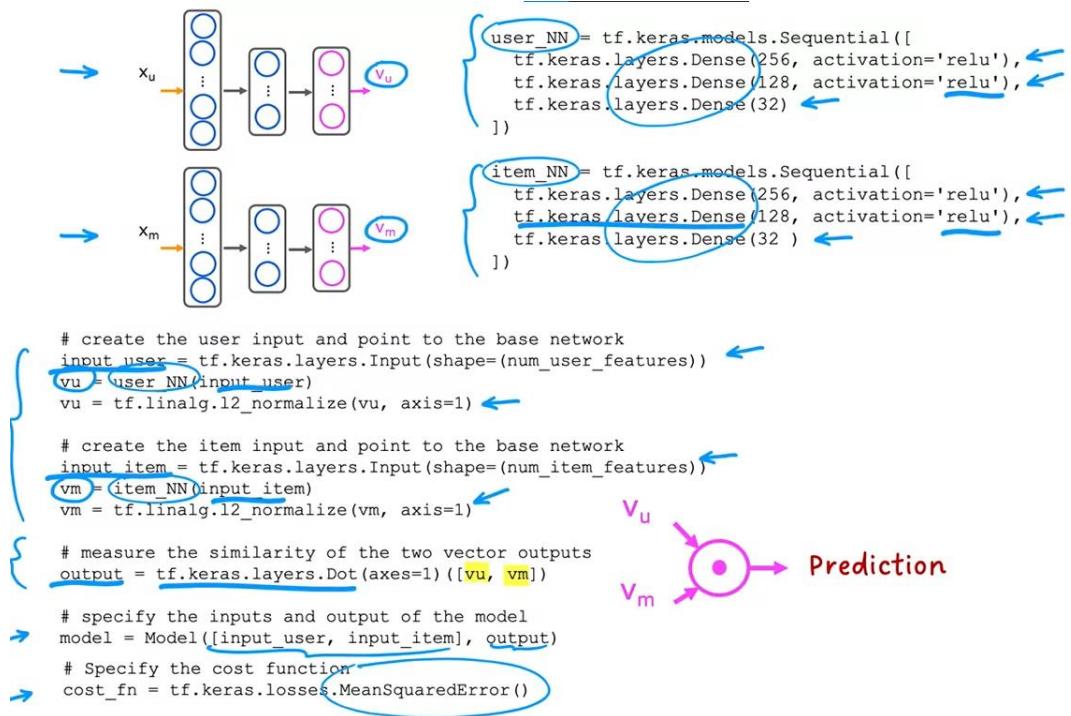
Amelioration: Do not accept ads from exploitative businesses

But how do you define what is

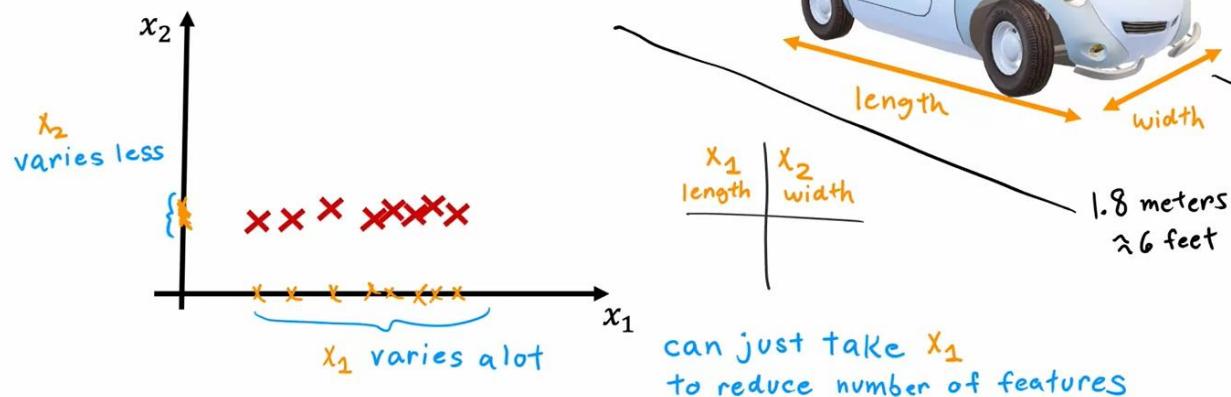


Other problematic cases:

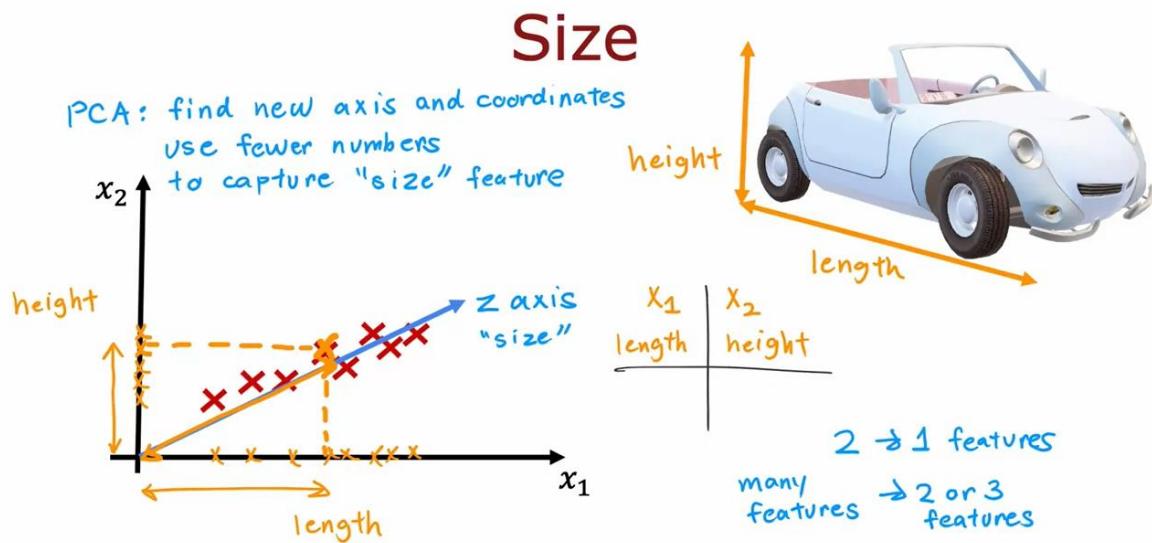
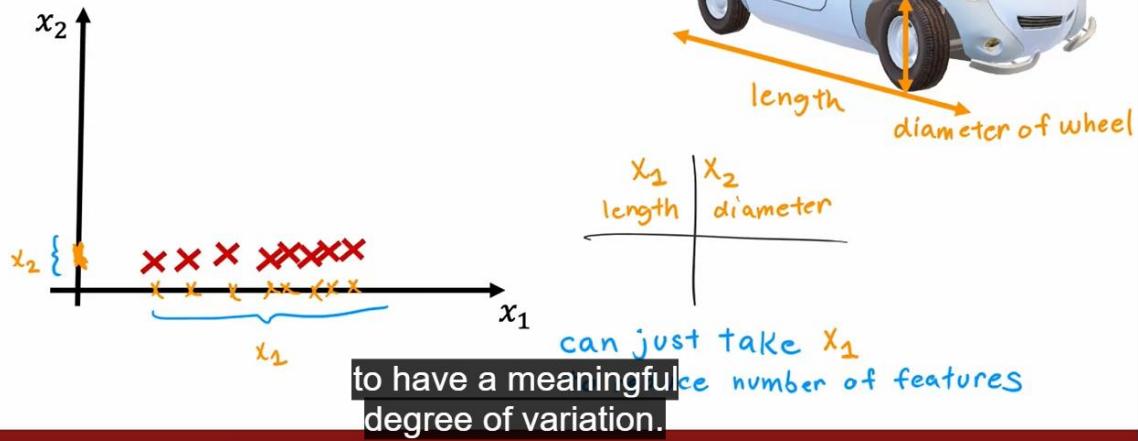
- • Maximizing user engagement (e.g. watch time) has led to large social media/video sharing sites to amplify conspiracy theories and hate/toxicity
- Amelioration : Filter out problematic content such as hate speech, fraud, scams and violent content
- • Can a ranking system maximize your profit rather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



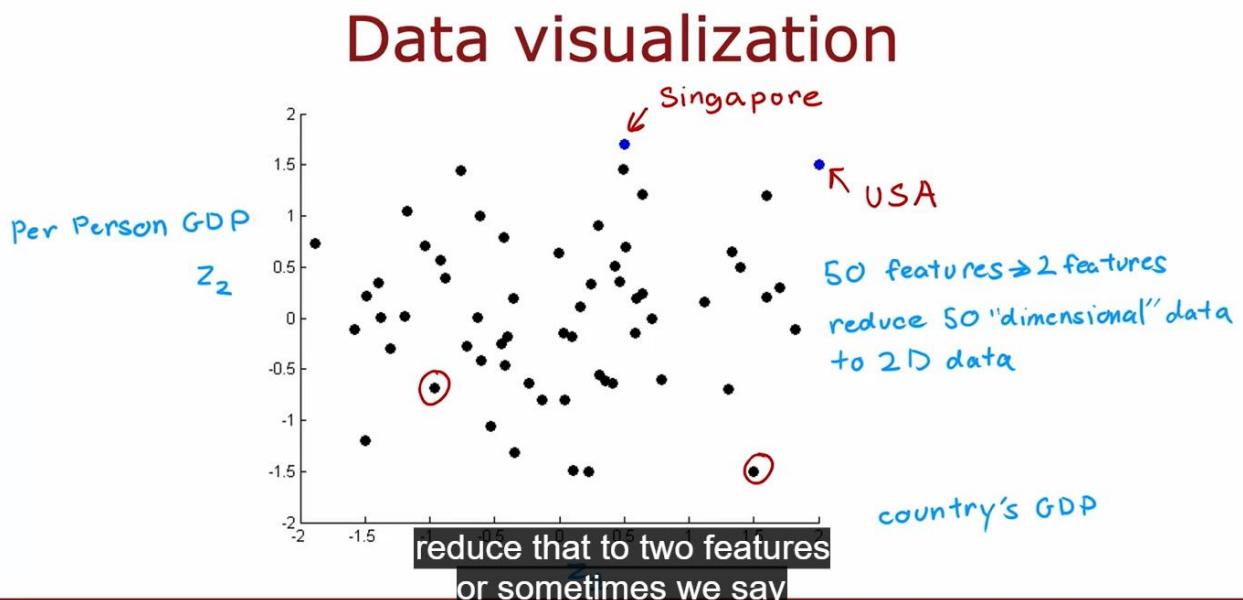
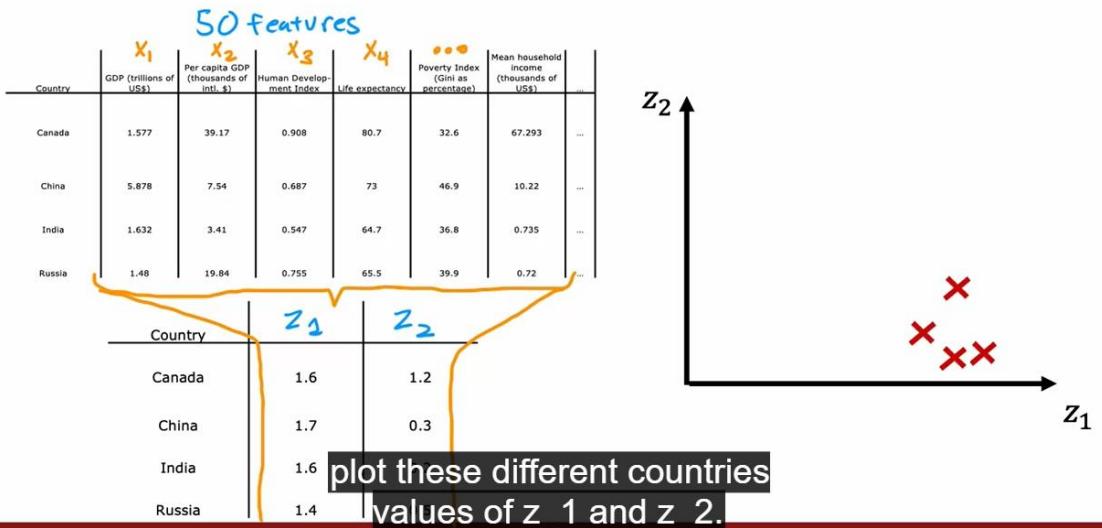
Car measurements



Car measurements



Let's look at one more example.

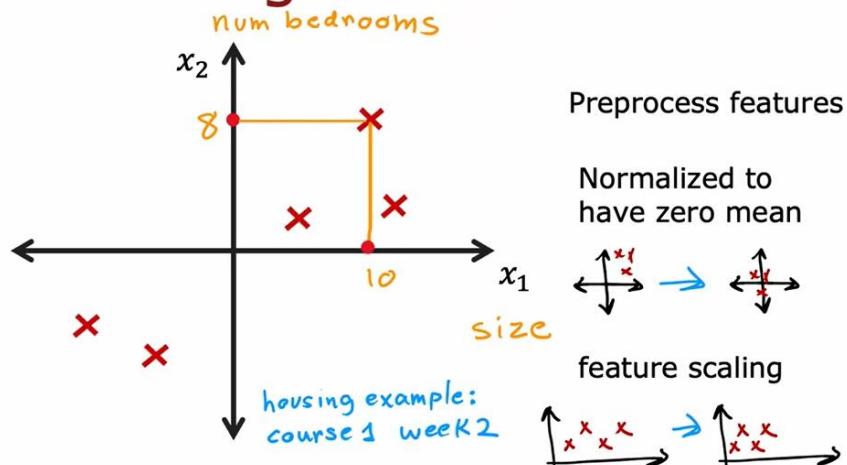


PCA algorithm

coordinates

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

Can we choose
a different axis?



To examine what PCA does,

DeepLearning.AI Stanford ONLINE Andrew Ng

Choose an axis

"project"
examples
onto the axis

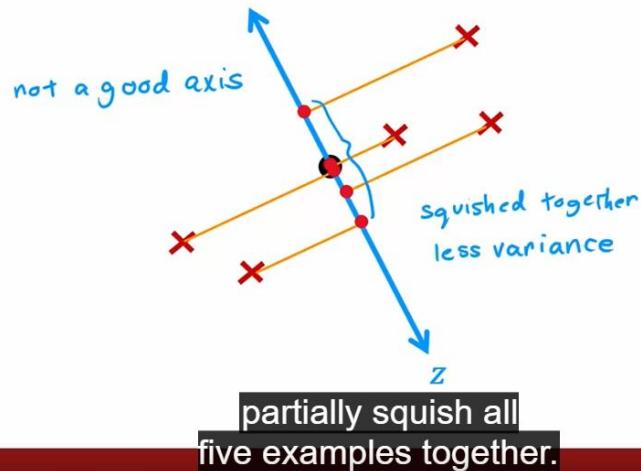


variance is large
capturing info
of original data

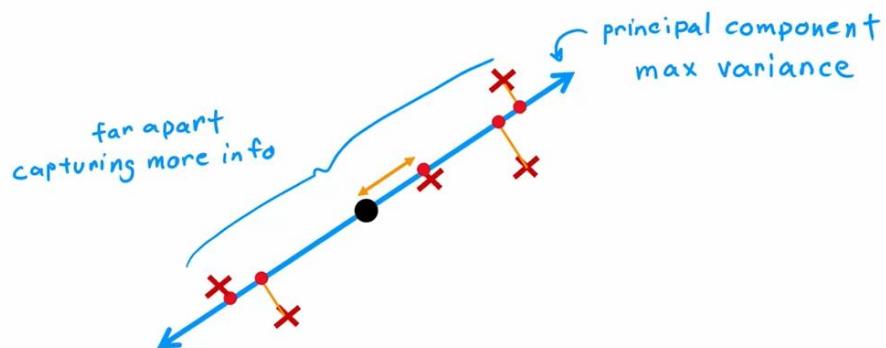
of the information in the
original five examples.

DeepLearning.AI Stanford ONLINE Andrew Ng

Choose an axis

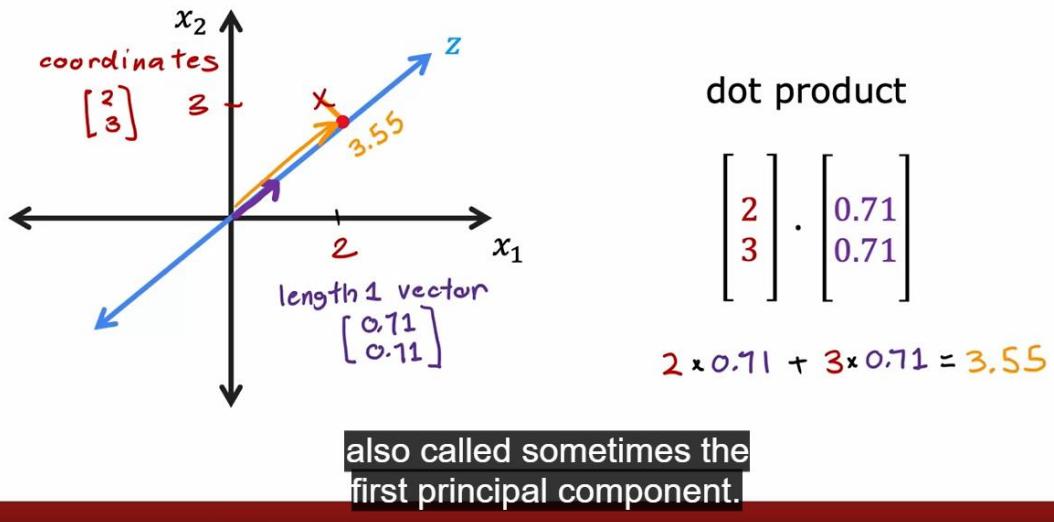


Choose an axis

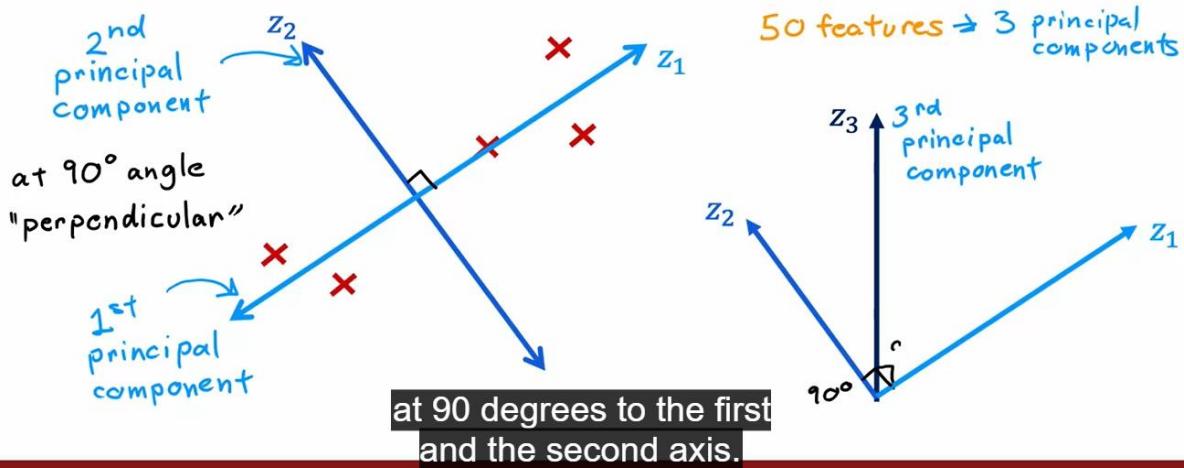


Let me show you a
visualization of how

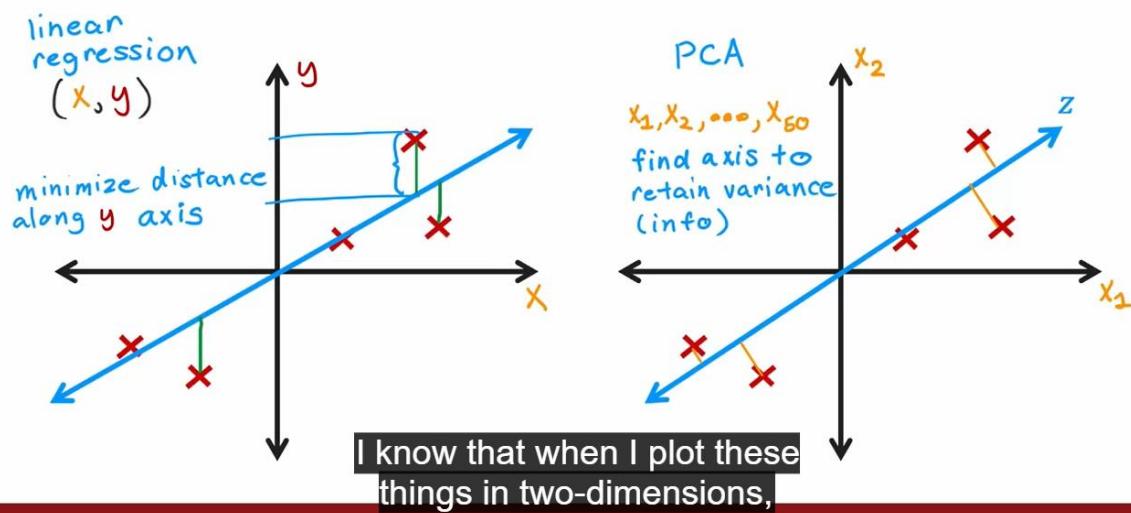
Coordinate on the new axis



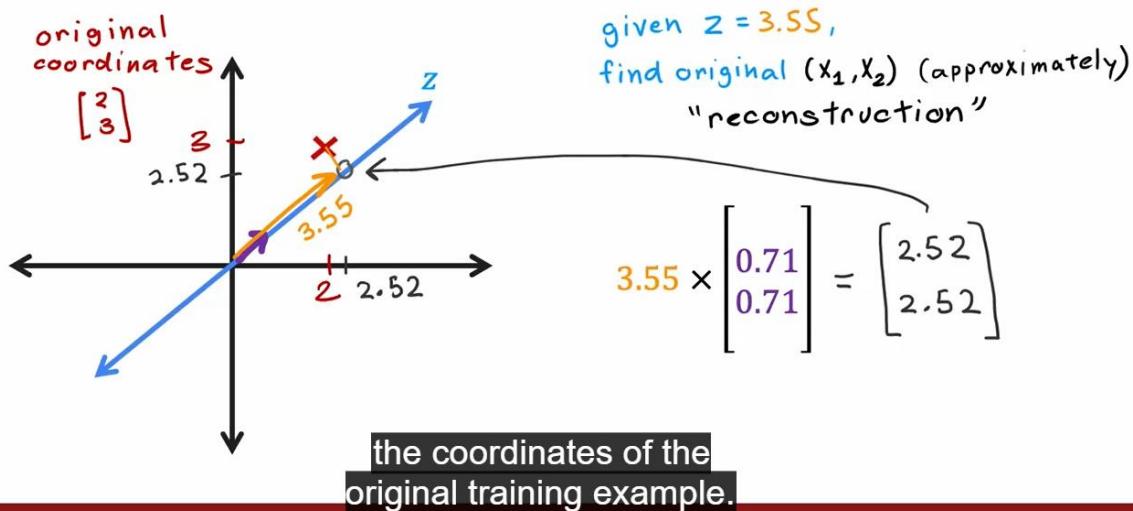
More principal components



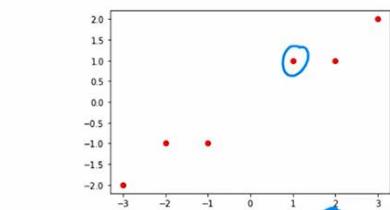
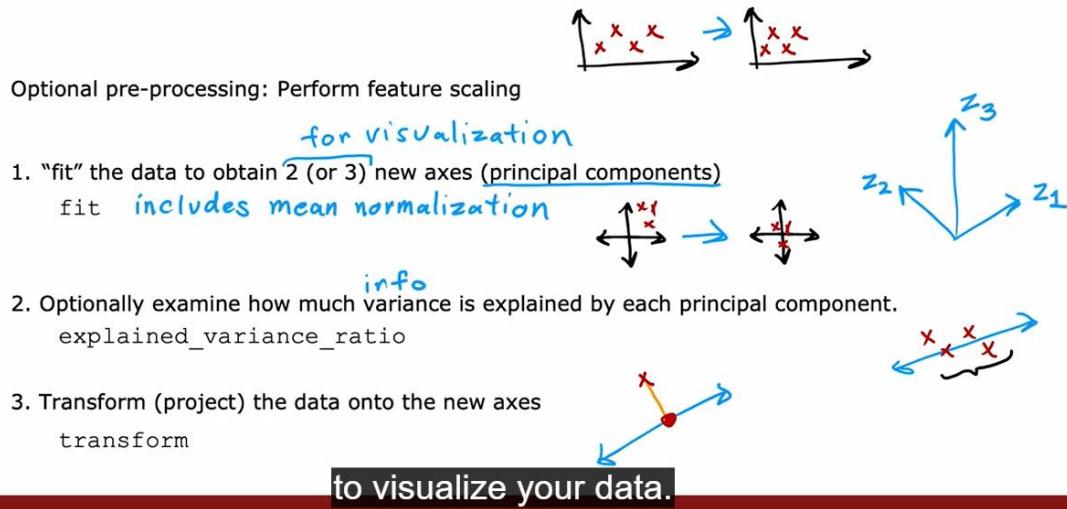
PCA is not linear regression



Approximation to the original data



PCA in scikit-learn

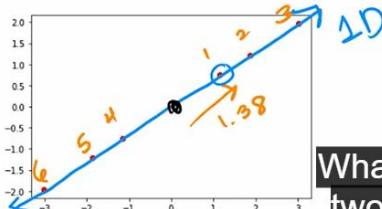


Example

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

2D

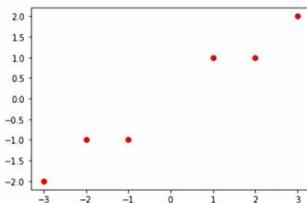
```
pca_1 = PCA(n_components=1)  
pca_1.fit(X)  
pca_1.explained_variance_ratio_ 0.992  
X_trans_1 = pca_1.transform(X)  
X_reduced_1 = pca_1.inverse_transform(X_trans_1)
```



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

What if you were to compute two principal components?



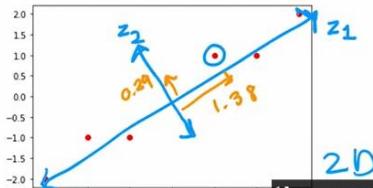


Example

```
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_  
X_trans_2 = pca.transform(X)  
X_reduced_2 = pca.inverse_transform(X_trans_2)
```



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]])

the optional lab where you can

Applications of PCA

Visualization *reduce to 2 or 3 features*

Less frequently used for:

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

$$n = 1000 \rightarrow 100$$

This used to make a difference
in the running time