

UnSupervised Learning

Clustering

Clustering
Anomaly
Recommender System
Reinforcement Learning

Clustering
Supervised

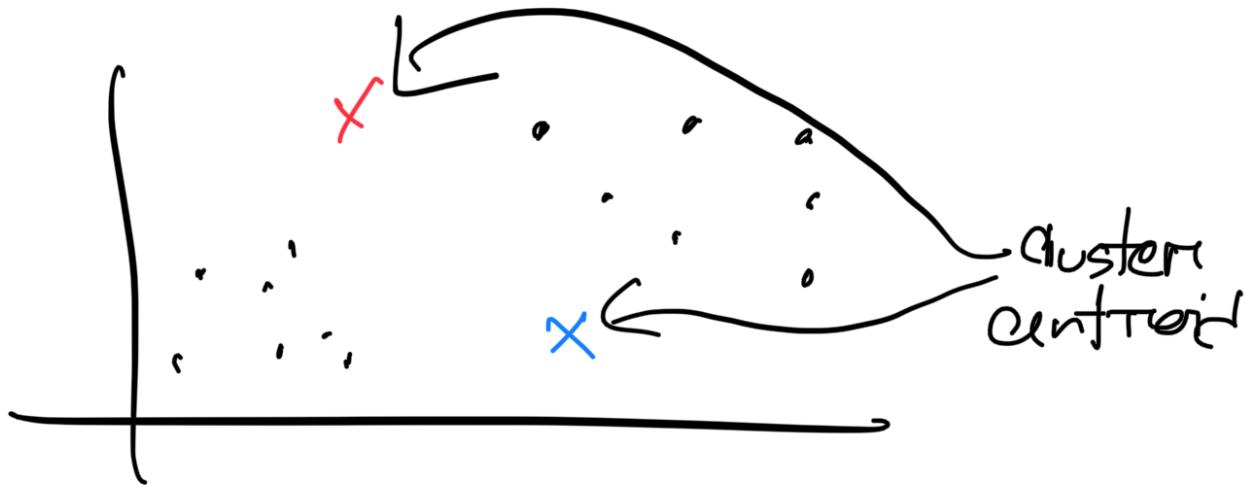


Training Set: $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$

Clustering

Growing Skills
Develop career

Astronomical data analysis



Step 1: Assign each point to its closest centroid

Step 2: Recompute the centroids

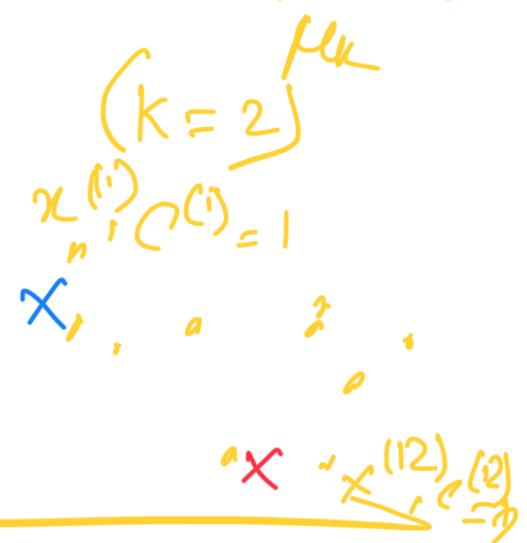
Cluster Centroids

Randomly initialize K clusters μ_1, μ_2, μ_k

Repeat {

Assign points to
cluster centroids

for $i = 1$ to m



$c^{(i)}$ = index (1 to B) of cluster
centroid cluster to $x^{(i)}$

$\min_{c=1}^B \|x^{(i)} - \mu_c\|^2$

$$\|\mathbf{m}_k\|^2 \rightarrow \|\mathbf{m}_k\|$$

Move Cluster centroids

for $k=1$ to K

μ_k = avg (mean) of points assigned to cluster k

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

}

Optimization Objective

$c^{(i)}$ = index of cluster

μ_k = cluster centroid k

μ_c = .. " of which example

Cost function

K -means algo

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

example a
 μ_k = cluster centroid is assigned
 $\mu_C(i)$ = " " of cluster
 to which example $x^{(i)}$ assigned
 Cost func

$$\bar{J}(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k) = \min \sum_{i=1}^m \|x^{(i)} - \mu_c^{(i)}\|^2$$

↑
distortion

Initializing K-means

Step 0: Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_k$

Repeat of

Step 1: Assign points to cluster centroids,

Step 2: Move cluster centers

Choose k from
Randomly pick k training
Set μ_1, μ_2, \dots all equal to the
 k example.

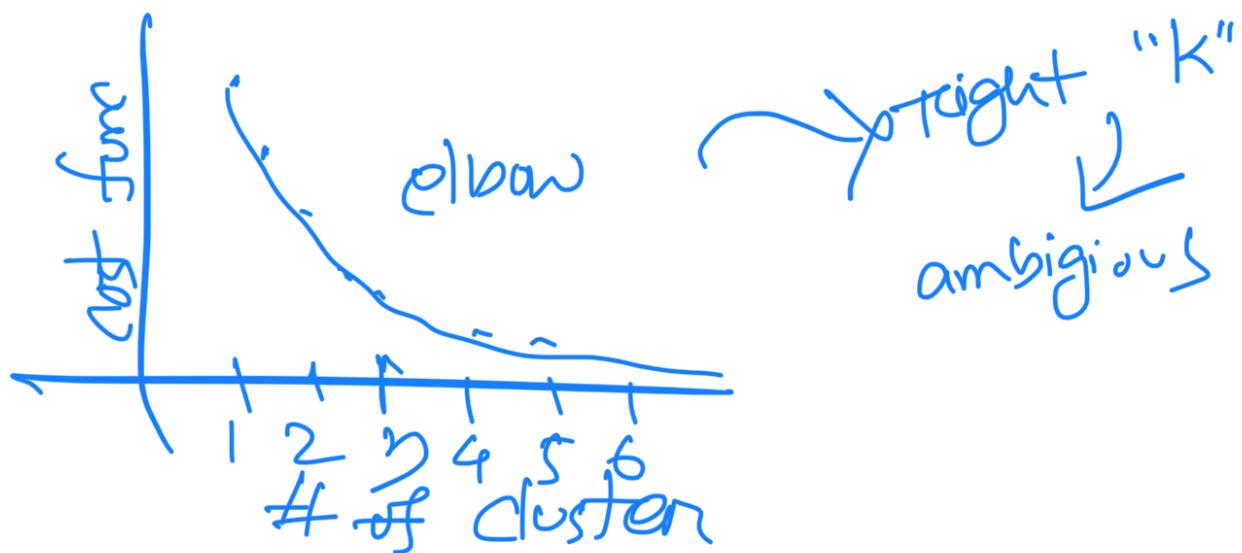


for $i = 1$ to 100 {

Randomly initialize k-means
Run k-means. Get $c^{(1)}, \dots, c^{(n)}$
 μ_1, \dots, μ_k
Compute cost func (distortion)
 $J(c^{(1)}, \dots, c^{(n)}, \mu_1, \dots, \mu_k)$

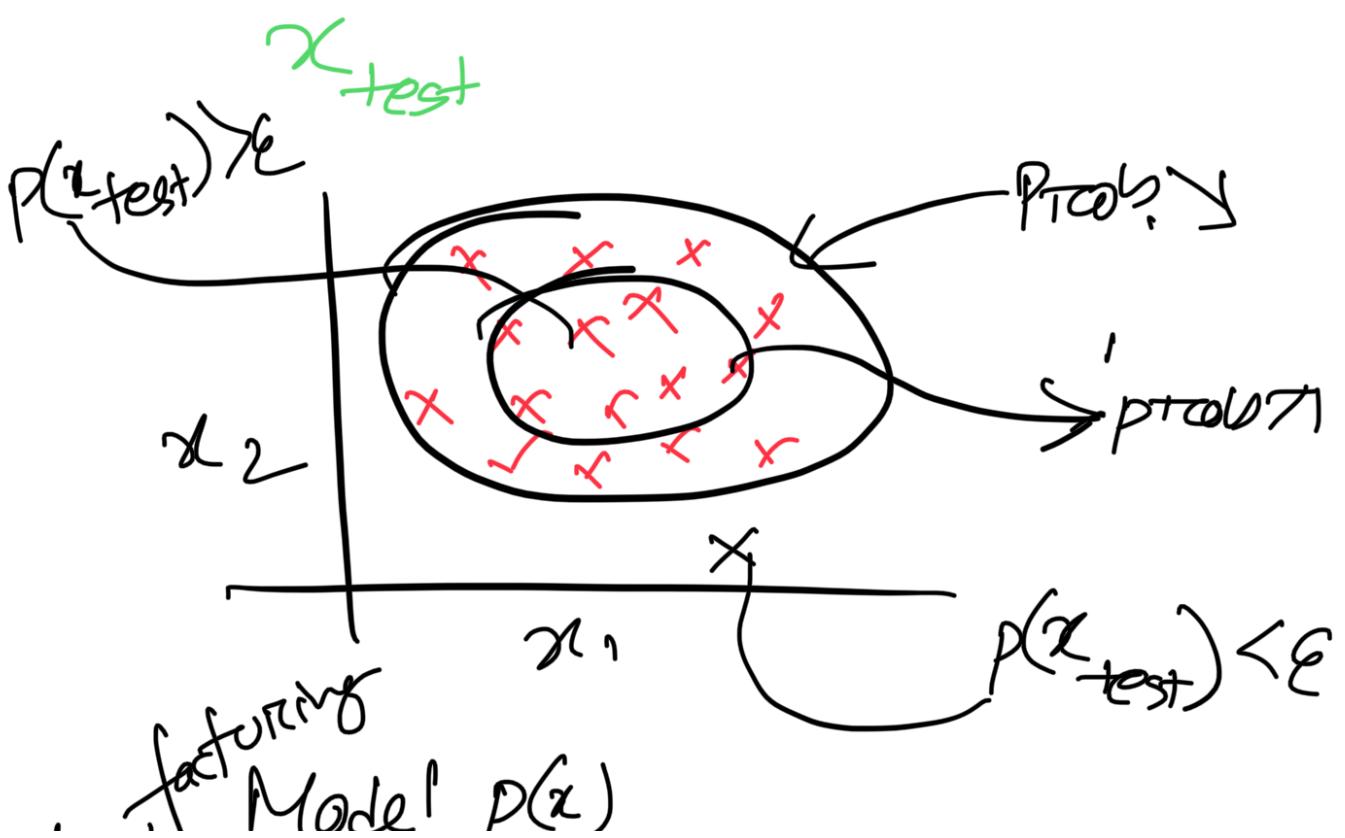
1

K means algo



Anomaly detection

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



Many
from C How often log in?
how many web pages visited?

Monitor

Gaussian distribution



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

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

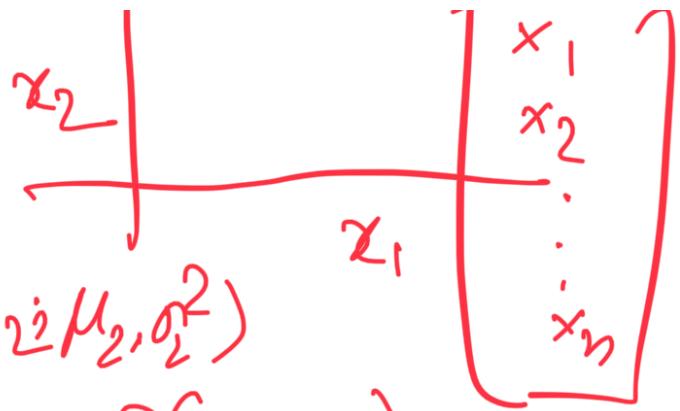
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

Density estimation

1

L

Each example



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

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

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

1. choose n features \rightarrow might have anomalous examples

② $\hat{\mu}_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$ $\hat{\sigma}_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \hat{\mu}_j)^2$

3. Give new examples

$$P(x) = \prod_{j=1}^n P(x_j; \hat{\mu}_j, \hat{\sigma}_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi \hat{\sigma}_j^2}} \exp \left(-\frac{(x_j - \hat{\mu}_j)^2}{2\hat{\sigma}_j^2} \right)$$

if $P(x) < \epsilon$

Evaluating a system
previously observed anomalous

$y=1$

non-anomalous
 $y=0$

include a few anomalous data

Train	6000	good	10K good engine	20 flawed
CV	2000	"	10 anomalies	
Test	"		10	"

fine ϵ fine ϵ

Algo evaluation

$$y = \begin{cases} 1 & P(x) \leq \epsilon \text{ (Anomaly)} \\ 0 & P(x) > \epsilon \text{ (Normal)} \end{cases}$$

metric use for performance

Anomaly
detection

Very Small \oplus

Supervised
Learning
large number $\oplus \&$

— ex.
very long \ominus ex.

diff types
of anomalies.
Hard of algo
to learn, future
anomalies are
different

from

Manufacturing

new unseen
defects

Monitoring
machines in
a data

Carefully choosing data

(+) Example

Enough \oplus examples
to get a sense of
what positive ex.
looks like.

Spam

Finding seen
defects
Weather prediction

Disease classification

Non-Gaussian features

plt.hist()

np.log(x)

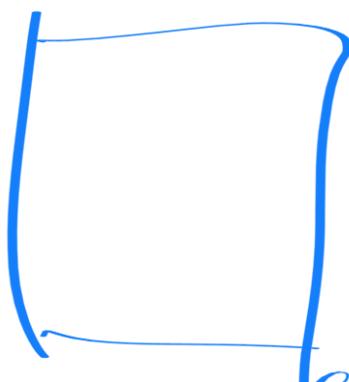
$$x \leftarrow \lg(x_1)$$

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

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

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

Recommended System
Predicting movie rating



$M \times U$ = # of users

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

$y(i, j)$ movie $i \rightarrow$ movie
 \sum rated $j \rightarrow$ # user

features of movie $n_v=4$
 $n_m=5$

$n=2$

$$\omega^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix} \quad b^{(1)} = 0 \quad x^{(3)} = \begin{bmatrix} 0.95 \\ 0 \end{bmatrix} \quad x^{(1)} = \begin{bmatrix} 0.95 \\ 0 \end{bmatrix}$$

$$x^{(2)} = \begin{bmatrix} 0.95 \\ 0 \end{bmatrix}$$

$$\omega^{(1)}, x^{(3)} + b^{(1)} = 4.95$$

$$\omega^{(2)}, x^{(1)} + b^{(2)}$$

↳ features

Cost function

~~$$J(\omega, b) = \frac{1}{m} \sum_{i: R(i, j) = 1} \left(\omega^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i, j)} \right)^2$$~~

$\sum_{i: R(i, j) = 1}$ $\sum_{k=1}^n (\omega_k^{(j)})^2$ ~~for~~ ~~feature~~

$$= \frac{1}{2} \sum_{i: R(i, j) = 1} \left(\omega^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i, j)} \right)^2$$

$$+ \gamma \sum_{k=1}^n (\omega_k^{(j)})^2$$

↪ K_{-1} -

$$= \frac{1}{2} \sum_{j=1}^m \sum_{\substack{i: \pi(i,j)=1}} (\omega^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \gamma \sum_{j=1}^m \sum_{k=1}^n (\omega_k^{(j)})^2$$

$$\omega^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$$

$$\omega^{(2)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$$

$$\omega^{(3)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$$

$$\omega^{(4)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$$

$$\omega^{(4)} \cdot x^{(i)} + b^{(j)}$$

$$\omega^{(1)} \cdot x^{(1)} \approx 5$$

$$\omega^{(2)} \cdot x^{(2)} \approx 5$$

$$\left. \begin{array}{l} x^{(1)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ x^{(2)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \end{array} \right\}$$

$$J(x) = \frac{1}{2} \sum_{j: \pi(i,j)=1} (\omega^{(j)} \cdot x^{(i)} + b^{(j)} - y^{(i,j)})^2 + \gamma \sum_{k=1}^m (x_k^{(i)})^2$$

$$J(x^{(1)}, x^{(2)}, \dots, x^{(m)}) = \frac{1}{2} \sum_{i=1}^m \sum_{\substack{j: \pi(i,j)=1}}$$

$$\frac{1}{2} \left(\omega^{(j)} x^{(i)} + b^{(j)} - y^{(i, j)} \right)^2 + \frac{\lambda}{2} \sum_{i=1}^m \sum_{k=1}^n (x_k^{(i)})^2$$

Put them together

$$J(\omega, b, x) = \frac{1}{2} \sum_{(i, j), R(i, j) = 1} \left(\omega^{(j)} x^{(i)} + b^{(j)} - y^{(i, j)} \right)^2$$

$$\frac{\lambda}{2} \sum_{j=1}^n \sum_{k=1}^n (\omega_k^{(j)})^2 + \frac{\lambda}{2} \sum_{i=1}^m \sum_{k=1}^n (x_k^{(i)})^2$$

Linear regression

Items liked by users
click on it

Predict $y^{(i, j)} \rightarrow \omega^{(j)} x^{(i)} + b^{(j)}$

$$g(\omega^{(j)}, x^{(i)} + b^{(j)})$$

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

binary

$$1 + e^{-z}$$

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

$$(f_{(\omega, b, x)}(x)) - (1 - y^{(i, j)})$$

$$\log(1 - f_{(\omega, b, x)}(x))$$

$$J(\omega, b, x) = \sum_{(i, j) \in \mathcal{N}(i, \hat{j})} L(f_{(\omega, b, x)}(x), y^{(i, j)})$$

Mean Normalization

$$\begin{aligned} \text{min } \omega^{(1)}, \dots, \omega^{(n_v)} \\ b^{(1)}, \dots, b^{(n_v)} \\ x^{(1)}, \dots, x^{(n_m)} \end{aligned} \quad \frac{1}{2} \sum_{(i, j) \in \mathcal{N}(i, \hat{j})} (\omega^{(\hat{j})} \cdot x^{(i)} + b^{(\hat{j})} - y^{(i, \hat{j})})^2 +$$
$$\lambda \frac{1}{2} \sum_{j=1}^{n_v} \sum_{k=1}^n (\omega_k^{(j)})^2 + \lambda \frac{1}{2} \sum_{i=1}^{n_m} \sum_{k=1}^{n_m} (x_k^{(i)})^2$$

mean normalization

$$(x_k^{(i)})^2$$

$$\begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & 2 & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \end{bmatrix} \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \end{bmatrix} \quad \mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 5 & 4 & ? \\ 0 & 5 & 0 & ? & ? \end{bmatrix} \begin{bmatrix} 2.25 \\ 1.25 \end{bmatrix} \quad \begin{bmatrix} 2.25 \\ 1.25 \end{bmatrix}$$

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

$$\omega^{(j)}, x^{(i)} + b^{(j)} + \text{fle};$$

$$\begin{aligned} \text{user-5} \quad \omega^{(5)} \begin{bmatrix} 0 \\ 0 \end{bmatrix} & \quad b^{(5)} = 0 \\ \omega^{(5)}, x^{(1)} + b^{(5)} + M_1 & = 2.25 \end{aligned}$$

Tensorflow implementation

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

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

$$J = (\hat{y} - y)^2$$

Gradient descent algorithm

Custom Training Loop

`w = tf.Variable(3.0)`

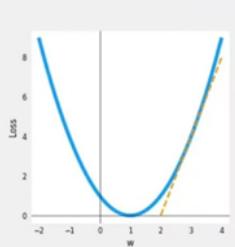
`x = 1.0`

Tf.variables are the parameters we want to optimize

Repeat until convergence

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

Fix $b = 0$ for this example



Optimize

```
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.
```

And if you write this with tf
our gradient tape as tape f.

Custom Training loop
↳ Auto diff / Autograd

Implementation in TensorFlow

Gradient descent algorithm

Repeat until convergence

$$\begin{aligned} w &= w - \alpha \frac{d}{dw} J(w, b, x) \\ b &= b - \alpha \frac{d}{db} J(w, b, x) \\ x &= X - \alpha \frac{d}{dx} J(w, b, x) \end{aligned}$$

```
# 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, ←
            num_users, num_movies, lambda) ←
            n_u n_m ←

        # 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

Collaborative filtering

$x_1 x_2 x_3$
 $\nwarrow \uparrow \nearrow$
 n

To find other items related to \sim

Find item k with $\alpha^{(k)}$
similar to $\alpha^{(i)}$

$$\sum_i^n (x_i^{(k)} - x_i^{(i)})^2 \quad x^{(k)} \sim x^{(i)}$$

Cold start

→ rank new items

→ show something reasonable

User side info. ..

Item: Genre, movie star, studio

User: Demographics (age, gender)

Collaborative filtering

Recommend item → based on user rating
Similar to you ↳

Content based filtering

Recommend items → features of users
& item → Good Match

New tools to
revisit the last
points if you go
back & rewatch
something

User Features:

- Age
- Gender
- Country
- Movie Watched
- Average rating

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

Movie Features

- Year
- Genre
- Reviews
- Average rating

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

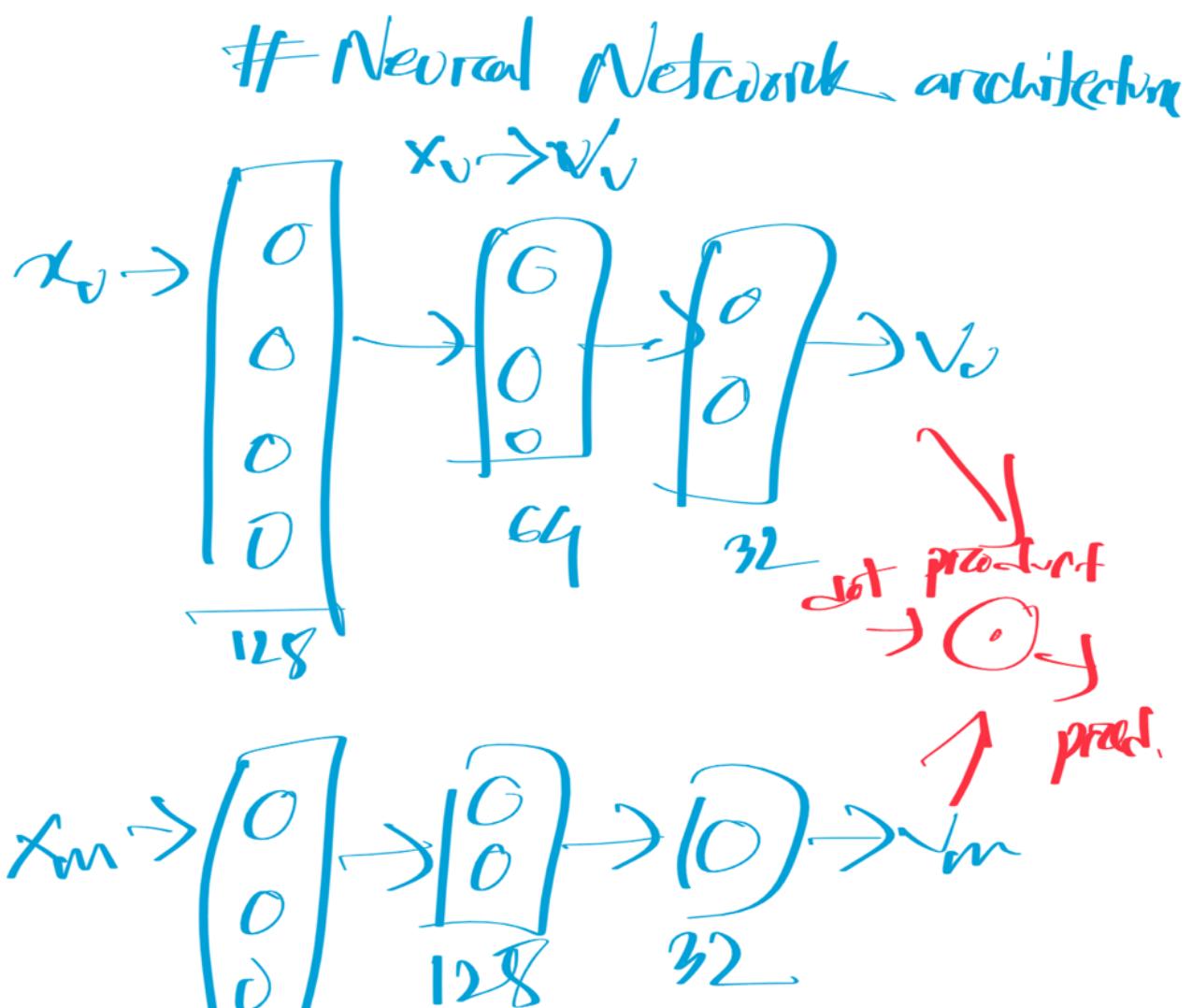
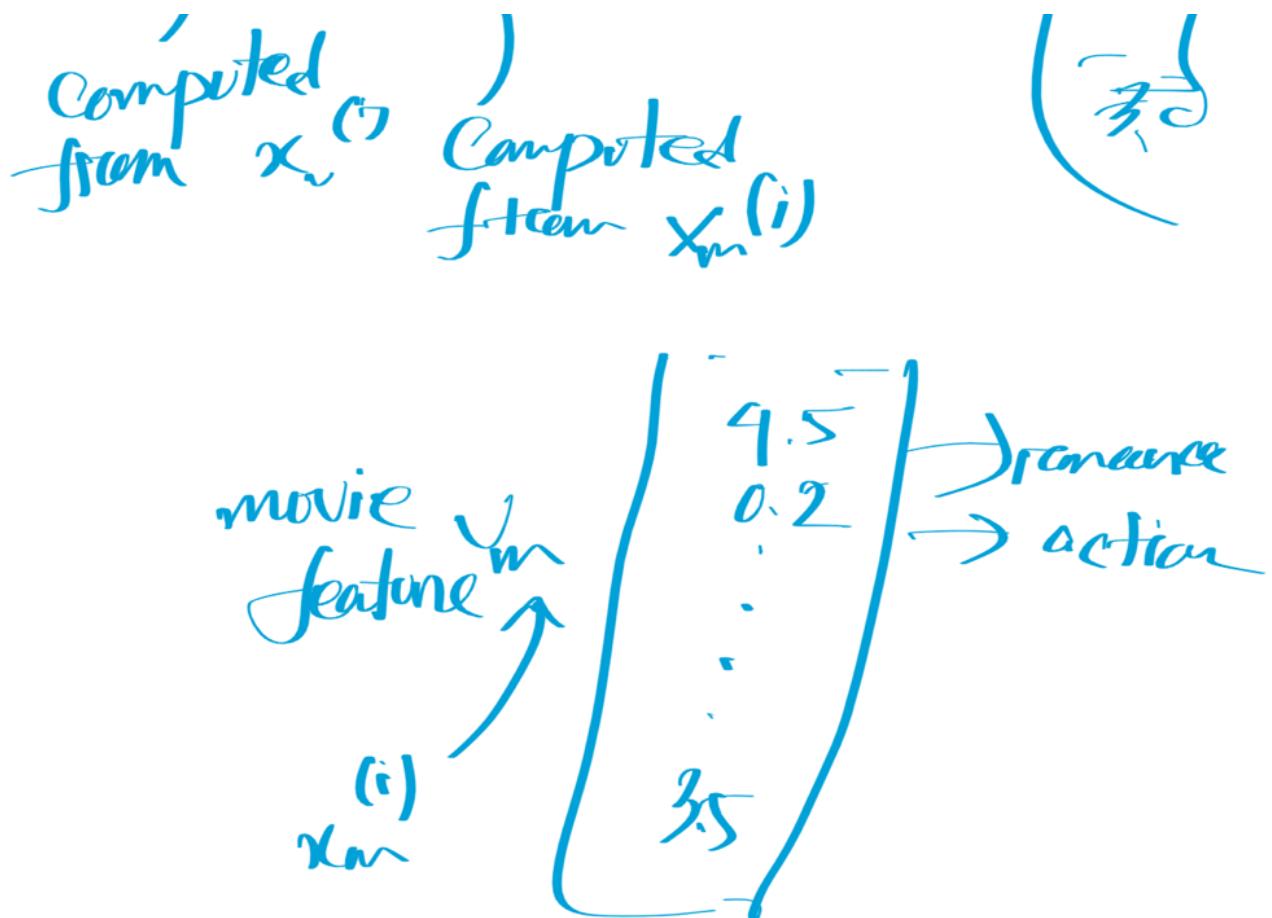
$\omega^j x^i + \beta^j$

$v_v^{(j)}$

$v_m^{(i)}$

Yeast's
preference

$v_v^{(j)} = 0.9$
 $v_m^{(i)} = 0.1$



256

$g(\mathbf{v}_i^{(i)}, \mathbf{v}_m^{(i)}) \rightarrow$ predict the probability $y^{(i,i)} = 1$

Cost

$$J = \sum_{(i,j) \in \mathcal{N}(i,j) \neq i} (\mathbf{v}_i^{(i)} \cdot \mathbf{v}_m^{(i)} - y^{(i,j)})^2 + \text{Reg} \sum_{i=1}^n \|\mathbf{v}_i^{(i)}\|^2$$

$\mathbf{v}_i^{(i)} \rightarrow$ 32 length \rightarrow user i

$\mathbf{v}_m^{(i)} \rightarrow$ 32 " \rightarrow movie i
feature $x_m^{(i)}$
features $x_m^{(i)}$

$$\|\mathbf{v}_m^{(k)} - \mathbf{v}_m^{(i)}\|^2$$

efficiently find recommendation
from a large set

2 Steps
Retrieval
large list of items

Ranking
rank the list so get very initial time

→ 10 movies → 10 most similar

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

→ 3 genre → Suggest 10
→ top 10 → in the country

Display the ranks

$$= \frac{1}{2} \sum_{i,j} (\omega^{(j)} x^{(i)} + b^{(j)})^2$$

$\rightarrow \sum_{i,j} (\omega^{(j)})^2$

$$\cdot \cdot \cdot \overset{1}{2} \underset{k=1}{\overset{1}{\sim}} \cdot \cdot \cdot$$

Collaborative filtering

Recommend items \rightarrow rating of users
gave similar rating
as you

Content-based \rightarrow features of users & item
to find good fit

$\pi(i,j) = 1 \rightarrow$ user j has rated item
 $\gamma(i,j)$ rating given by user j on item i

User features

$x_j^{(G)}$

for user j

Age/gender/Country/Movie watched/
Avg. rating per genre

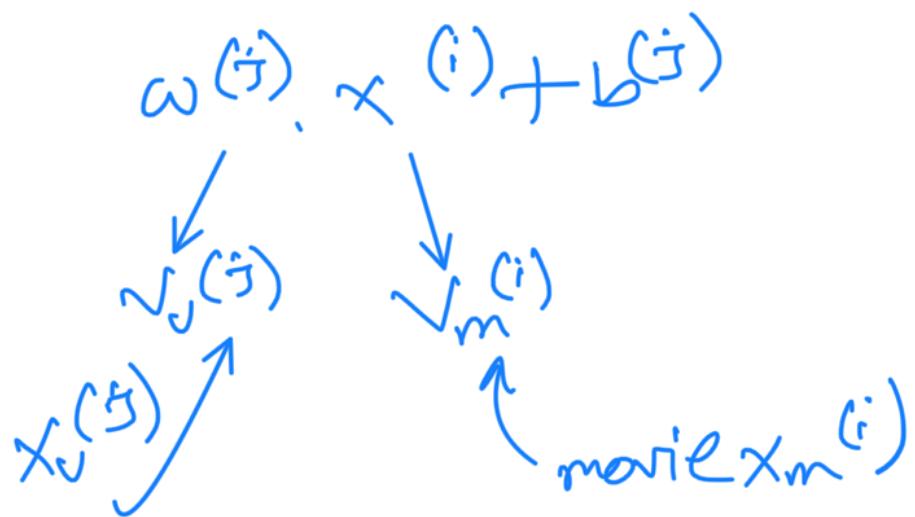
Movie features

Year/genre/genres/Review/Avg rating

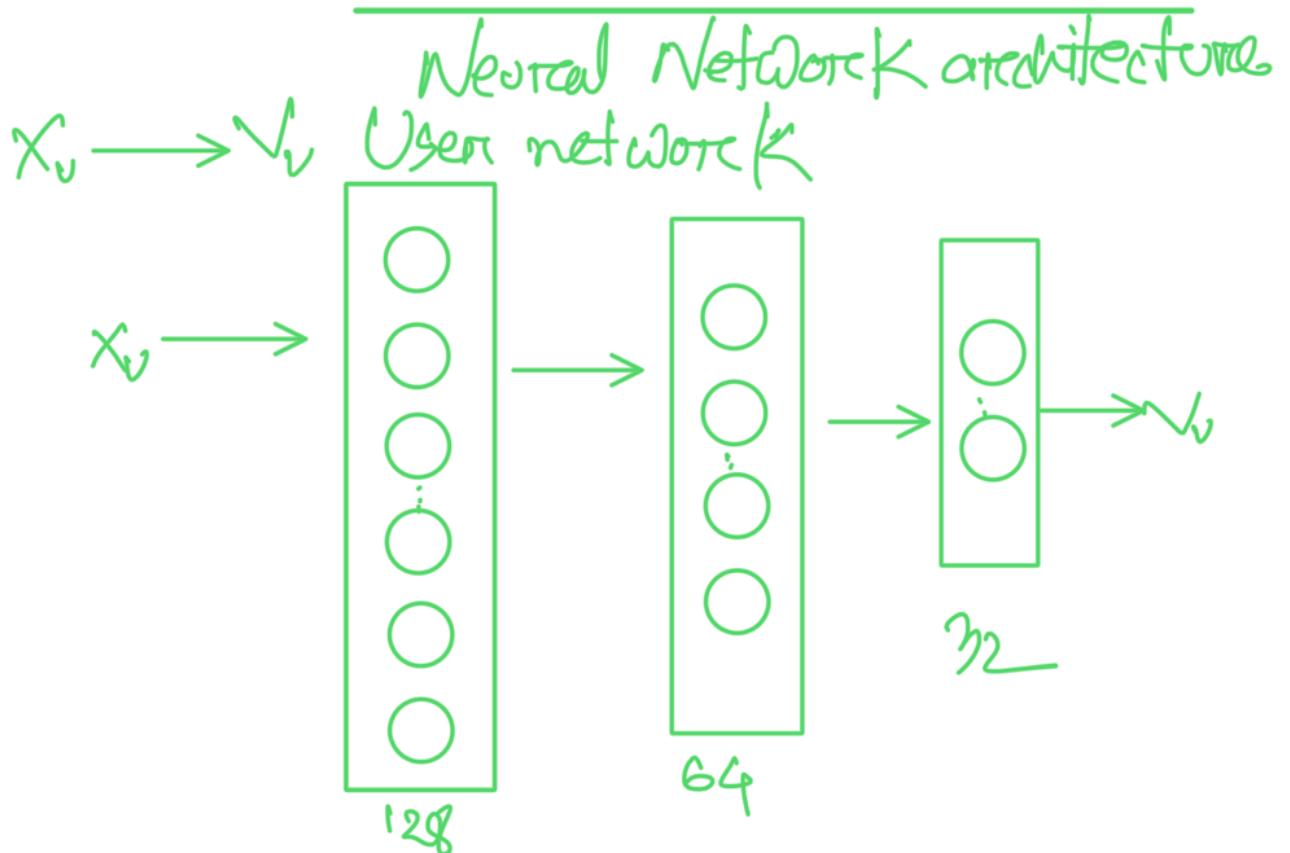
$V(i)$ for movie i
 x_m

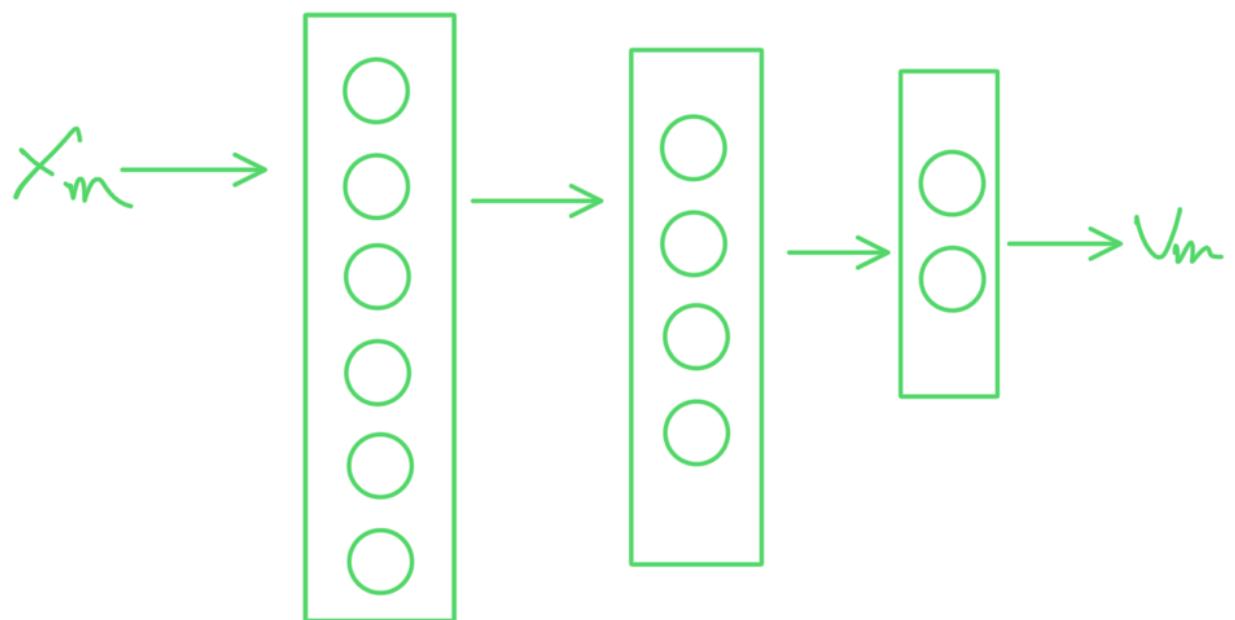
$x_m \neq x_u$ (Size not be same)

Learn to match



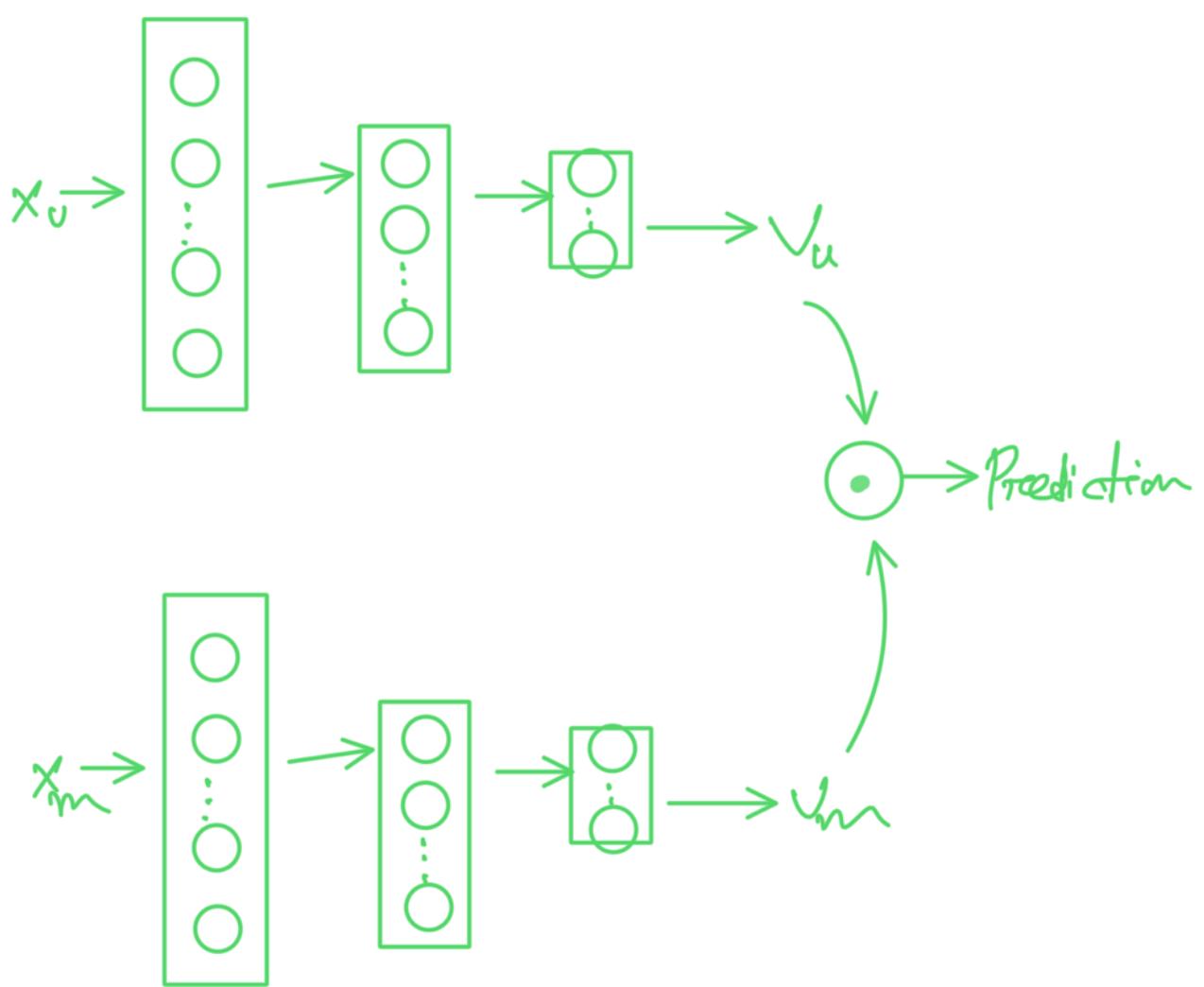
Content Based filtering





$g(v_u, v_m) \rightarrow \text{Probability } y^{(ij)}$

①



Cost Func: $\bar{J} = \sum_{\substack{(i,j) : r(i,j) \\ = 1}} (v_u^{(j)} \cdot v_m^{(i)} - y^{(i,j)})^2 + \lambda \|\mathbf{v}\|_2^2$ Regularization term

To find movies similar to movie i :

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

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

Recommending from
large Dataset

Retrieval & Ranking

① for each 10 watched movie

$$\text{find 10} \\ \|v_m^{(k)} - v_m^{(i)}\|^2$$

② most viewed 3 genre \rightarrow top 10

③ Top 20 movies in the country

Combined Retrieved item

remove duplicates or
already watched

Ranking



Based on predictions

Display ranks

Another procedure \rightarrow do VM then
take VM & do of
real time product

Retrieval steps

Retrieve more \rightarrow performance \uparrow
recommended \downarrow

to optimize \rightarrow carry offline exp.

retrieve additional items
more relevant to recommend

$$P(Y^{(r,i)}) = 1$$

Movies \rightarrow 5★

products most likely purchased

Ads \rightarrow clicked \oplus high bid

Products \rightarrow largest profit

Video → max watched time

Travel industry



Payday loans

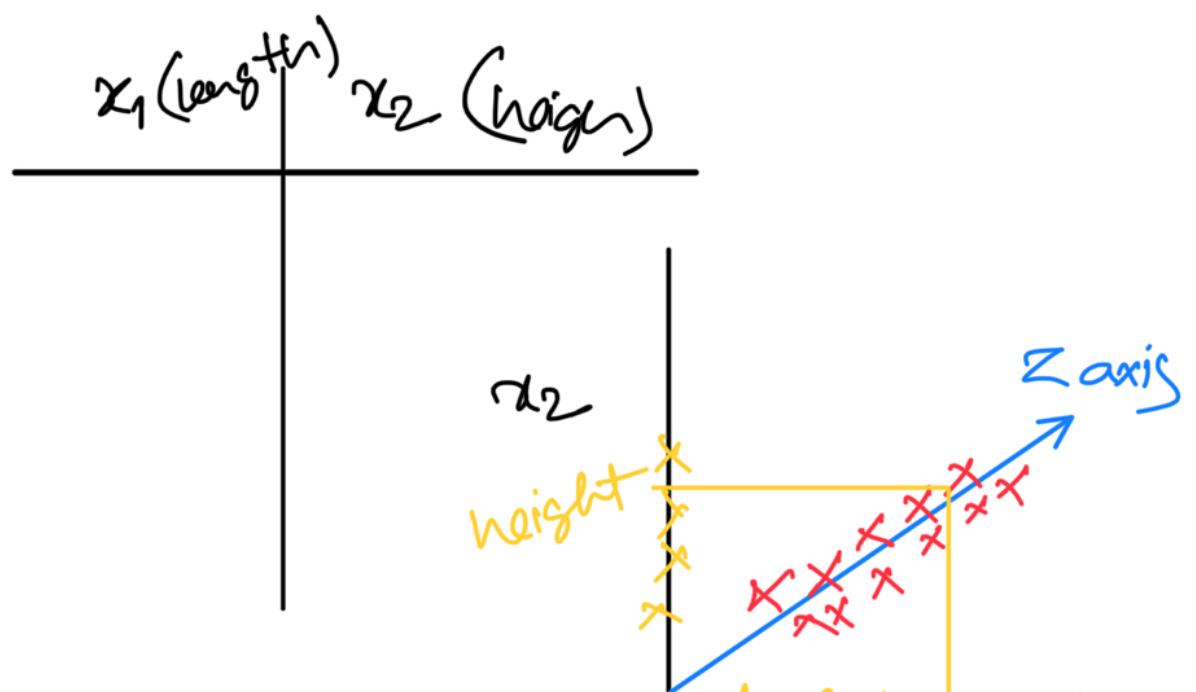
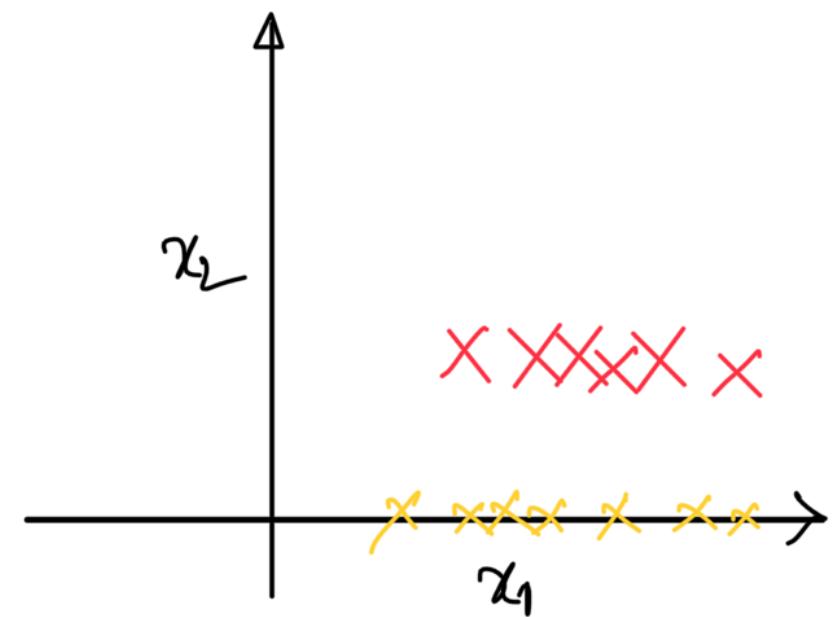
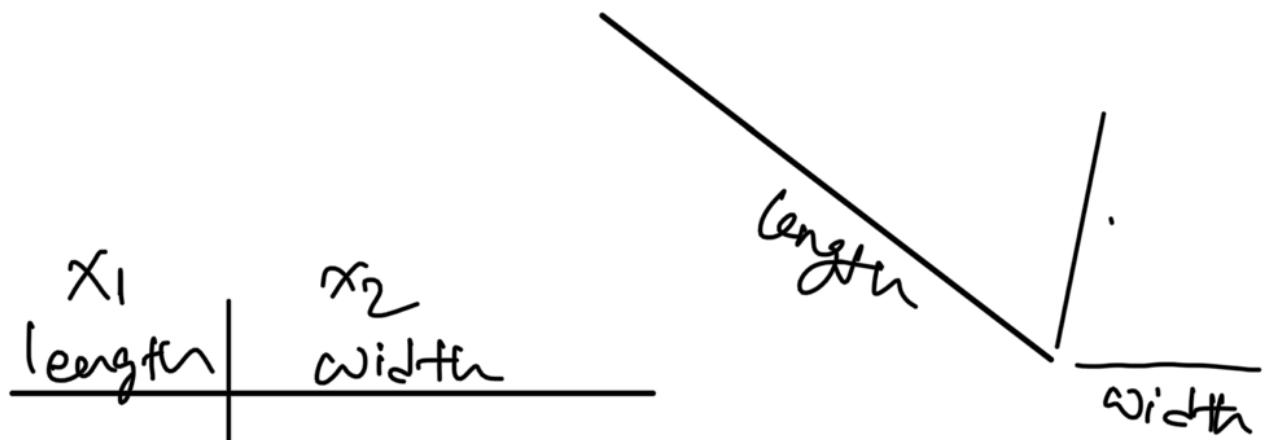
squeeze customers

Bid higher
for adds

more
more
profit

Principal Component
Analysis

Core Measurement



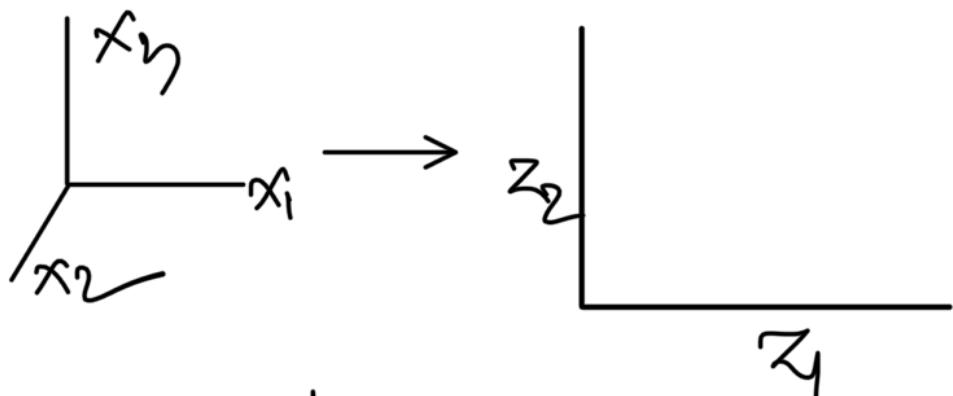


 find new axis & coordinates
 use fewer numbers
 to capture "size" features

$$2 \rightarrow 1$$

$$1000 \rightarrow 1$$

3D \rightarrow 2D



life expectancy

50 features \rightarrow 2 features
 50 D \rightarrow 2 D

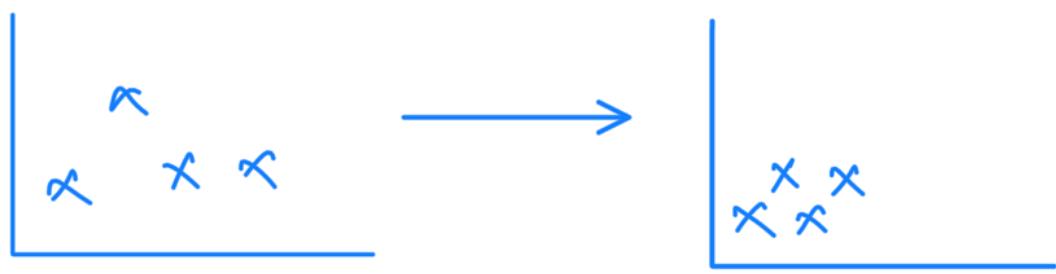
2 \rightarrow axis

PCA Algorithm

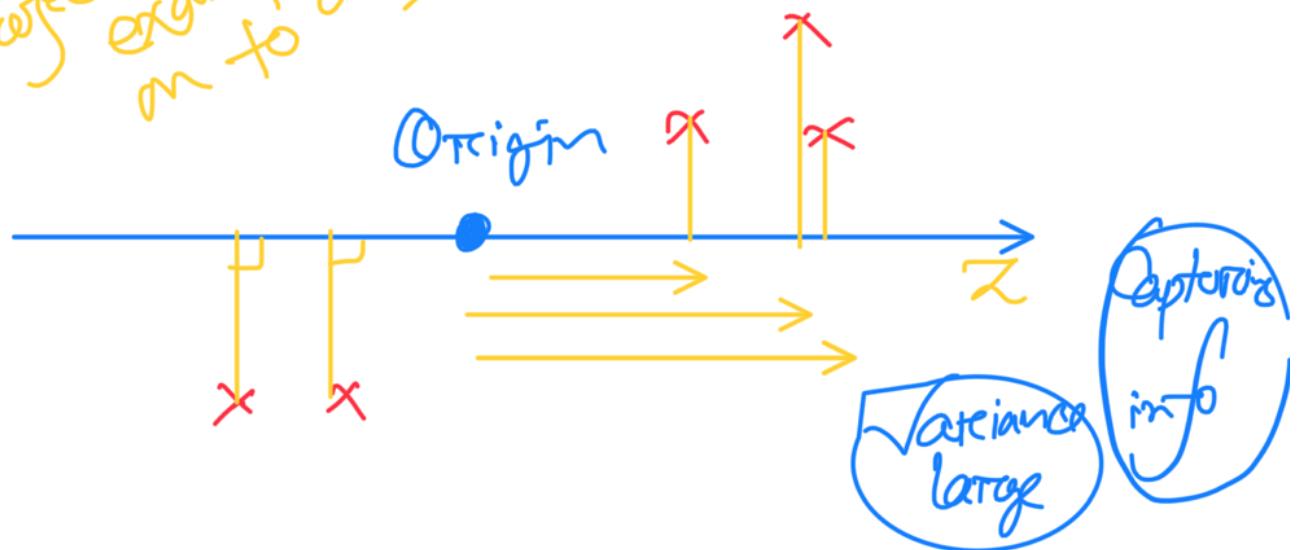
Normalized to have zero mean



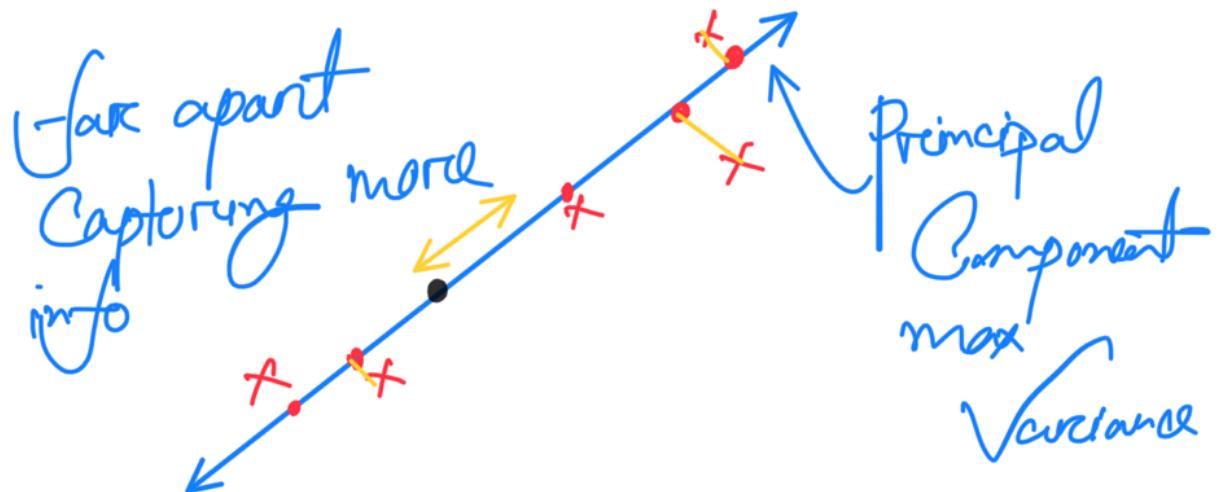
Feature Scaling



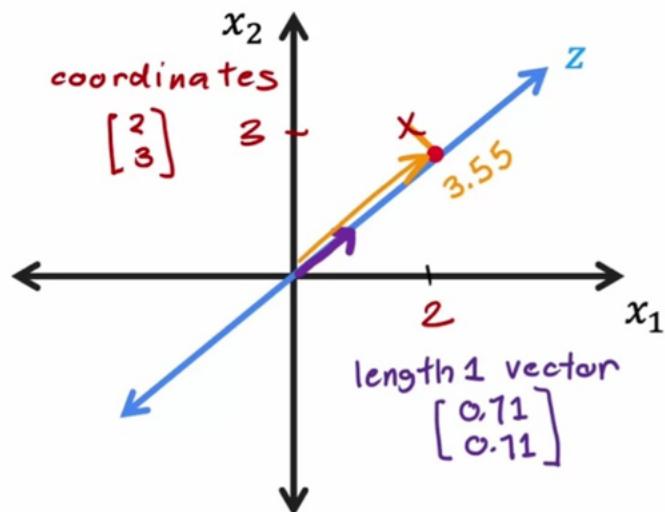
Project examples on to axis



Squished z axis



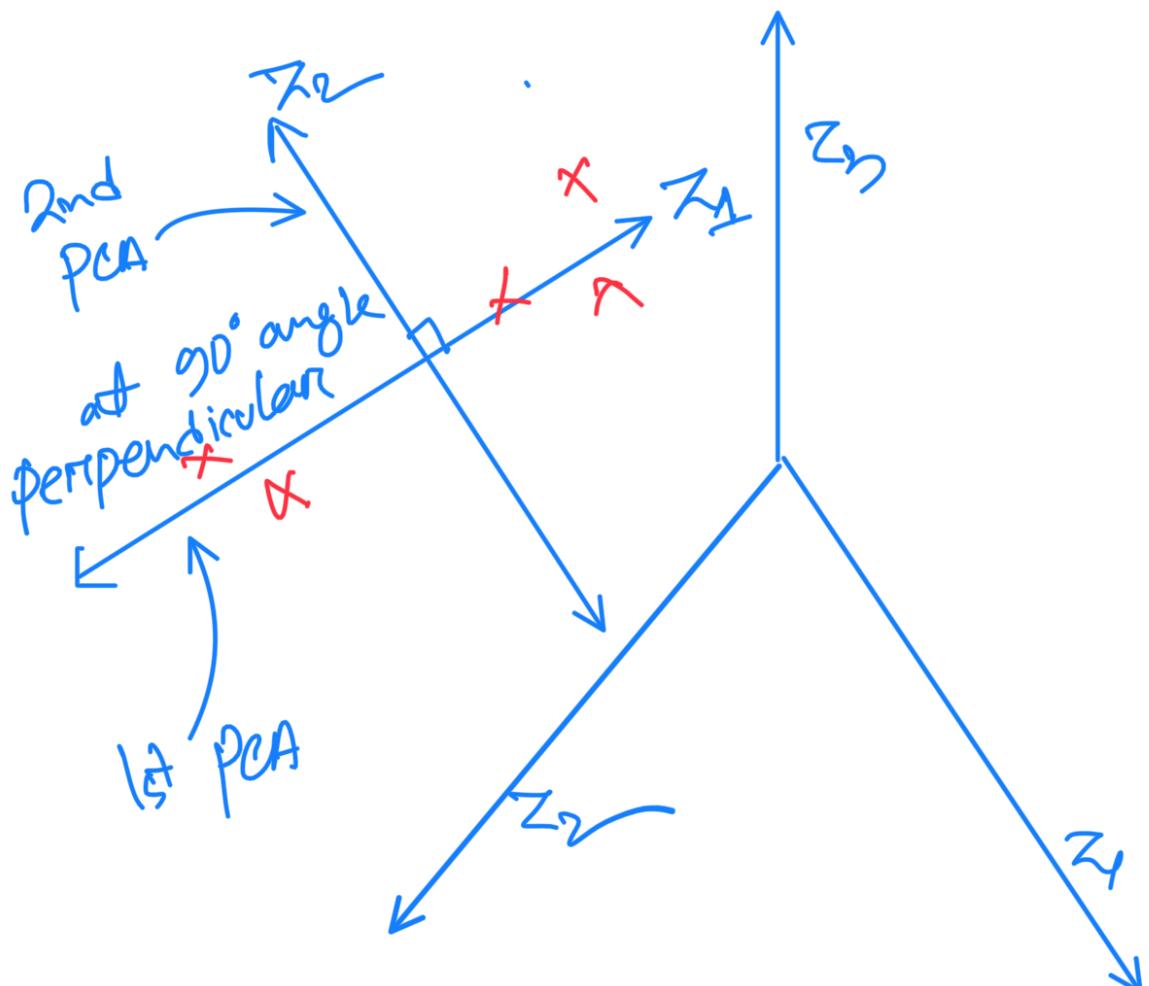
Coordinate on the new axis



dot product

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

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



linear regression
minimize distance
along y axis

PCA
 x_1, x_2, \dots, x_50
 find axis to
 retain variance
 (info)

TC construction

$$z = 3.55$$

find original (x_1, x_2)
 approximately
 $z \approx 0.71$ $r \approx 7$

$$\begin{bmatrix} 0.72 \\ 0.71 \end{bmatrix} = \begin{bmatrix} 0.52 \\ 2.53 \end{bmatrix}$$

preprocessing : perform feature scaling

① fit the data to obtain 2 (or 3) new axis
 (principle components)
 include mean normalization

② Originally examine how much variance
 is explained by each PCA
explained variance ratio

③ Transform (project) the data onto
 new axis
transform

Application
 visualization

Data compression

(to reduce storage or transmission costs)

Speed of training data (Supervised learning)

