

Day-96, Folgun 22, 2082 (Mar 6, 2026)

- ① Using Pop-item features
- ② Collaborative filtering Algorithm
- ③ Supervised vs UnSupervised ML Algorithm
- ④ Recommenders $j(w, b, x)$ using even $|x|$ as Parameters

Source: Coursera Course - Unsupervised Learning, Recommenders, Reinforcement

Learning :
Week(2)

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

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

$$x^{(1)} \Rightarrow \begin{bmatrix} 0.9 \\ 0 \end{bmatrix}$$

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

For user i : Predict rating for movie j as: $w \cdot x^{(i)} + b$
 $w^{(i)} \Rightarrow \begin{bmatrix} 5 \\ 0 \end{bmatrix}$ $b^{(1)} \Rightarrow 0$ $x^{(3)} = \begin{bmatrix} 0.99 \\ 0 \end{bmatrix}$ \rightarrow just linear regression:
 $w^{(1)} \cdot x^{(3)} + b^{(1)} \Rightarrow 4.95$

Now for each user we can predict
 for user j . Predict j 's rating for movie (i)

$$\sum_{j=1}^n \delta(i,j) = 1.$$

(0st function):

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

Use only user j has rated the movie i is the δ .

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
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

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

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

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

DeepLearning.AI

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

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

Regulation terms:

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

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

- $m^{(j)}$ = no. of movies rated by user j
- To learn $w^{(j)}$, $b^{(j)}$

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

number of features

The diagram shows the cost function equation for user j . A yellow box surrounds the summation term $\sum_{i: r(i,j)=1}$. Inside this box, a blue circle highlights the index i and the condition $r(i,j)=1$. Two blue arrows point from the handwritten note "number of features" to the n in the term $(w_k^{(j)})^2$.

COST function:

To learn Parameters $w^{(j)}, b^{(j)}$ for user j :

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

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

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

$f(x)$

Under the Hood we can use Linear Regression'

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 :

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

Collaborative filtering algorithm

→ linear Regression needs x_1 and x_2 . But what if you don't have these x_1 and x_2 → then we need to extract or find from the dataset itself.

ignoring $b^{(i)}$.

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

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

$$w^{(3)} \cdot x^{(3)} \approx 0$$

$$w^{(4)} \cdot x^{(4)} \approx 0$$

What value must

$$x^2 \text{ hold?} \\ x^{(1)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

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	?	?	?

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

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

Here's the data
that we had before.

$b^{(4)} = 0 \leftarrow w^{(4)} \cdot x^{(4)} + b^{(4)}$
Using

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

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

Goal: To learn meaningful features for movies based on user ratings, can accurately predict missing ratings -

① Find the best x^i (movie features) and w^j (user preferences) that minimize prediction.

② Adjust feature values (x^i) to best fit the given ratings.

$x^1 [1, 0]$
 (High Romance / Not Action)

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

$w^{(1)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}, w^{(2)} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}, w^{(3)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}, w^{(4)} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}$ }
 $b^{(1)} = 0, b^{(2)} = 0, b^{(3)} = 0, b^{(4)} = 0 \leftarrow$ }
 using $w^{(j)} \cdot x^{(i)} + b^{(j)}$
 $\{ w^{(1)} \cdot x^{(1)} \approx 5$
 $w^{(2)} \cdot x^{(1)} \approx 5$
 $w^{(3)} \cdot x^{(1)} \approx 0$
 $w^{(4)} \cdot x^{(1)} \approx 0$ }
 $\rightarrow x^{(1)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

Cost function

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

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

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

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

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

In this algorithm we can actually learn features from the movies and ratings.

Given $w^1, b^1, w^2, b^2,$

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.

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

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

this term here is exactly

the same as this term here.

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	Cow
$\rightarrow i=1$ Movie 1	5	?
$\rightarrow i=2$ Movie 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:

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

the regularization term from

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:

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

as a function of w , b , and x ?

Gradient Descent

collaborative filtering

Linear regression (course 1)

repeat {

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

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

}

parameters

w, b ~~for~~

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

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

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

w and b here to denote
the parameters for all

What Are We Trying to Do?

We are trying to predict movie ratings for users and, at the same time, learn features for movies (e.g., how much romance or action they contain).

The Problem:

- Normally, in machine learning, we need predefined features (like romance/action scores) to train a model.
- But here, we don't know these features in advance!
- So, we need to learn them from user ratings instead.

How We Solve It:

1. User Preferences (w^j)

- Each user has a preference vector (w^j) showing how much they like different movie features.

2. Movie Features (x^i)

- We adjust movie features (x^i) so that when multiplied by user preferences, they match the given ratings.

3. Cost Function (Error Minimization)

- We use squared error to measure the difference between predicted and actual ratings.
- We optimize both user preferences and movie features using gradient descent.

4. Collaborative Filtering

- Since multiple users rate the same movies, we can infer good features by looking at shared ratings.
- This helps predict ratings for users who haven't watched a movie yet.

1

What Are We Trying to Do?

We want to predict movie ratings for users while learning meaningful features for movies from data.

2

Step 1: Predicting Movie Ratings

Given a user j and a movie i , the predicted rating is computed as:

$$\hat{y}_{i,j} = w^j \cdot x^i + b^j$$

where:

- w^j is the preference vector for user j (how much they like different genres).
- x^i is the feature vector for movie i (how much action, romance, etc.).
- b^j is the user's bias (how lenient or strict they are in rating).
- \cdot represents the dot product.

To simplify, we ignore b^j (set it to 0) for now.

$$w^j := w^j - \alpha \frac{\partial J(w^j)}{\partial w^j}$$

3

Step 4: Combining Everything (Collaborative Filtering Algorithm)

The final cost function includes both user preferences and movie features:

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

where:

- n_u is the number of users, and n_m is the number of movies.
- The first term minimizes prediction error.
- The last two terms regularize w and x .

We optimize this using gradient descent by updating both w and x simultaneously.

Step 2: Learning Movie Features (x^i)

We assume we already have user parameters w^j , and we want to determine good values for x^i by minimizing the error.

For a given movie i , the cost function is:

$$J(x^i) = \frac{1}{2} \sum_j (w^j \cdot x^i - y_{i,j})^2$$

Key Takeaways

- We don't need predefined features; we learn them from user ratings.
- Collaborative filtering allows multiple users' ratings to guide feature learning.
- We optimize both user preferences (w) and movie features (x) using gradient descent.

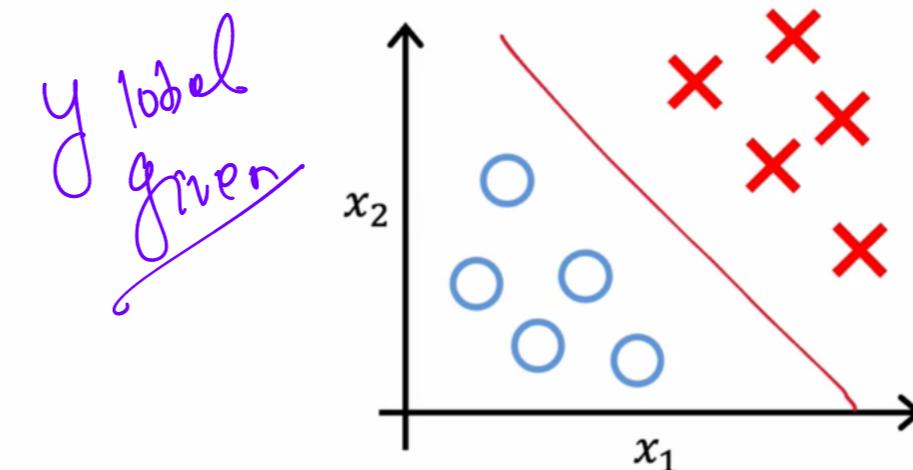
Would you like a worked-out example to see it in action? 😊

Allows automatically deriving recommendations for don't have predefined movie categories -

Supervised
→ Classification
→ Regression

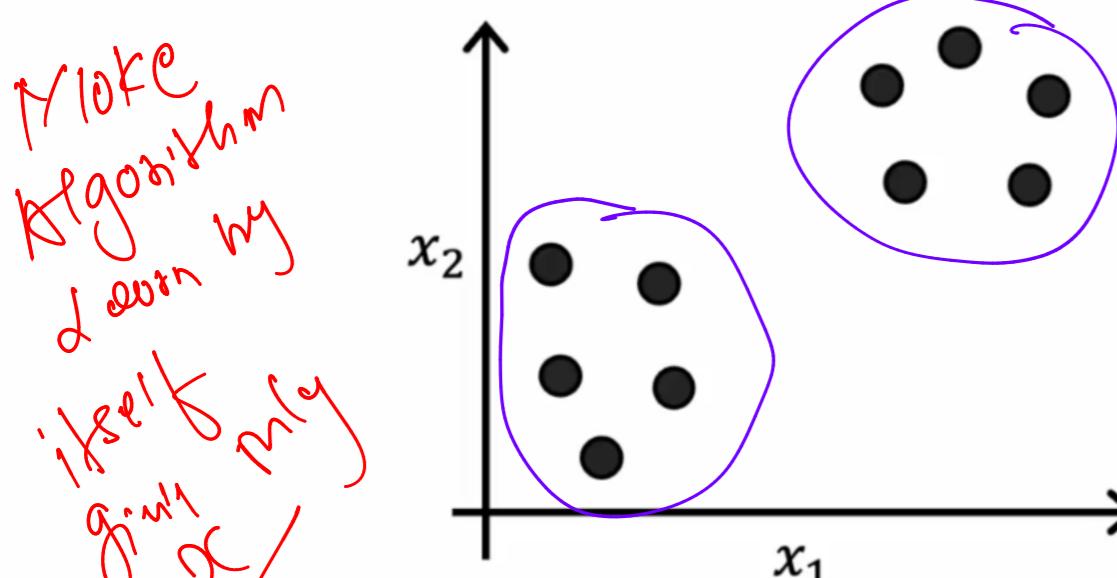
Teach the Algorithm
giving X and y

Supervised learning



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

Unsupervised learning



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

Doesn't have label y
Allow algorithm to learn by itself

- Clustering
- Applications of clustering
 - ① News Grouping
 - ② Market Segmentation (Trends)
 - ③ DNA Analysis
 - ④ Astronomical Data Analysis