

LECTURE 8

WHAT WE HAVE LEARNT?

- Regression
- Classification

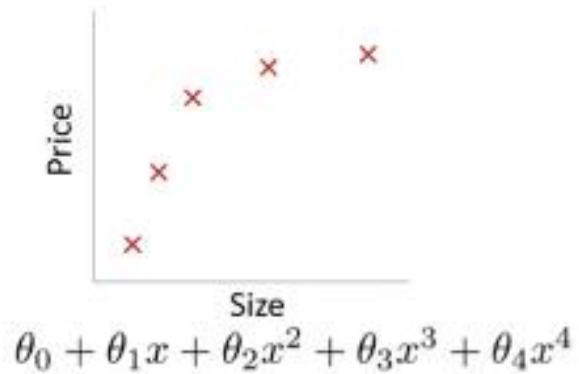
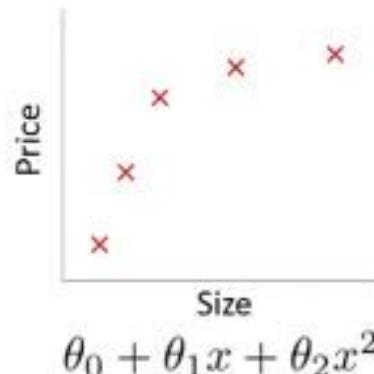
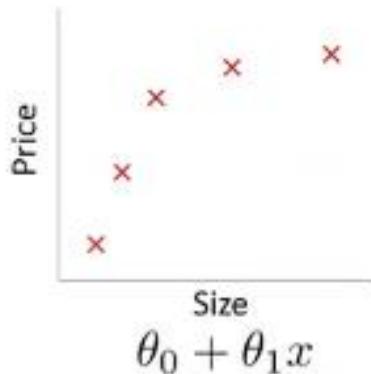
AGENDA FOR TODAY

- Problem of Overfitting & Regularization
- Deciding what to try next
- Model selection
- Bias vs Variance
- Learning Curves
- Performance Metrics
- Confusion Matrix
- Stochastic Gradient Descent and Mini-Batch Gradient Descent
- Unsupervised Learning

PROBLEM OF OVERFITTING

REGRESSION EXAMPLE

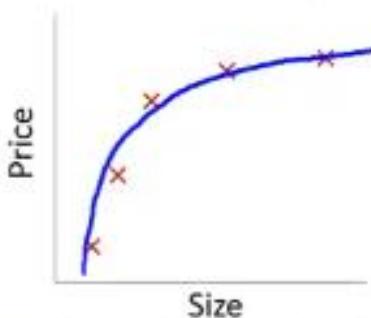
Example: Linear regression (housing prices)



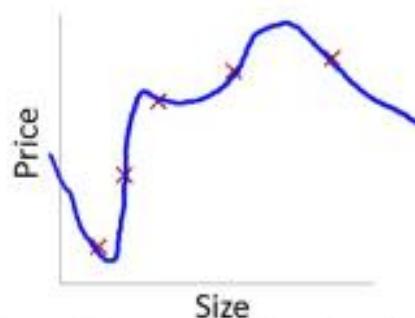
Example: Linear regression (housing prices)



$\rightarrow \theta_0 + \theta_1 x$
"Underfit" "High bias"



$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$
"Just right"

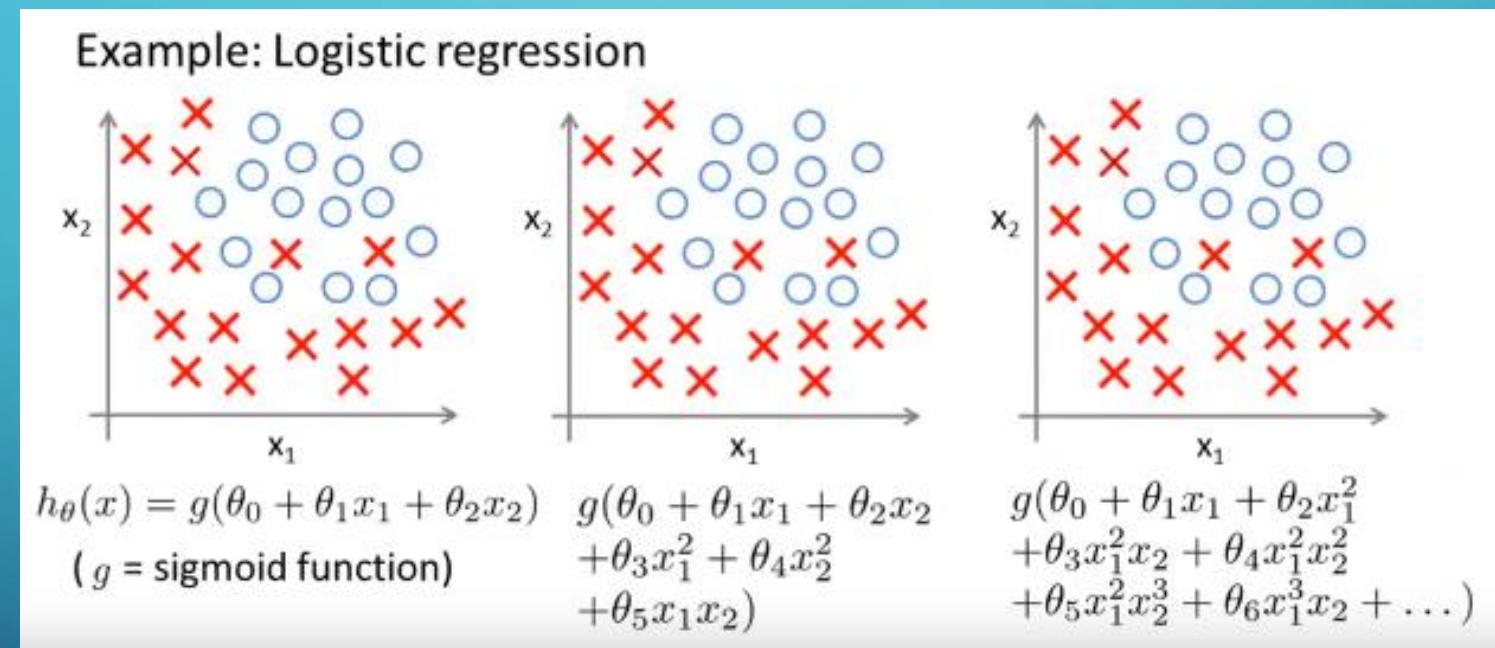


$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
"Overfit" "High variance"

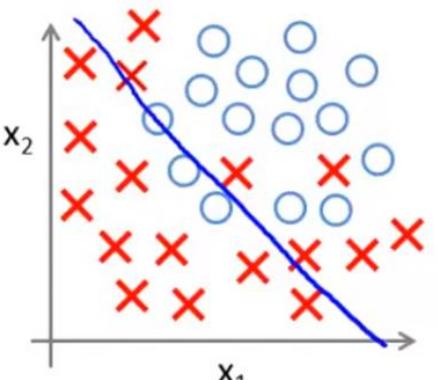
OVERFITTING

- If we have too many features, the learned hypothesis may fit the training set very well, but fail to generalise to new examples

CLASSIFICATION EXAMPLE



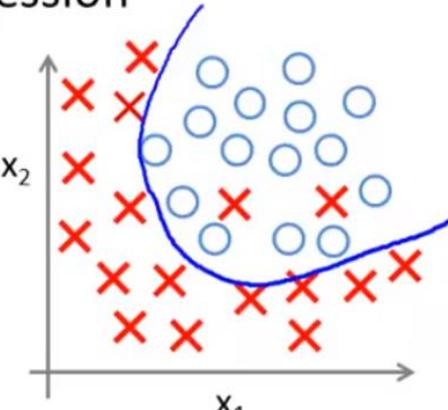
Example: Logistic regression



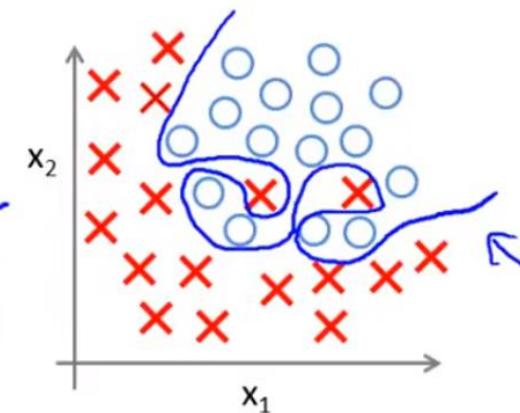
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

(g = sigmoid function)

↖ "Under-fit"



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 \underline{x_1 x_2})$$



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 \underline{x_1^2 x_2^3} + \theta_6 \underline{x_1^3 x_2} + \dots)$$

↖ "Over-fit"

Usually occurs due to high number of features

Addressing overfitting:

x_1 = size of house

x_2 = no. of bedrooms

x_3 = no. of floors

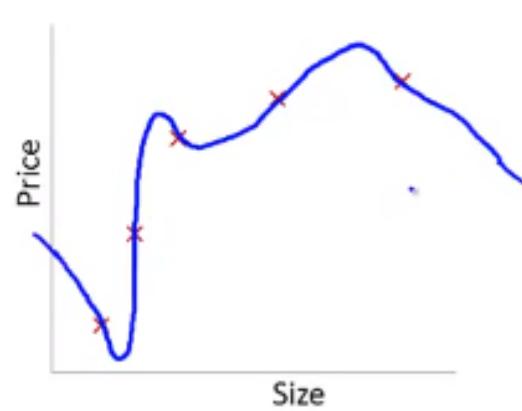
x_4 = age of house

x_5 = average income in neighborhood

x_6 = kitchen size

:

x_{100}

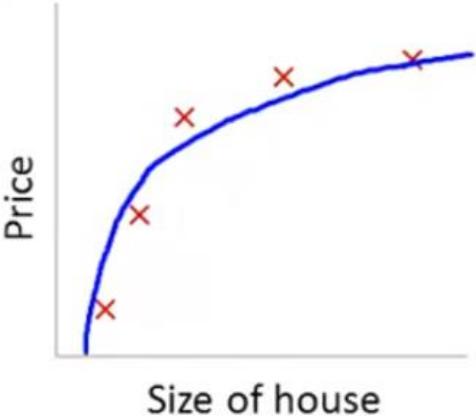


ADDRESSING OVERFITTING

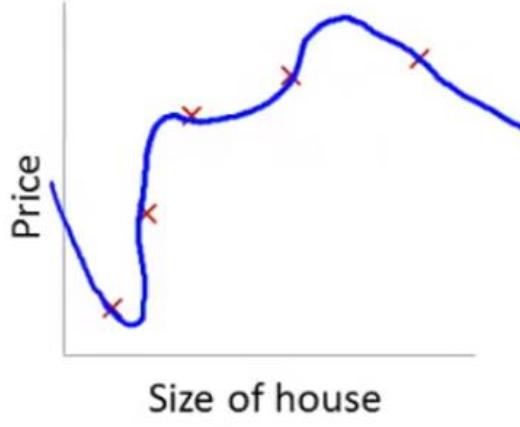
- Option1: Reduce number of features
 - ✓ Usually not advisable. Some complex techniques can be used for feature engineering
- Option2: Regularization
 - ✓ Keep all the features, but reduce magnitude/values of parameters
 - ✓ Works well when we have a lot of features, each of which contributes a bit to predicting output

REGULARIZATION AND COST FUNCTION

Intuition



$$\theta_0 + \theta_1x + \theta_2x^2$$



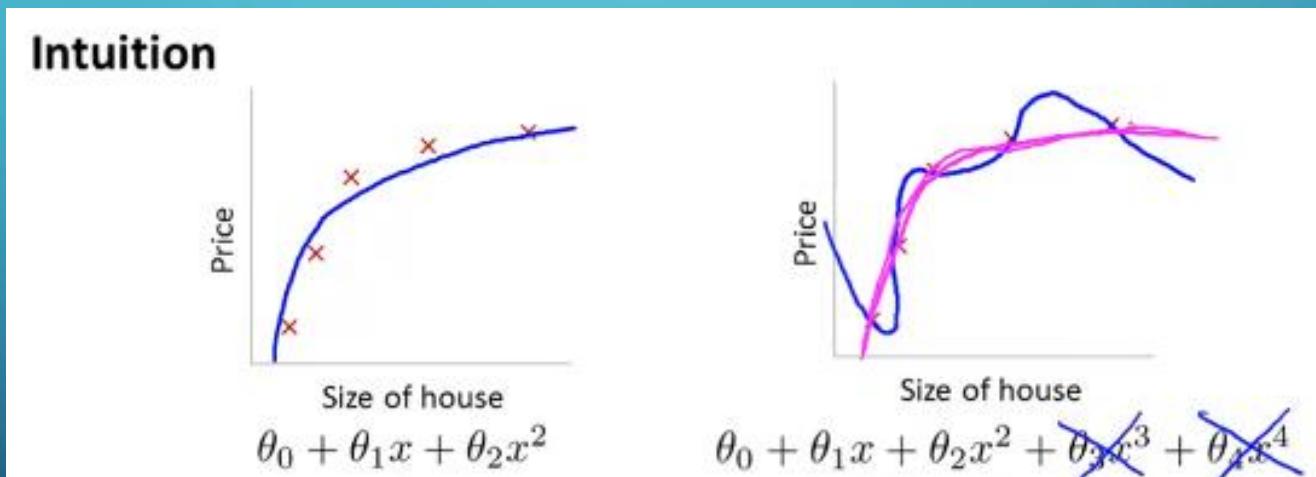
$$\theta_0 + \theta_1x + \theta_2x^2 + \theta_3x^3 + \theta_4x^4$$

Suppose we penalize and make θ_3, θ_4 really small.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \theta_3^2 + 1000 \theta_4^2$$

- Then θ_3 and θ_4 will have to be reduced to almost zero in order to minimise this function



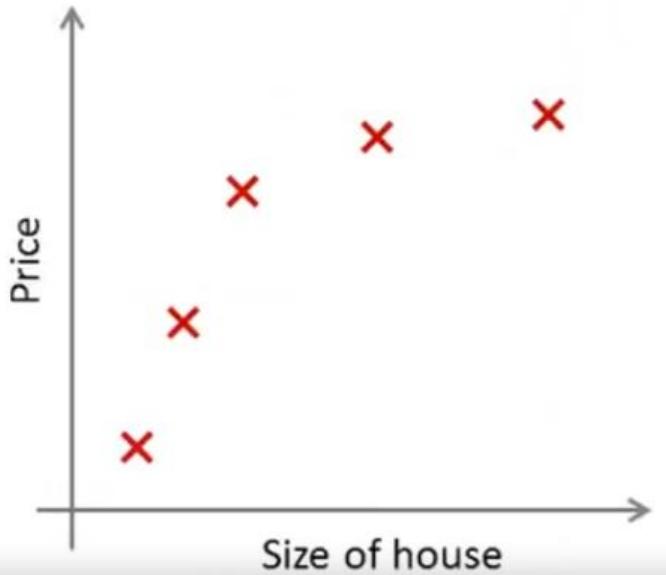
REGULARIZATION

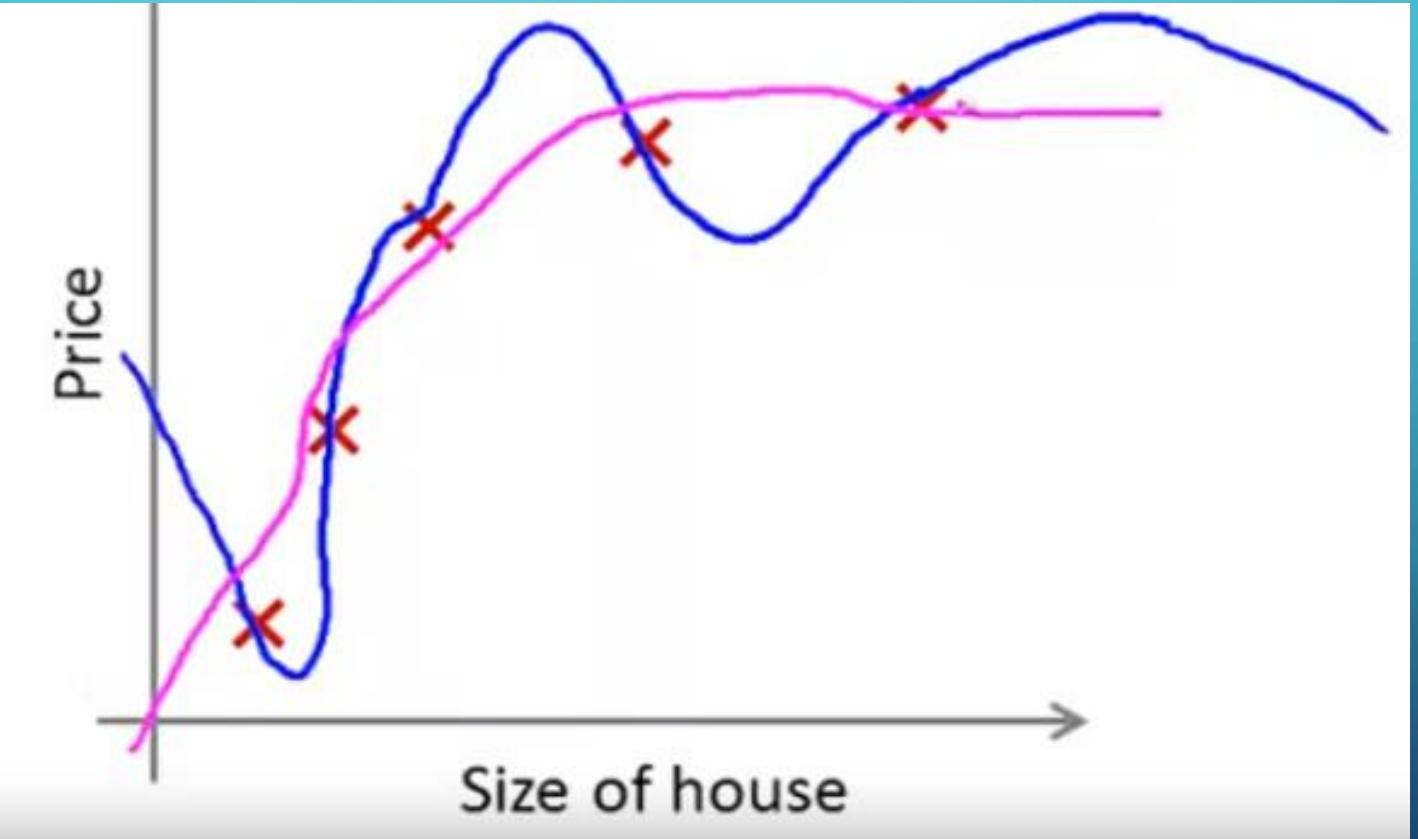
- “Simpler” hypothesis
- Less prone to overfitting

Regularization.

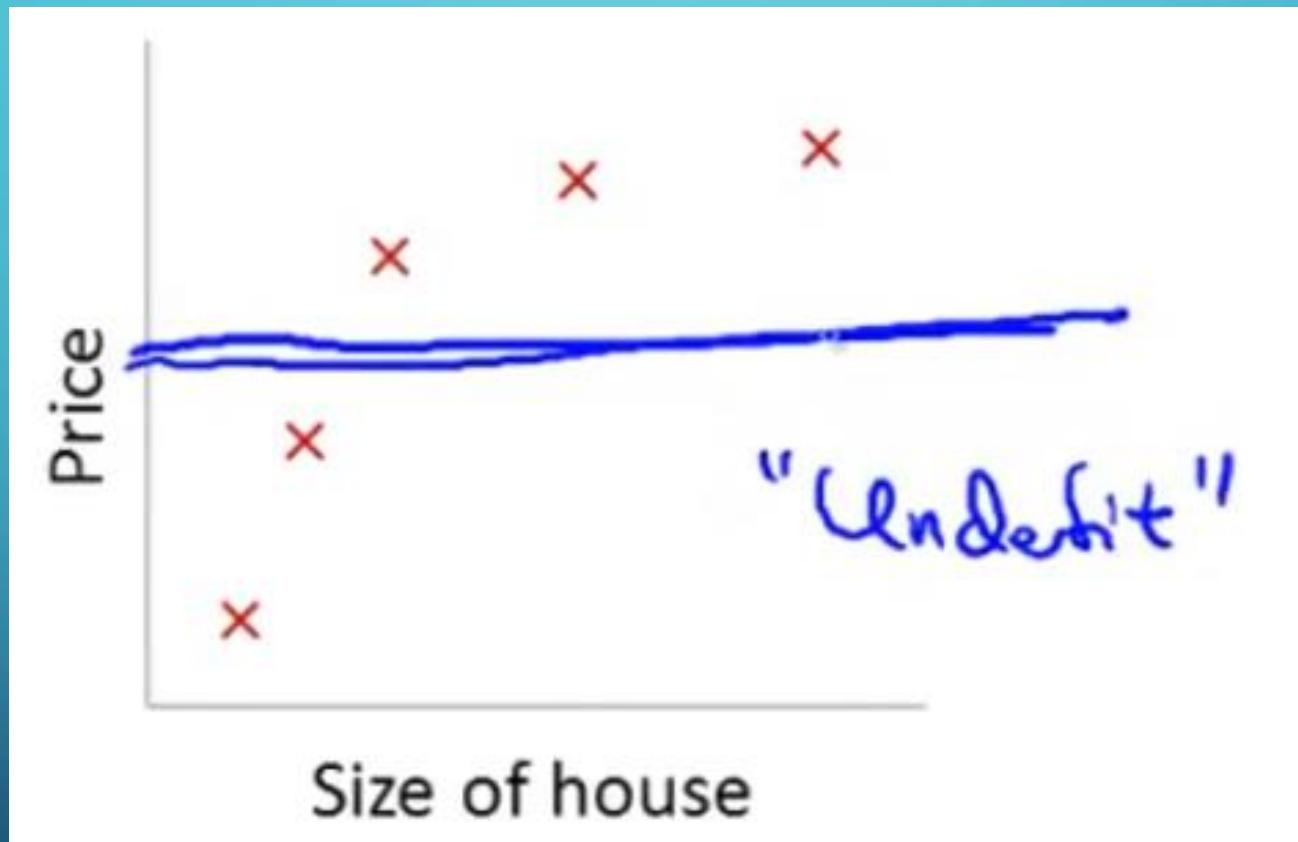
$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

$$\min_{\theta} J(\theta)$$





$h\theta(x) = \theta_0$ and other parameters are reduced to zero



REGULARIZATION IN LINEAR REGRESSION

COST FUNCTION

Regularized linear regression

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$
$$\min_{\theta} J(\theta)$$

GRADIENT DESCENT

EARLIER:

Gradient descent

Repeat {

$$\theta_j := \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (j = 0, 1, 2, 3, \dots, n)$$

}

NOW: STEP 1

Gradient descent

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\rightarrow \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (j = \textcolor{red}{X}, 1, 2, 3, \dots, n)$$

}

NOW: STEP2

Repeat {

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\rightarrow \theta_j := \theta_j - \alpha \left[\underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{(j = \cancel{0}, 1, 2, 3, \dots, n)} + \underbrace{\frac{\lambda}{m} \theta_j}_{\text{regularized}} \right]$$

}

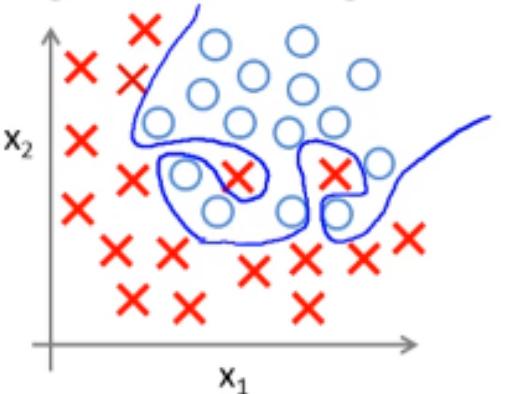
$$\frac{\partial}{\partial \theta_j} \underline{J(\theta)}$$

NOW: STEP3

$$\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

REGULARIZATION IN LOGISTIC REGRESSION

Regularized logistic regression.

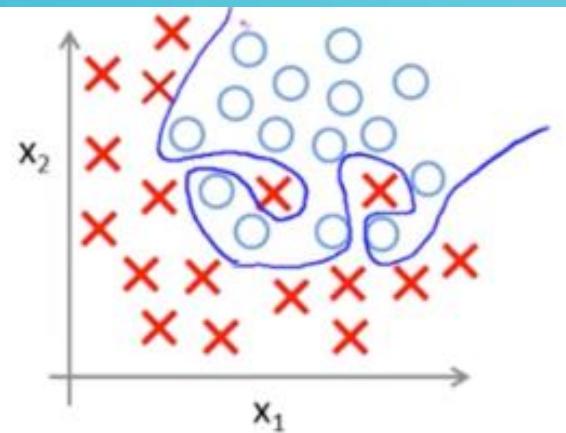


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

Cost function:

$$\rightarrow J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

COST FUNCTION



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

Cost function:

$$\Rightarrow J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

$\theta_1, \theta_2, \dots, \theta_n$

GRADIENT DESCENT

- As we saw earlier, it is exactly same Linear Regression:

$$\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

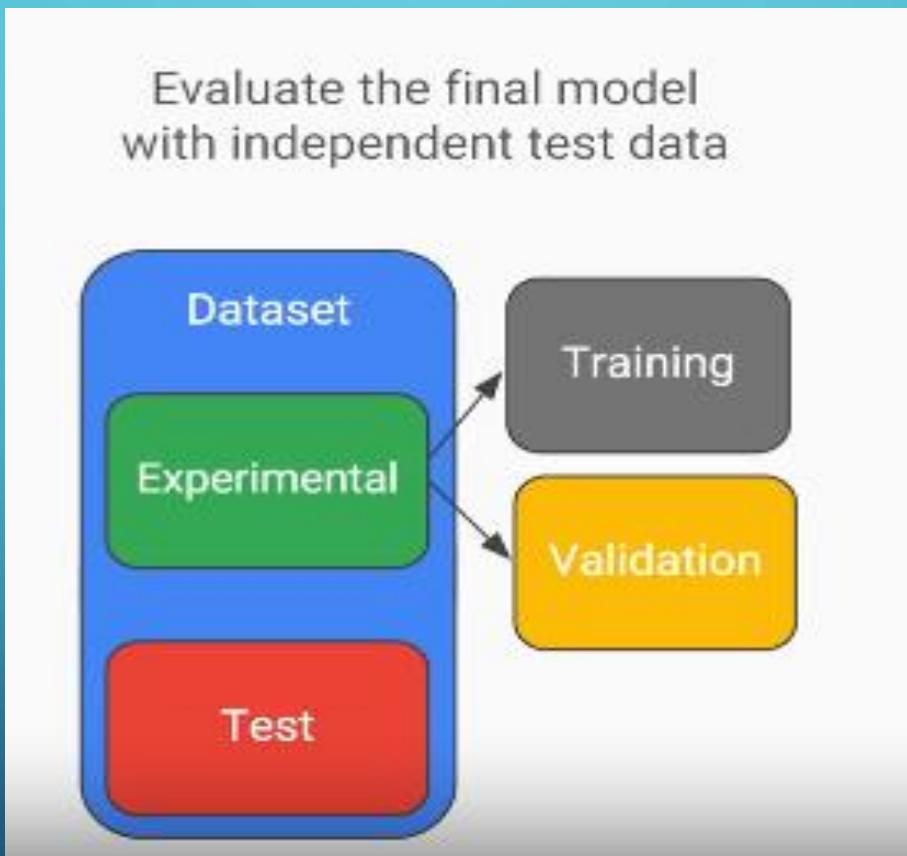
DECIDING WHAT TO TRY NEXT

- Suppose you have implemented linear regression to predict housing prices
- But you find out that it makes unacceptably large errors in its predictions. What to try next?
- Options:
 1. Get more training examples
 2. Try smaller set of features
 3. Try getting additional features
 4. Try adding polynomial features
 5. Increase/Decrease Regularization term

MACHINE LEARNING DIAGNOSTIC

- A test that you can run to gain insight what is/isn't working with a learning algorithm, and gain guidance as to how best to improve its performance
- Diagnostic can take time to implement, but doing so can be a very good use of your time

REVISIT DATA SPLIT



MODEL SELECTION

Model selection

1. $h_{\theta}(x) = \theta_0 + \theta_1 x$
2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \cdots + \theta_3 x^3$
⋮
10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10}$

Train/validation/test error

Training error:

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Cross Validation error:

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

Test error:

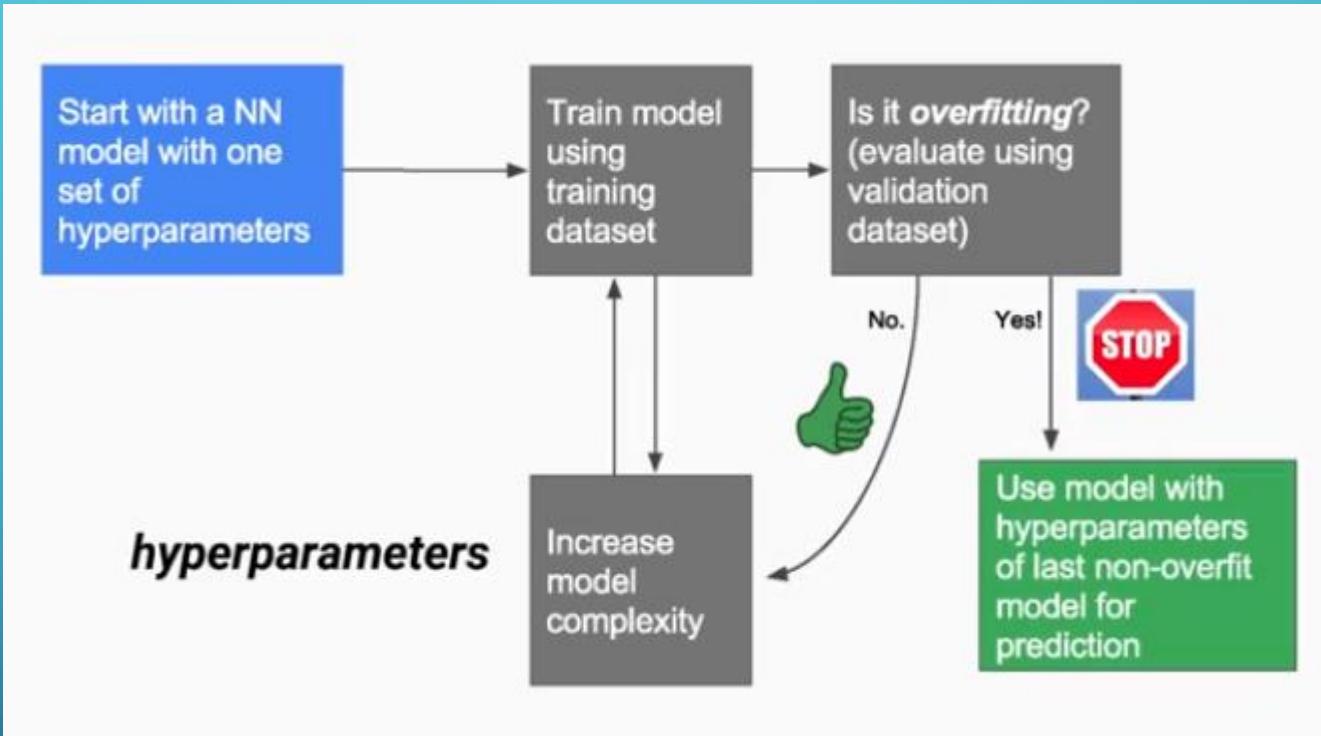
$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_\theta(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

Model selection

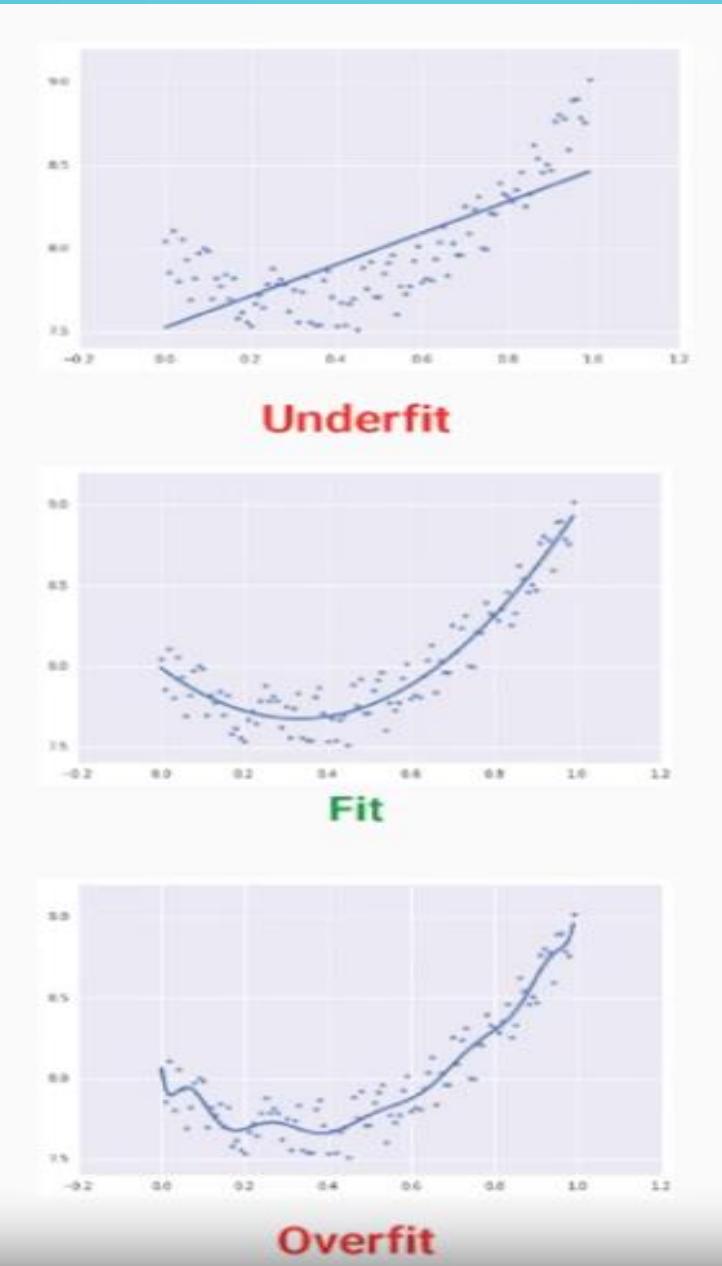
1. $h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)}$
⋮
 $J_{cv}(\theta^{(4)})$
10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)} \rightarrow J_{cv}(\theta^{(10)})$

Pick $\theta_0 + \theta_1 x_1 + \dots + \theta_4 x^4 \leftarrow$

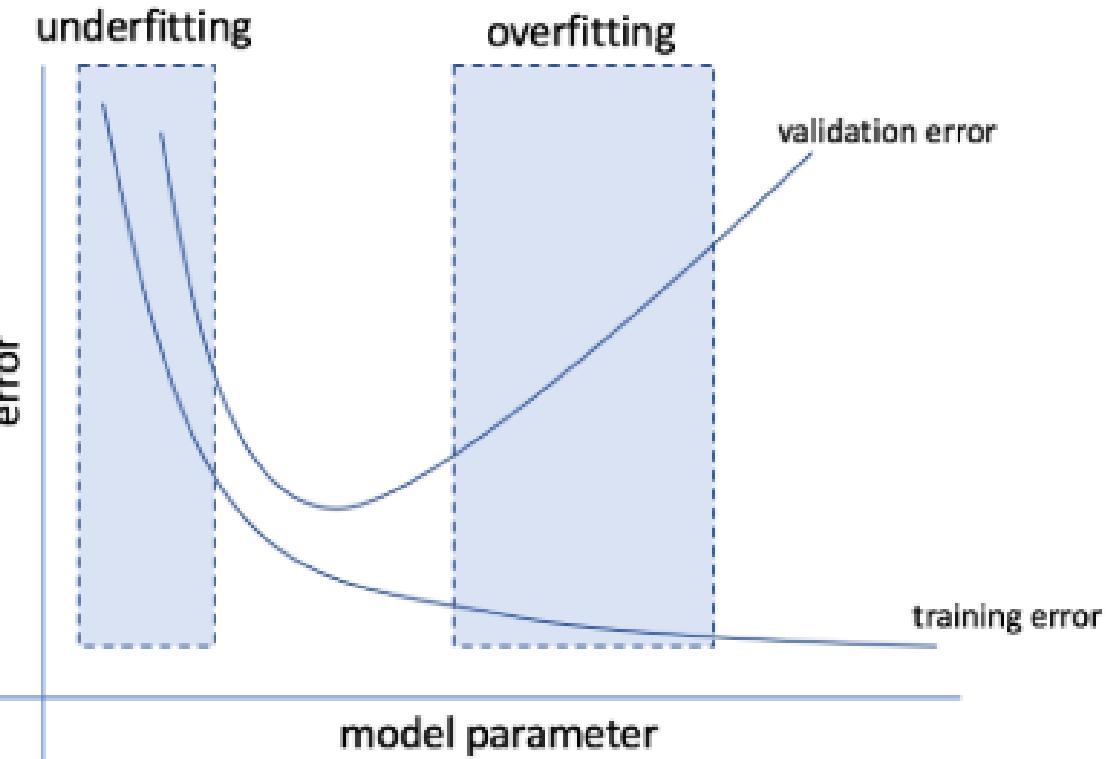
Estimate generalization error for test set $J_{test}(\theta^{(4)})$



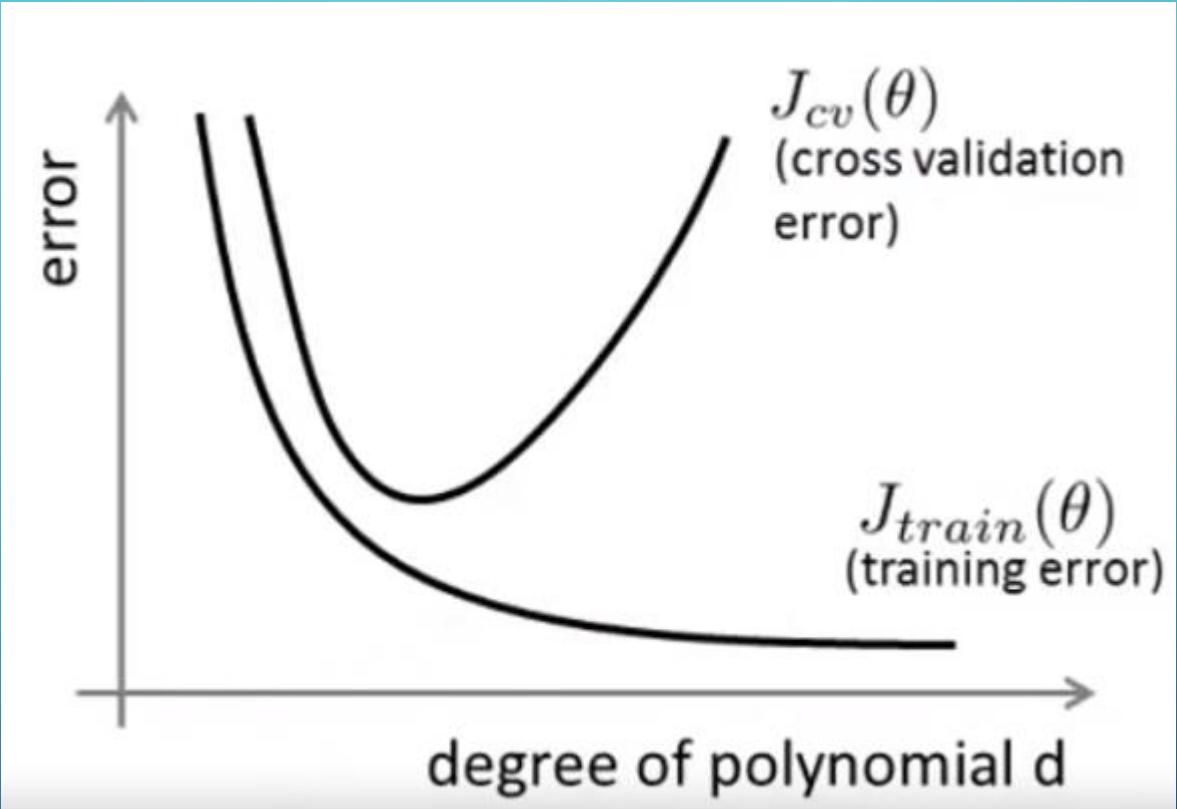
BIAS VS VARIANCE



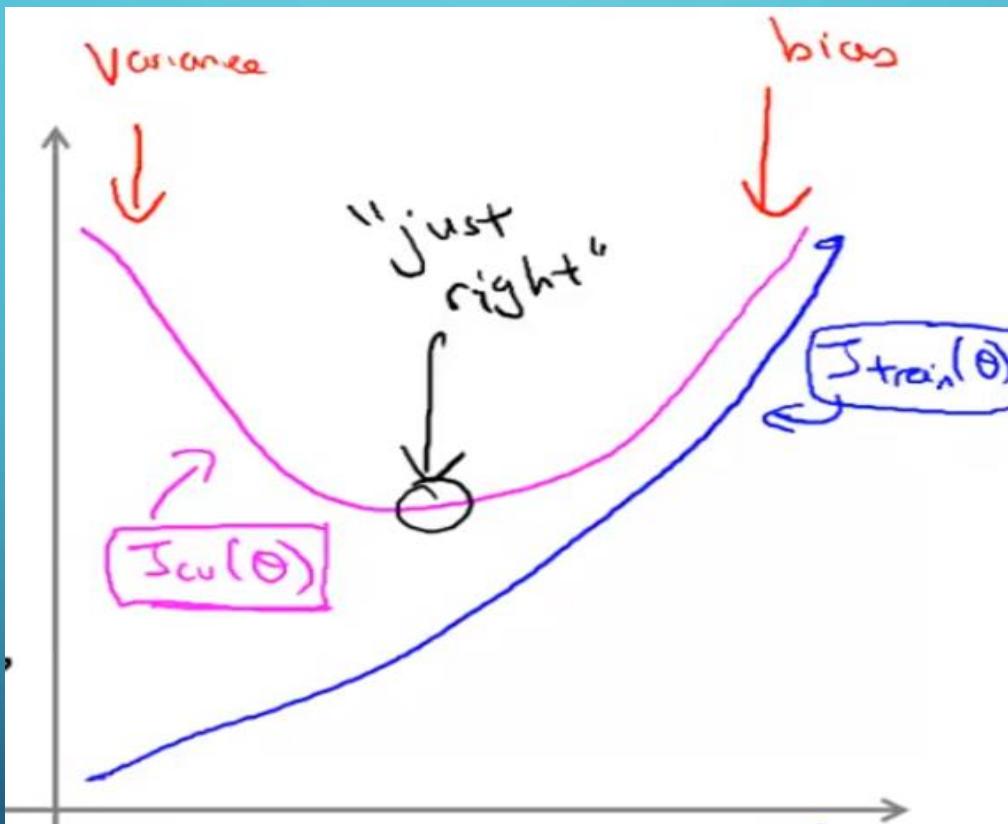
DEGREE POLYNOMIAL



Degree Polynomial can be viewed as Model parameter

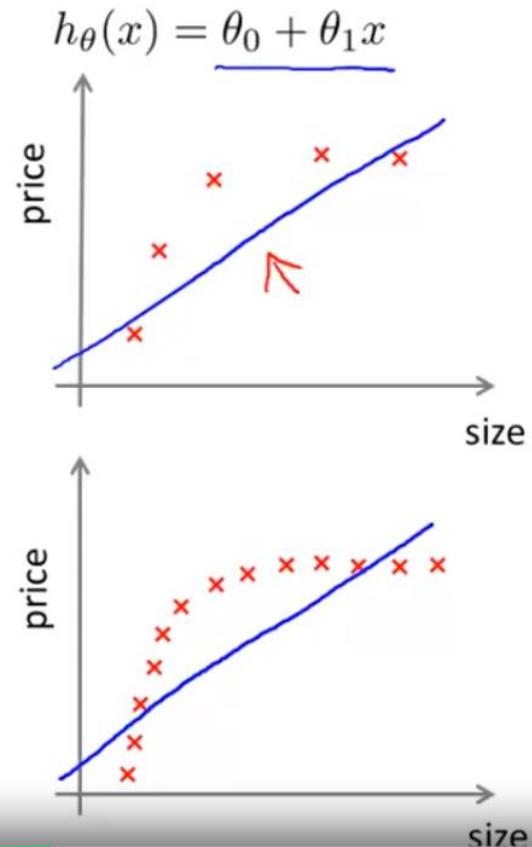
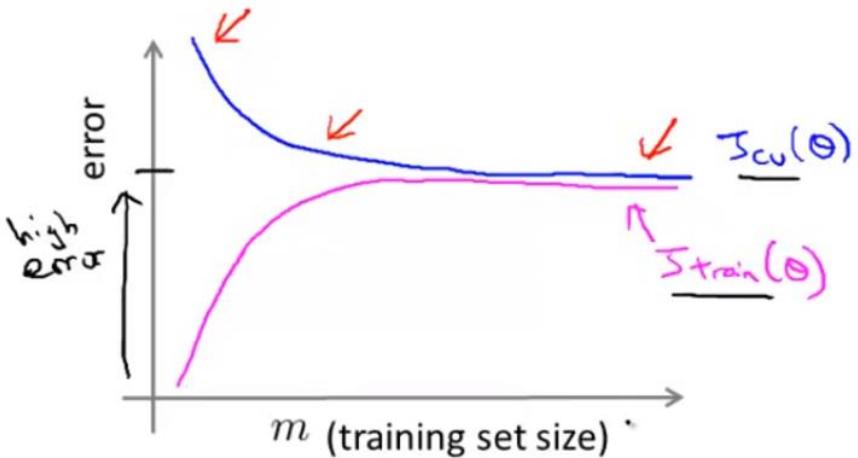


REGULARIZER



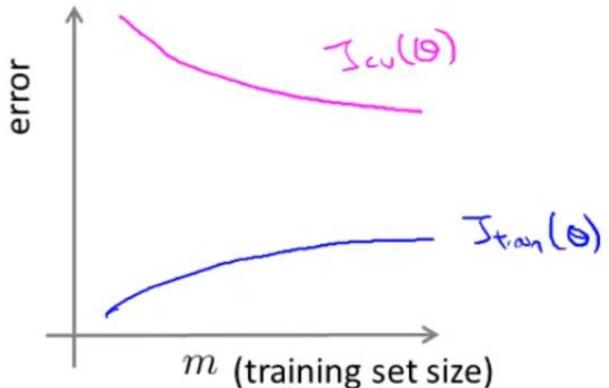
TRAINING SET SIZE

High bias



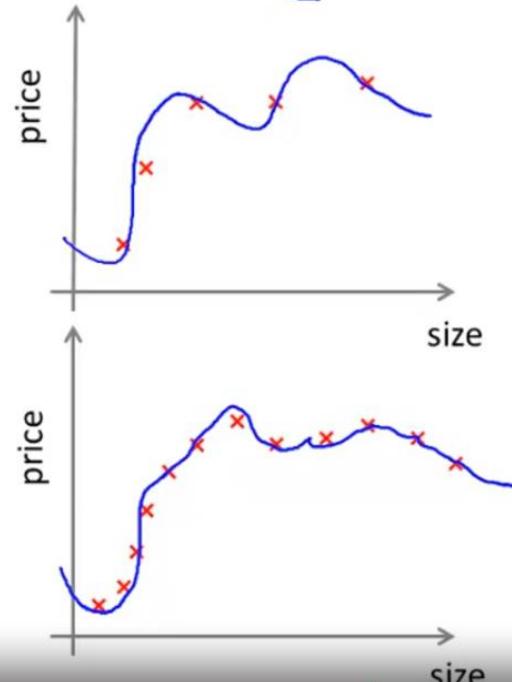
No use of getting more data

High variance



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{100}$$

(and small λ) ↗



Getting more data is likely to help

TIPS TO RESOLVE ML PROBLEMS

Underfitting (Highly bias)	Overfitting (High Variance)
Try getting additional features	Get more training examples
Try adding polynomial features	Try smaller set of features
Try decreasing λ	Try increasing λ

PERFORMANCE METRICS

PROBLEM WITH LOSS FUNCTION



1000 parking spaces
990 of them are **taken**
10 are **available**

An ML model that
always reported that a
space was occupied
would be right 99/100
times.

PERFORMANCE METRICS > LOSS FUNCTIONS

Performance metrics allow us to measure what matters

Loss Functions	Performance Metrics
<ul style="list-style-type: none">• During training• Harder to understand• Indirectly connected to business goals	<ul style="list-style-type: none">• After training• Easier to understand• Directly connected to business goals

INTUITION OF CONFUSION MATRIX

Use a confusion matrix to assess classification model performance

		Model Predictions	
		Positive	Negative
Labels	Positive	True Positives (TP) 	False Negatives (FN) <i>Type II Error</i> 
	Negative	False Positives (FP) <i>Type I Error</i> 	True Negatives (TN) 

Model says Yes

Model says No

PRECISION

Precision: true positives / total classified as positive

		Model Predictions	
		Positive	Negative
References	Positive	Available parking space exists Model predicts it is available	Available parking space exists Model doesn't predict it
	Negative	Available parking space doesn't exist Model predicts it is available	Available parking space doesn't exist Model correctly doesn't predict it
Precision		True Positives	False Negatives <i>Type II Error</i>
		False Positives <i>Type I error</i>	True Negatives

RECALL

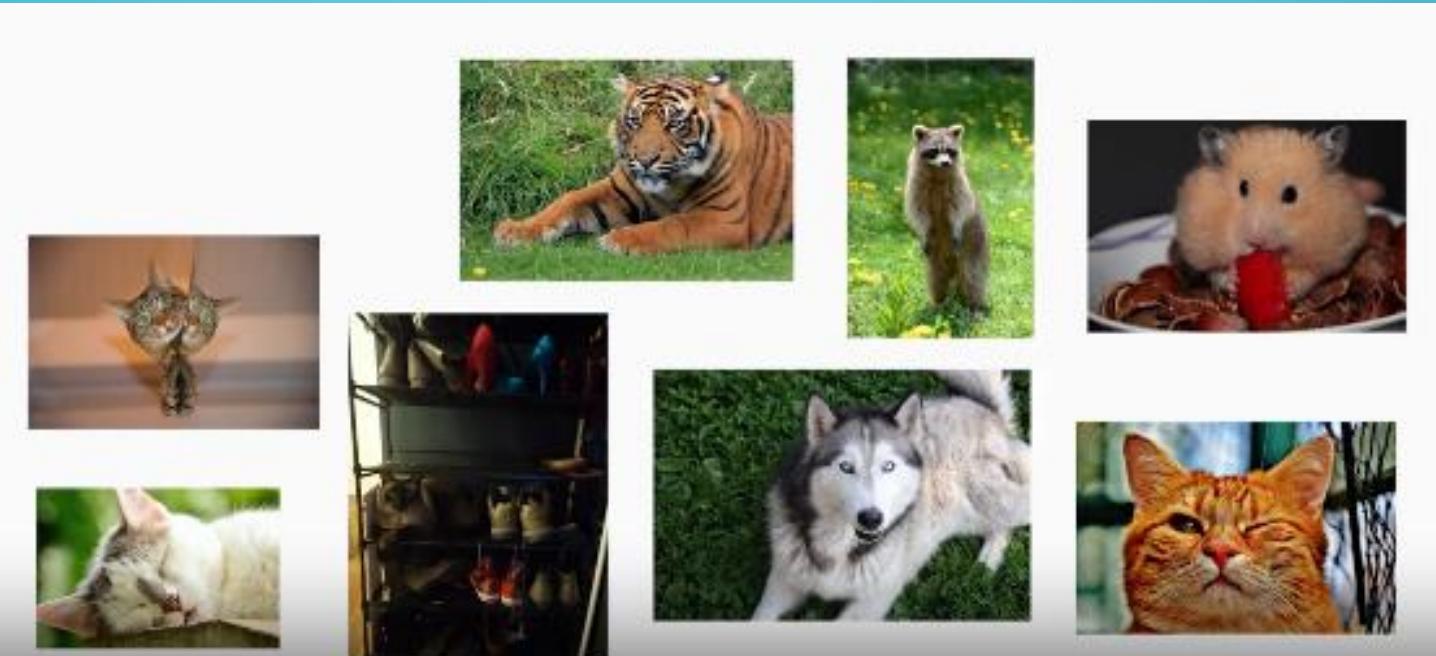
Recall: true positives / all actual positives in our reference

		Model Predictions	
		Positive	Negative
References	Negative	Available parking space exists Model predicts it is available	Available parking space exists Model doesn't predict it
	Positive	True Positives	False Negatives <i>Type II Error</i>
	Negative	Available parking space doesn't exist Model predicts it is available	Available parking space doesn't exist Model correctly doesn't predict it
	Positive	False Positives <i>Type I Error</i>	True Negatives

Recall

PRACTICE TIME

IDENTIFY CATS IN THE IMAGE AS CAT / NOT CAT





HYPOTHETICAL RESULT FROM OUR ML MODEL



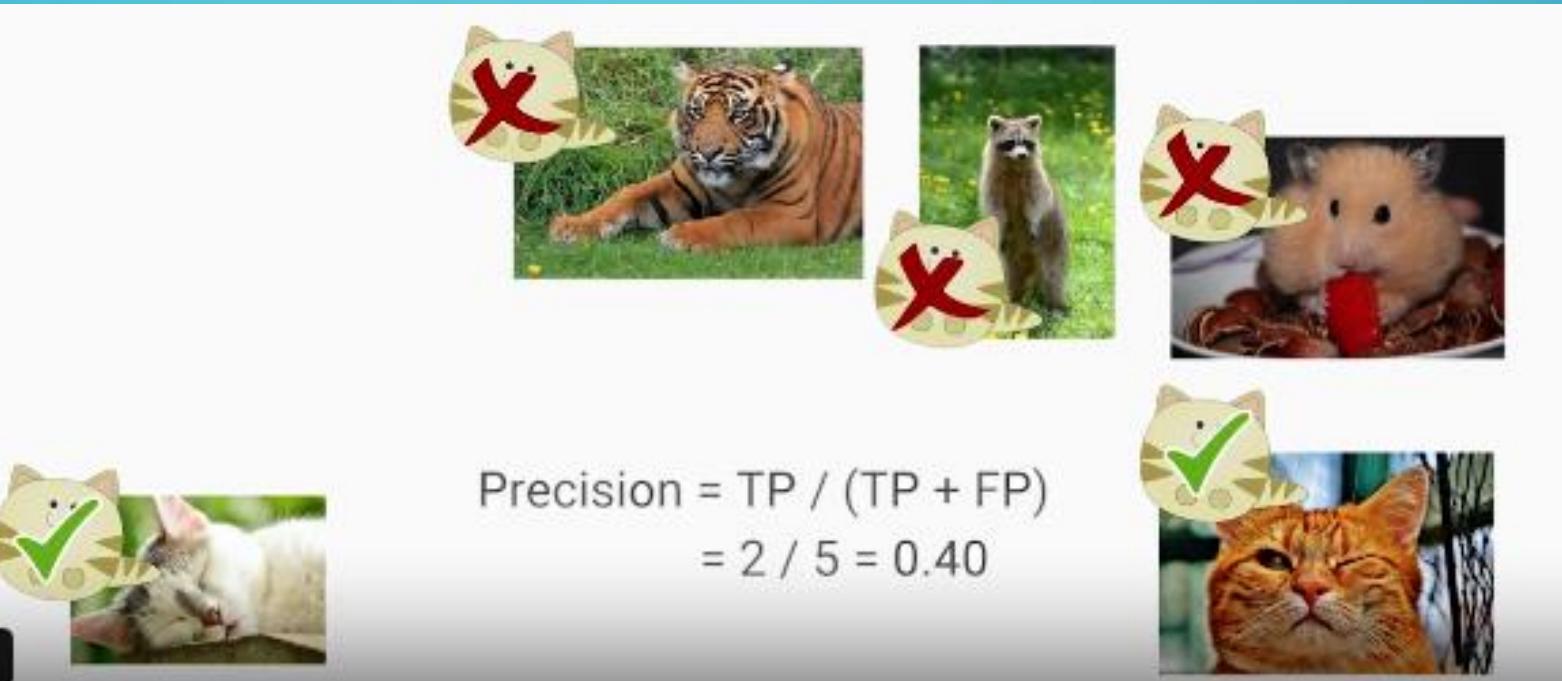
COMPARISON



ACCURACY



PRECISION



RECALL

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN}) \\ = 2 / 4 = 0.50$$



SUMMARY

- **Precision:** Makes sure whatever model predicts is correct, even if that means it misses out on some positives
- **Recall:** Makes sure that model correctly predicts maximum number of positives, even if that means it has more false positives.

TRADE OFF BETWEEN PRECISION AND RECALL

How to compare precision/recall numbers?

	Precision(P)	Recall (R)
Algorithm 1	0.5	0.4
Algorithm 2	0.7	0.1
Algorithm 3	0.02	1.0

AVERAGE WONT WORK!!

	Precision(P)	Recall (R)	Average
Algorithm 1	<u>0.5</u>	<u>0.4</u>	0.45
Algorithm 2	<u>0.7</u>	<u>0.1</u>	0.4
Algorithm 3	0.02	1.0	0.51

F1 SCORE: WHAT AN IDEA SIRJEE!!!

- F1 score is given by:
$$\frac{2PR}{P+R}$$
..... where P=Precision, R=Recall

	Precision(P)	Recall (R)	Average	F ₁ Score
Algorithm 1	<u>0.5</u>	<u>0.4</u>	0.45	0.444
Algorithm 2	<u>0.7</u>	<u>0.1</u>	0.4	0.175
Algorithm 3	<u>0.02</u>	1.0	0.51	0.0392

TYPES OF GRADIENT DESCENT

- Batch Gradient Descent
- Mini-Batch Gradient Descent
- Stochastic Gradient Descent

WHAT WE KNOW?

- We have already seen the Batch Gradient Descent

BATCH GRADIENT DESCENT

Linear regression with gradient descent

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

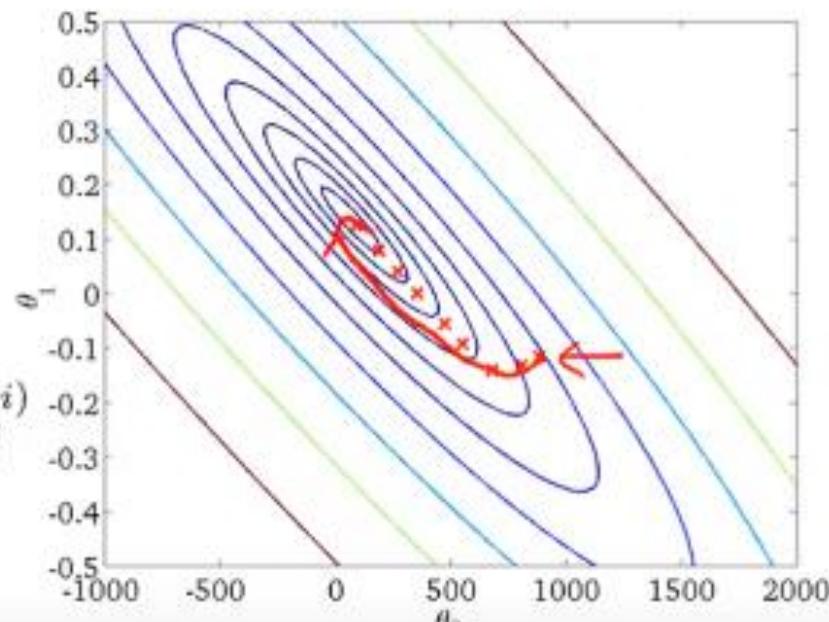
$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(for every $j = 0, \dots, n$)

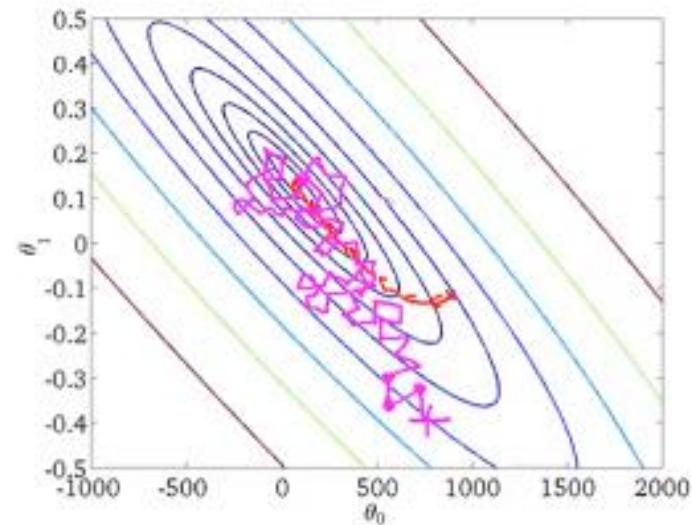
}



STOCHASTIC GRADIENT DESCENT

Stochastic gradient descent

1. Randomly shuffle (reorder) training examples
2. Repeat {
 - for $i := 1, \dots, m\{$
 - $\rightarrow \theta_j := \theta_j - \alpha(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$
 - (for every $j = 0, \dots, n$)}



Batch gradient descent: Use all m examples in each iteration

Stochastic gradient descent: Use 1 example in each iteration

MINI-BATCH GRADIENT DESCENT

- Use b examples in each iteration ... where $b = \text{mini-batch size}$

Mini-batch gradient descent

Say $b = 10, m = 1000$.

Repeat {

 for $i = 1, 11, 21, 31, \dots, 991$ {

$$\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

 (for every $j = 0, \dots, n$)

 }

.

}

UNSUPERVISED LEARNING

APPLICATIONS

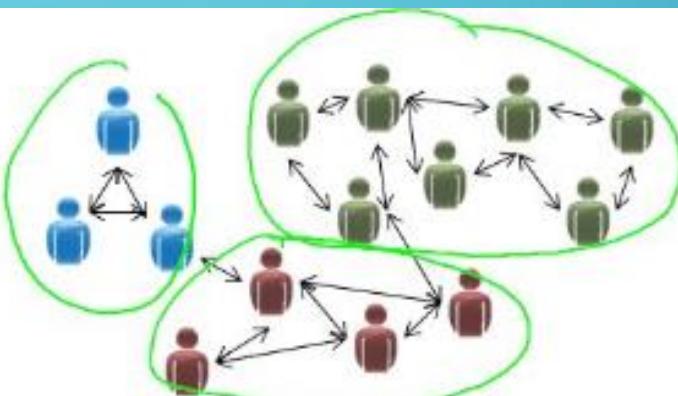
Applications of clustering



→ Market segmentation



→ Organize computing clusters

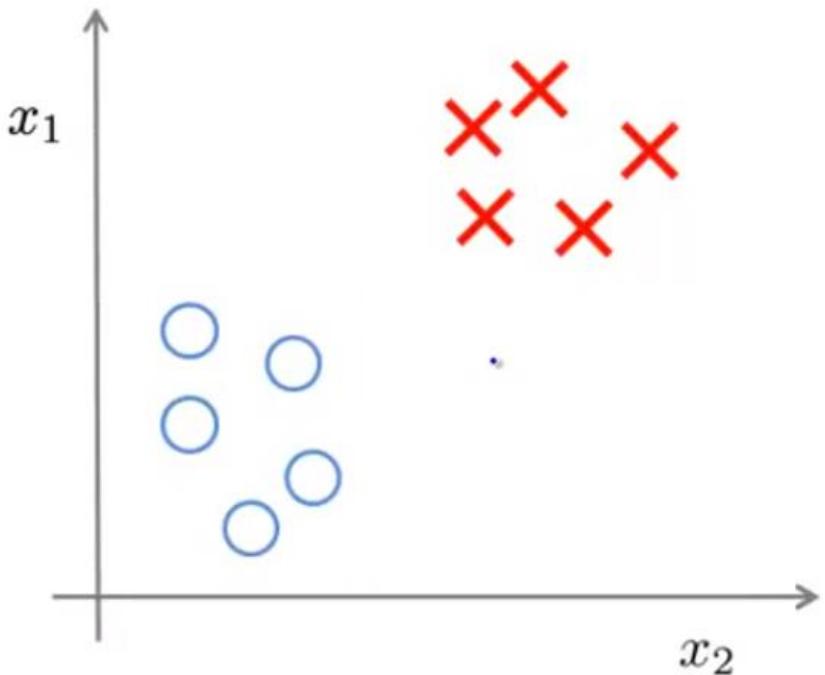


→ Social network analysis



→ Astronomical data analysis

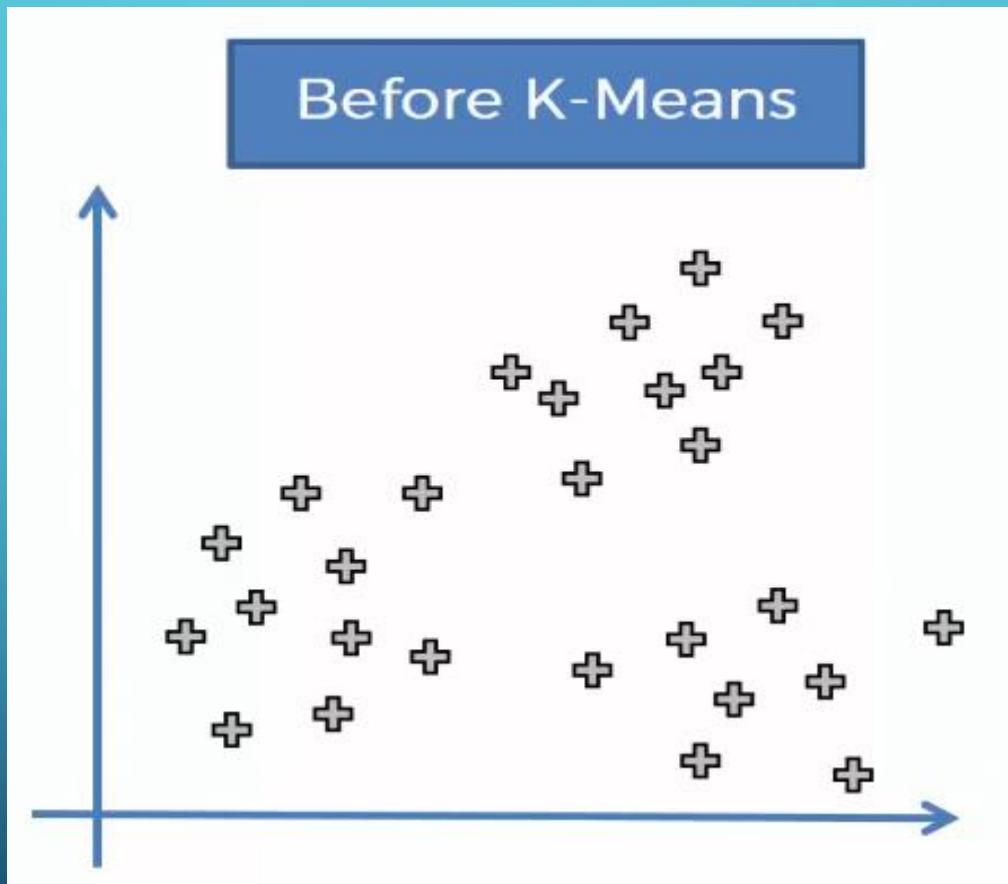
WHAT WE KNOW?



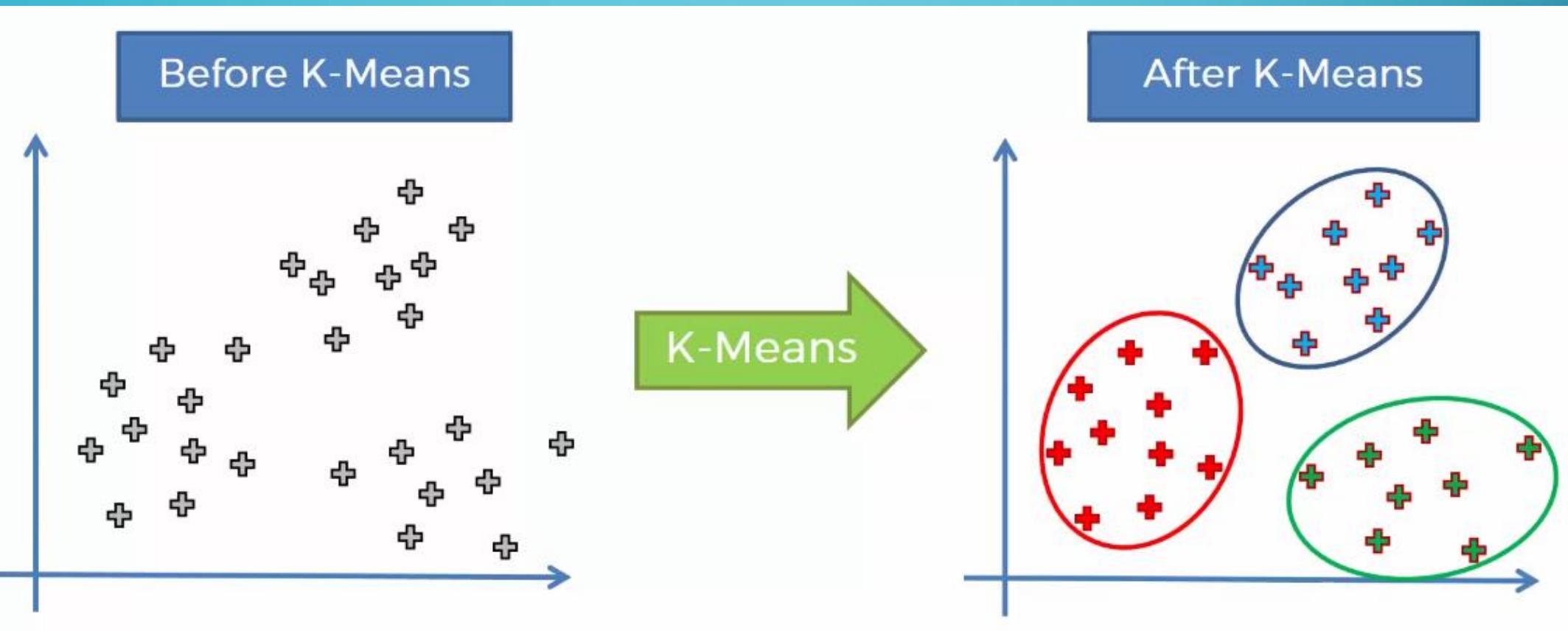
Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$

NO LABELS

Before K-Means



CLUSTERS



K-MEANS ALGORITHM

STEP 1: Choose the number K of clusters



STEP 2: Select at random K points, the centroids (not necessarily from your dataset)



STEP 3: Assign each data point to the closest centroid → That forms K clusters



STEP 4: Compute and place the new centroid of each cluster

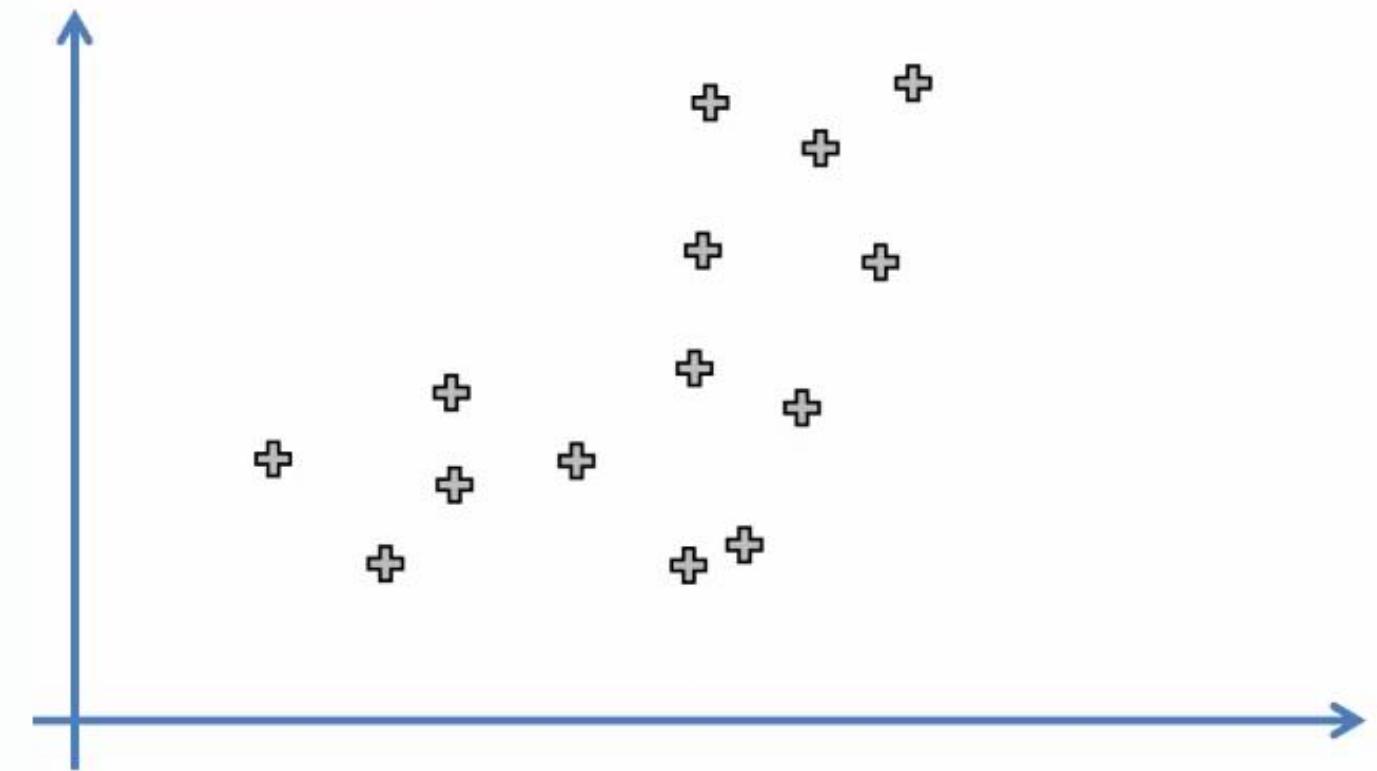


STEP 5: Reassign each data point to the new closest centroid.
If any reassignment took place, go to STEP 4, otherwise go to FIN.

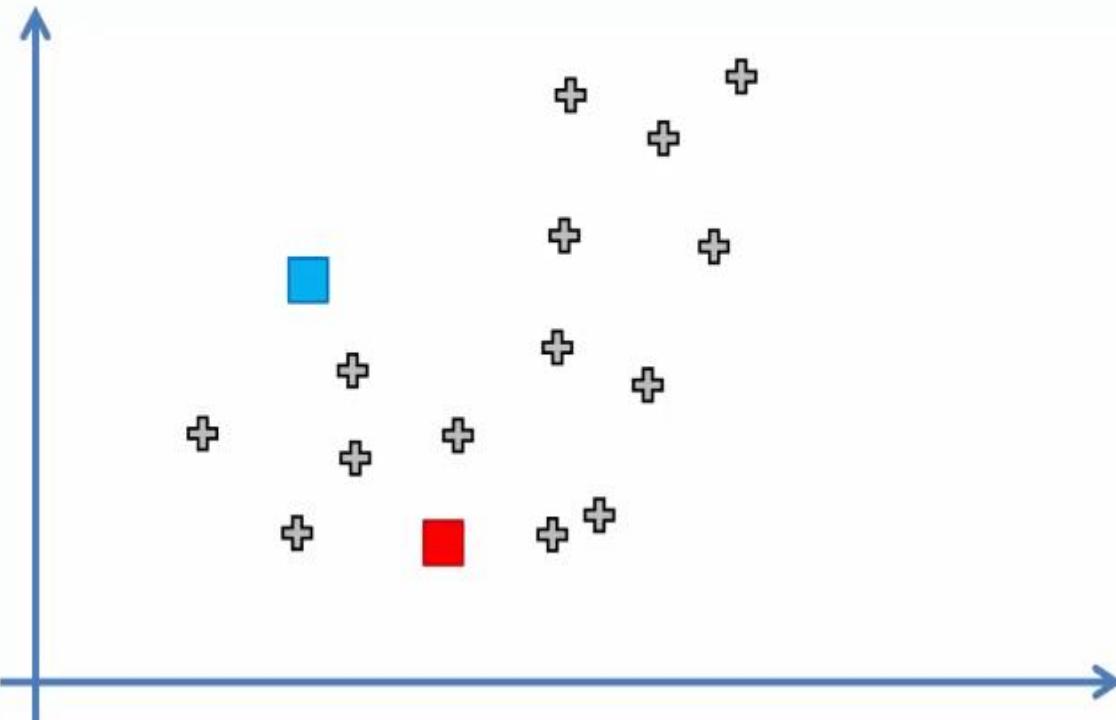


Your Model is Ready

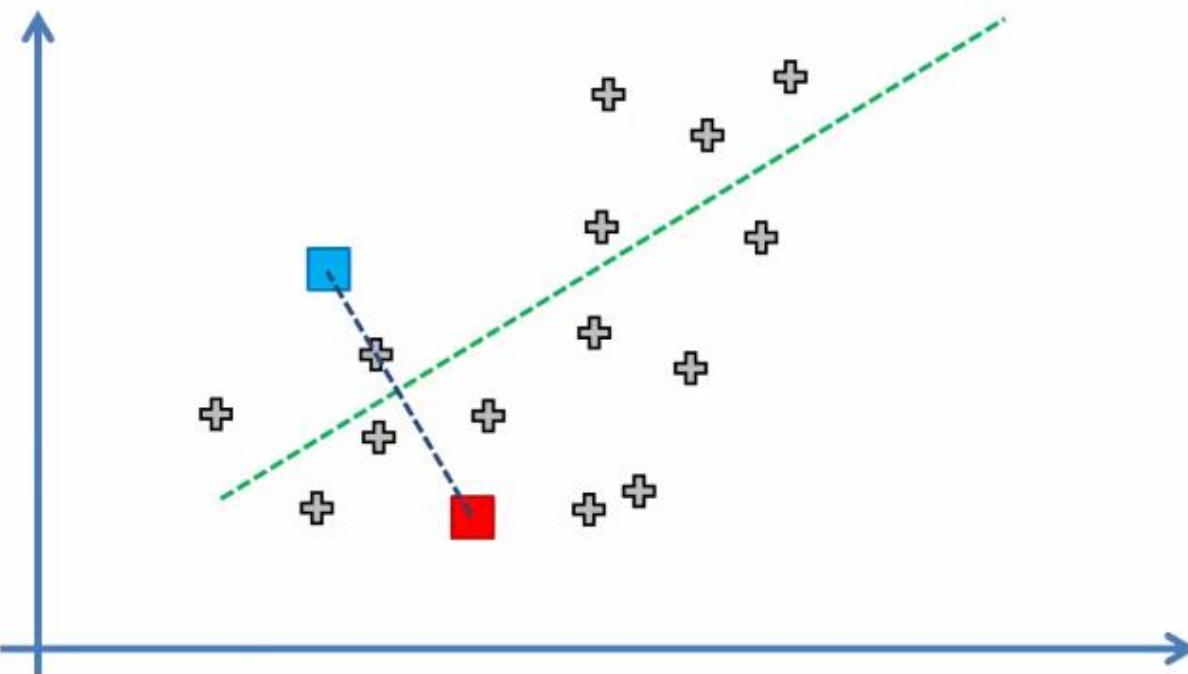
STEP 1: Choose the number K of clusters: $K = 2$



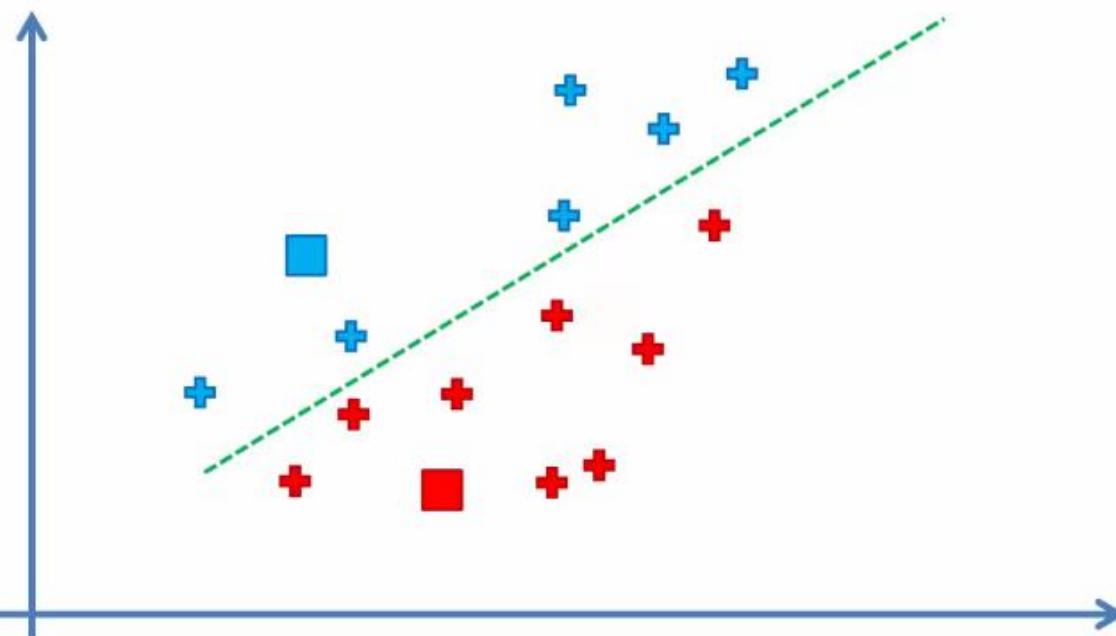
STEP 2: Select at random K points, the centroids (not necessarily from your dataset)



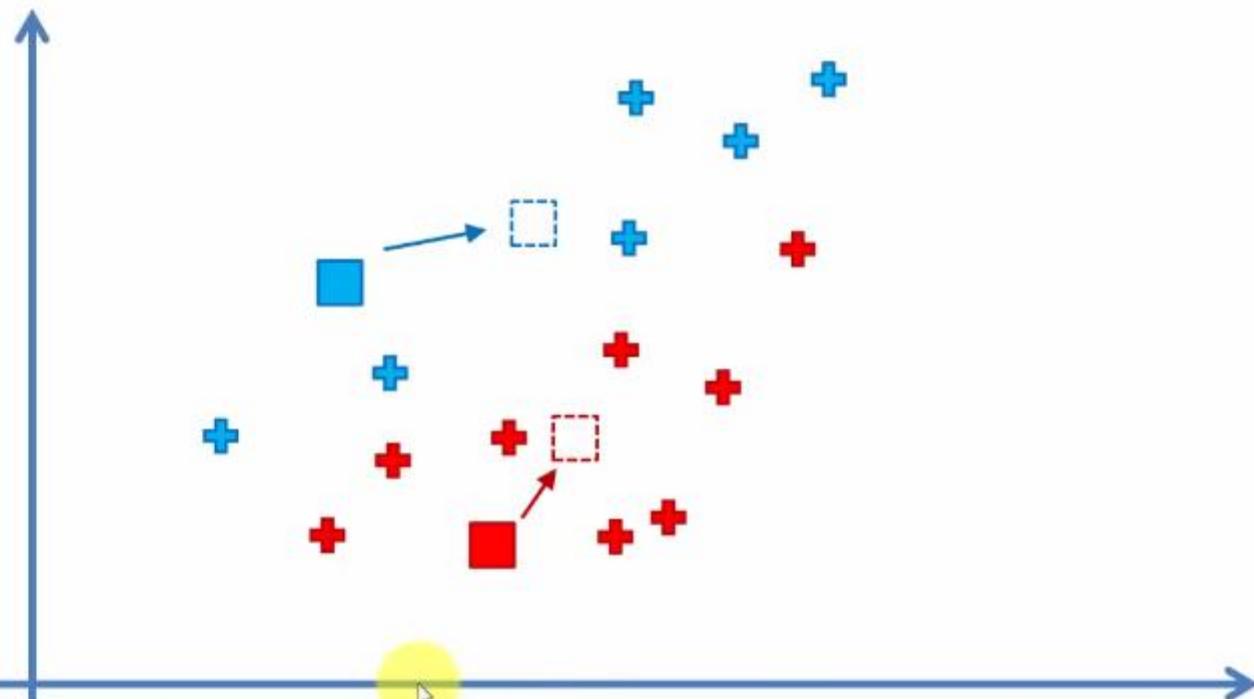
STEP 3: Assign each data point to the closest centroid \rightarrow That forms K clusters



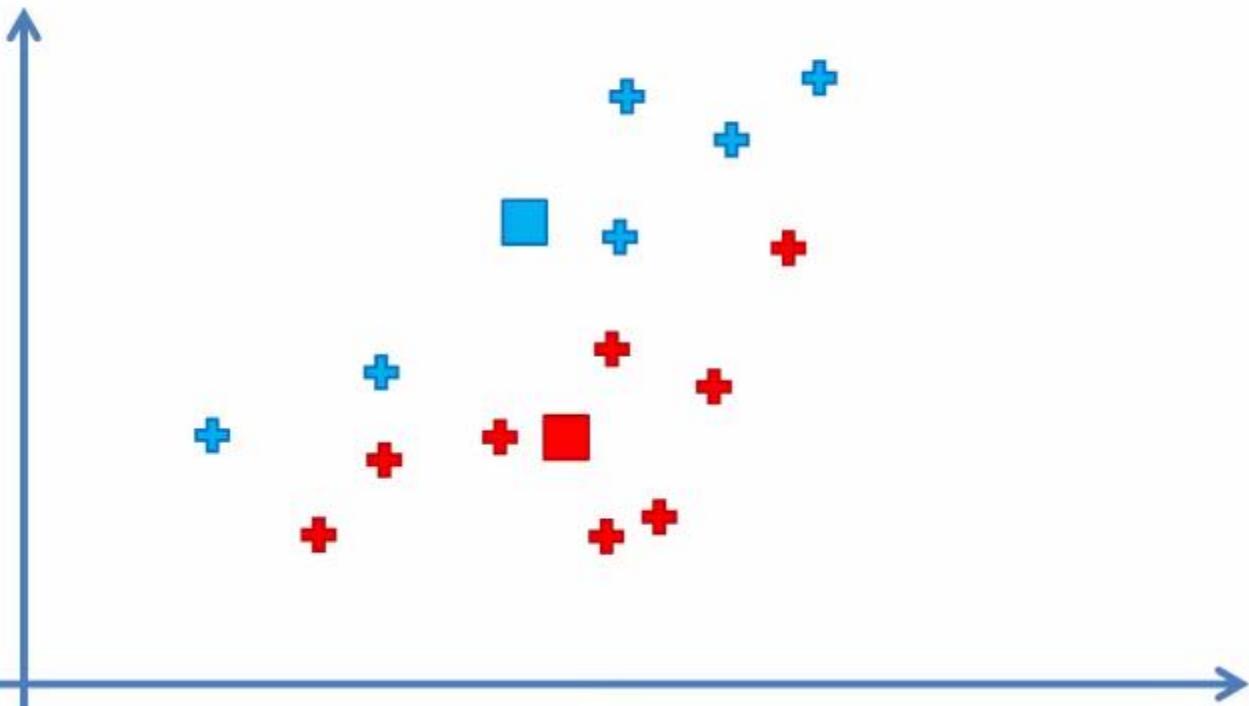
STEP 3: Assign each data point to the closest centroid \rightarrow That forms K clusters



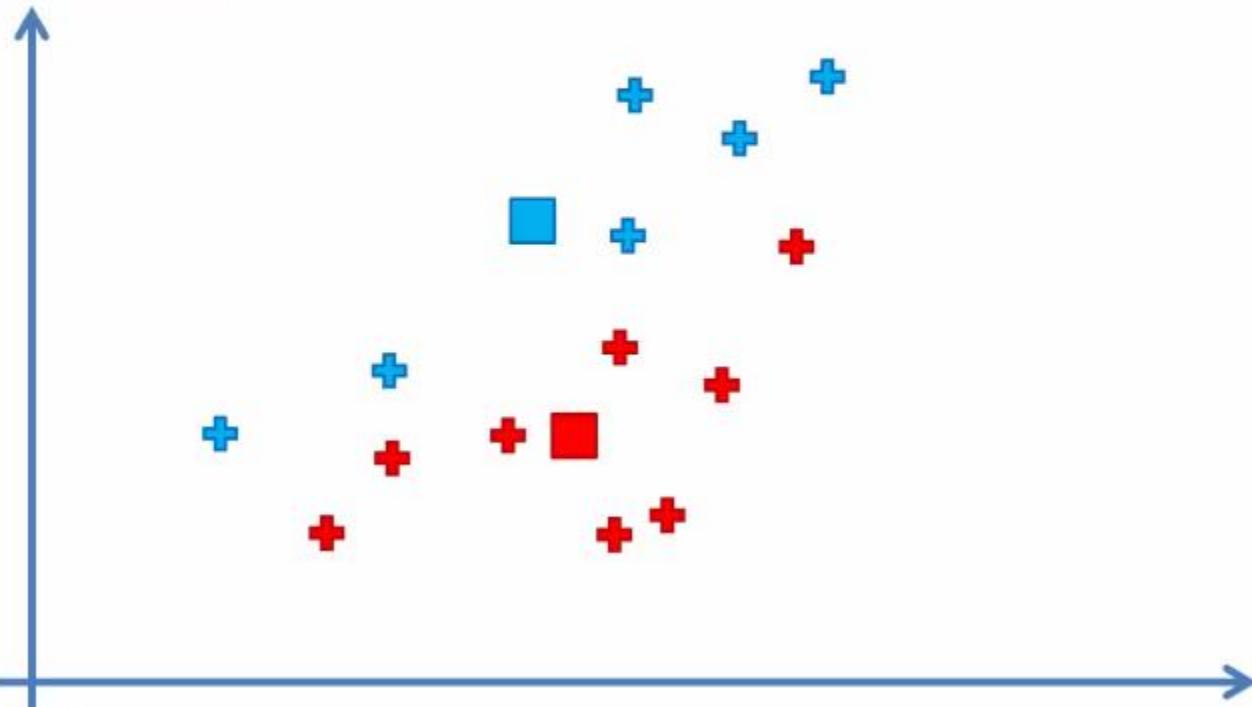
STEP 4: Compute and place the new centroid of each cluster



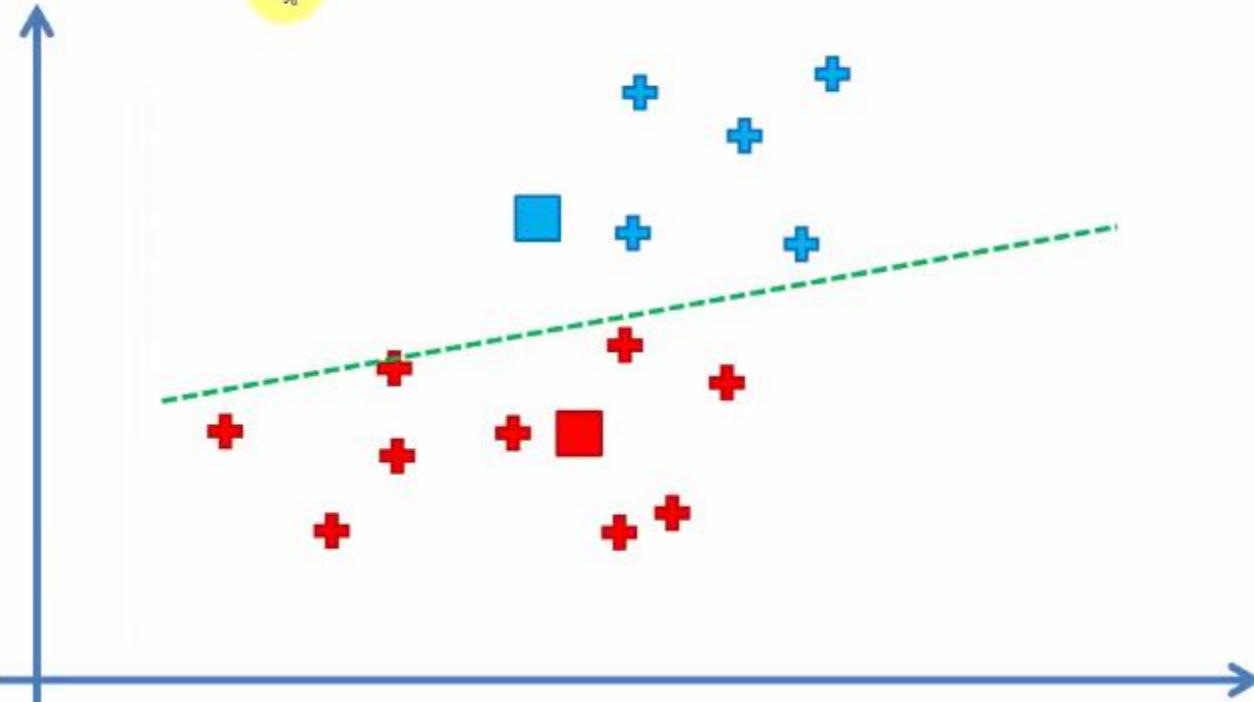
STEP 4: Compute and place the new centroid of each cluster



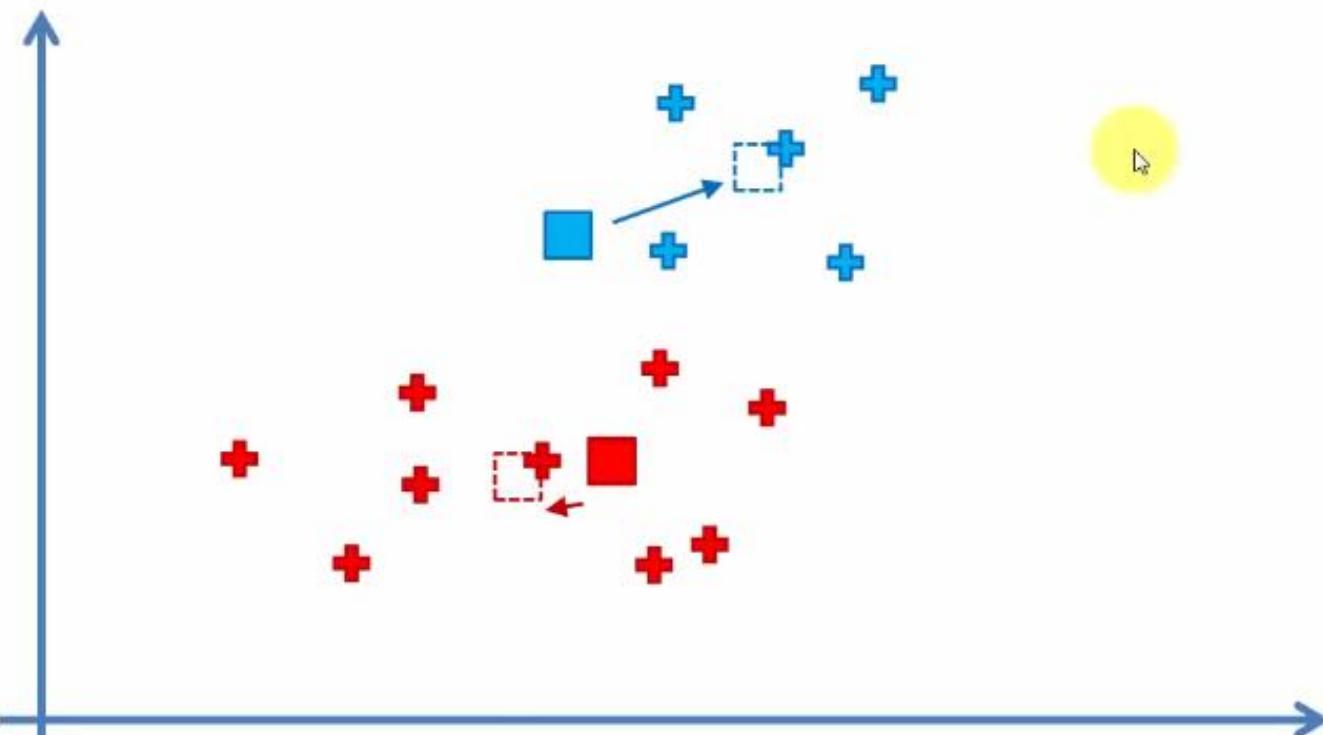
STEP 5: Reassign each data point to the new closest centroid.
If any reassignment took place, go to STEP 4, otherwise go to FIN.



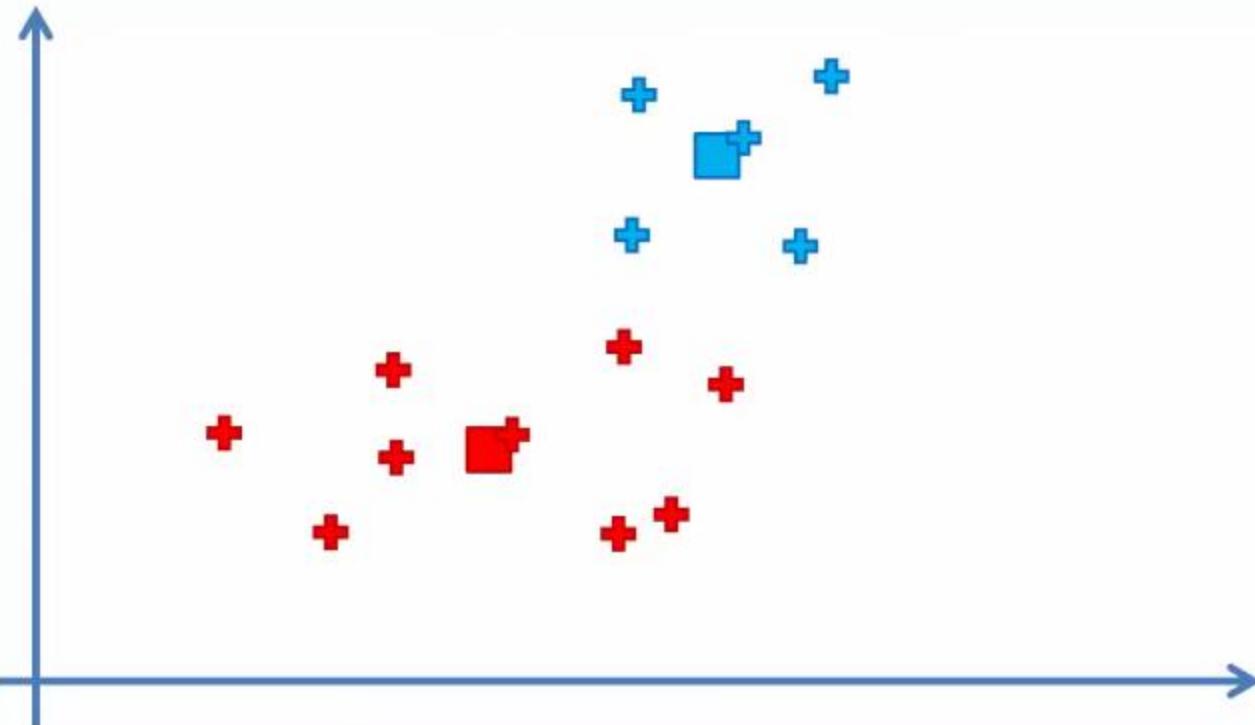
STEP 5: Reassign each data point to the new closest centroid.
If any reassignment took place, go to STEP 4, otherwise go to FIN.



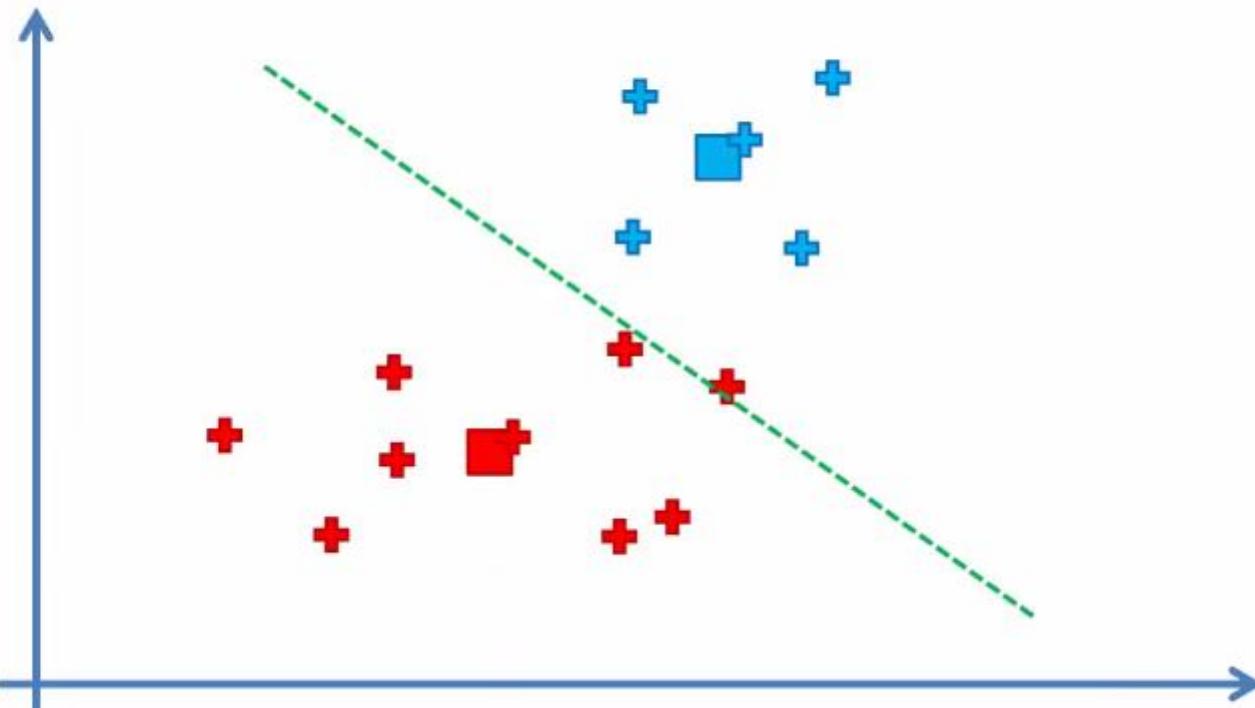
STEP 4: Compute and place the new centroid of each cluster



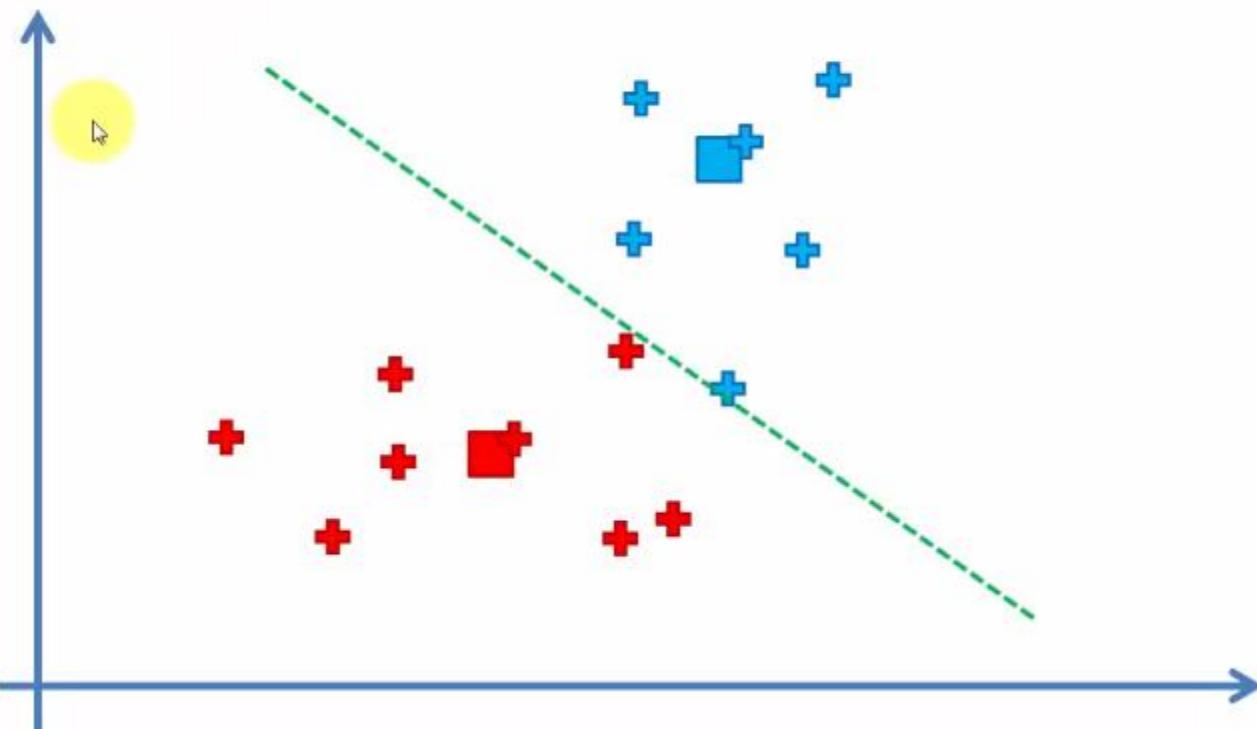
STEP 4: Compute and place the new centroid of each cluster



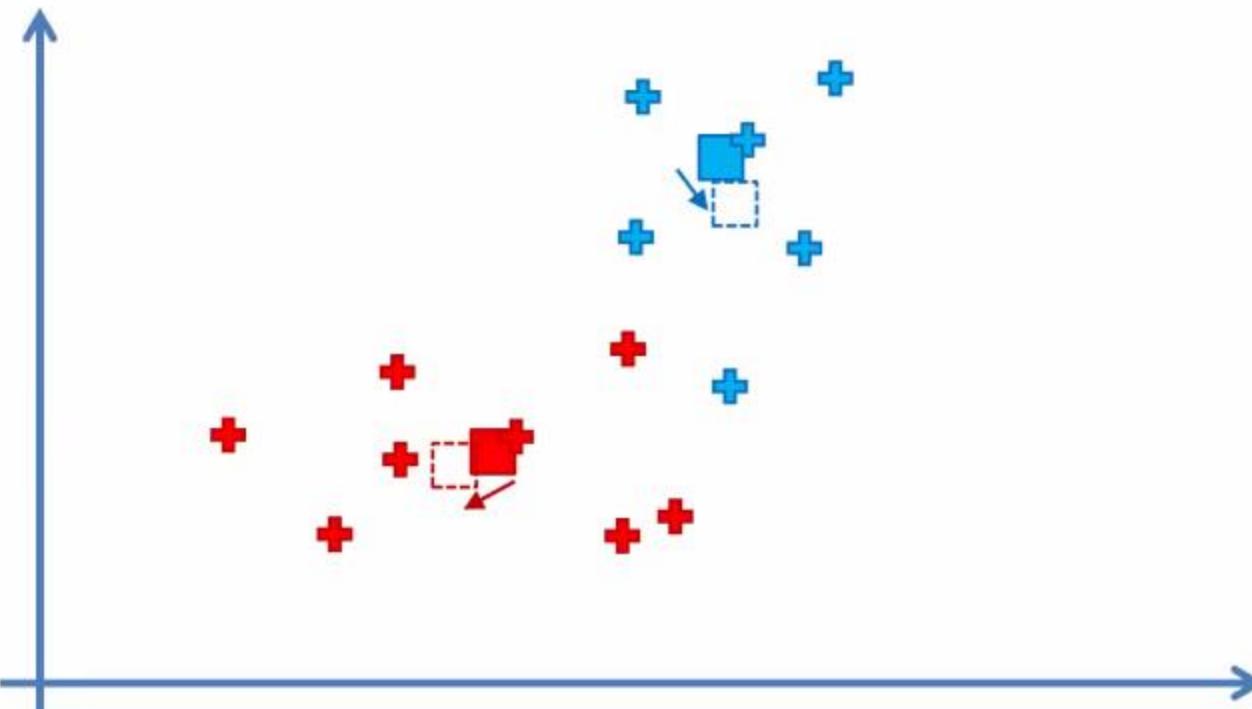
STEP 5: Reassign each data point to the new closest centroid.
If any reassignment took place, go to STEP 4, otherwise go to FIN.



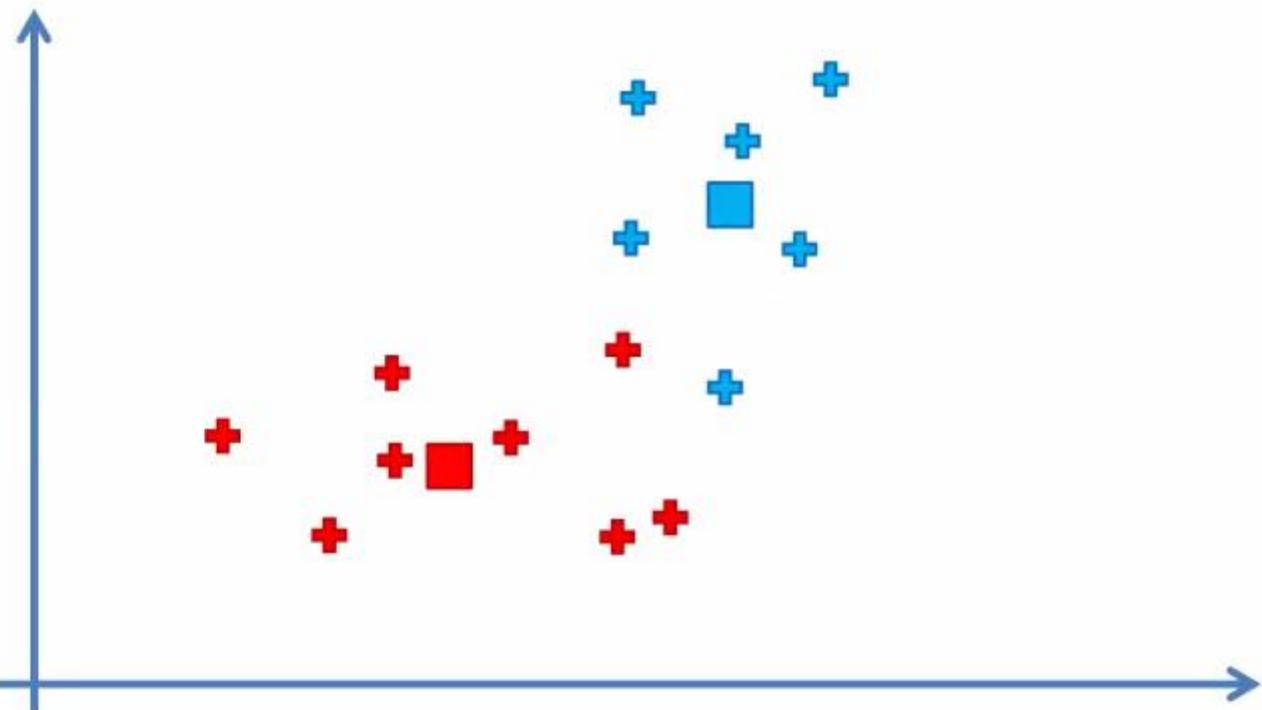
STEP 5: Reassign each data point to the new closest centroid.
If any reassignment took place, go to STEP 4, otherwise go to FIN.



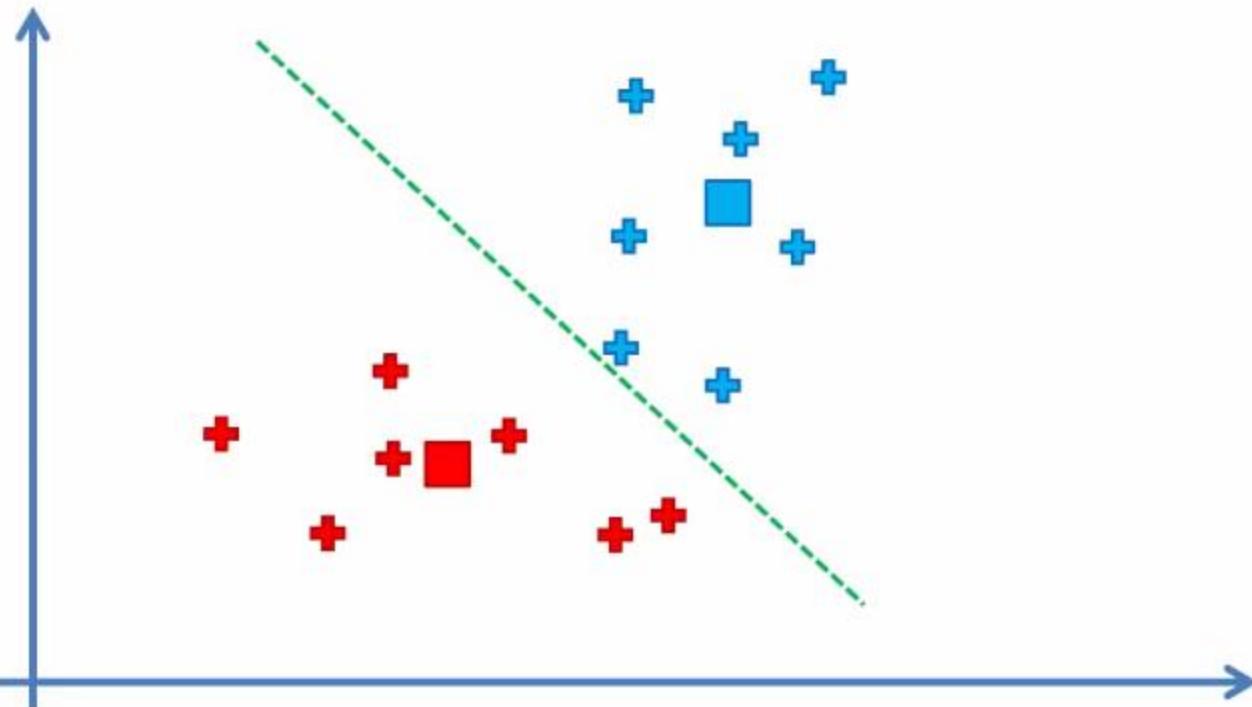
STEP 4: Compute and place the new centroid of each cluster



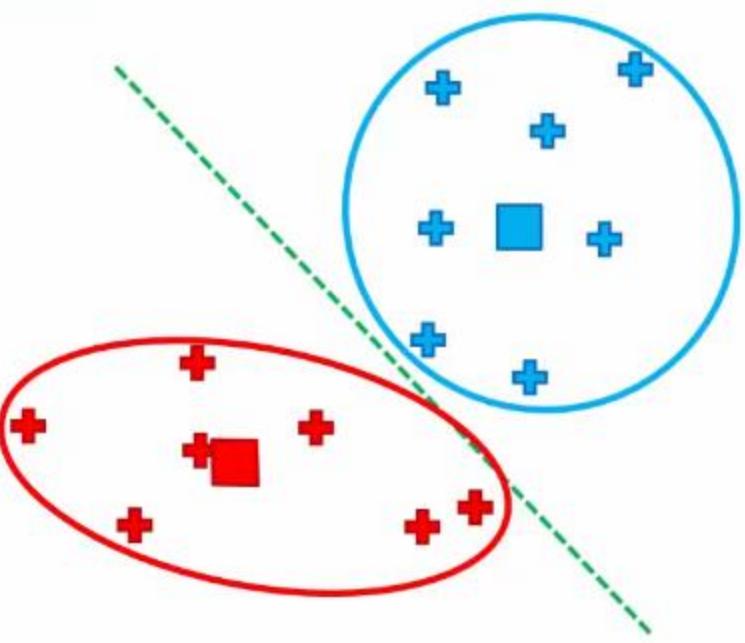
STEP 4: Compute and place the new centroid of each cluster



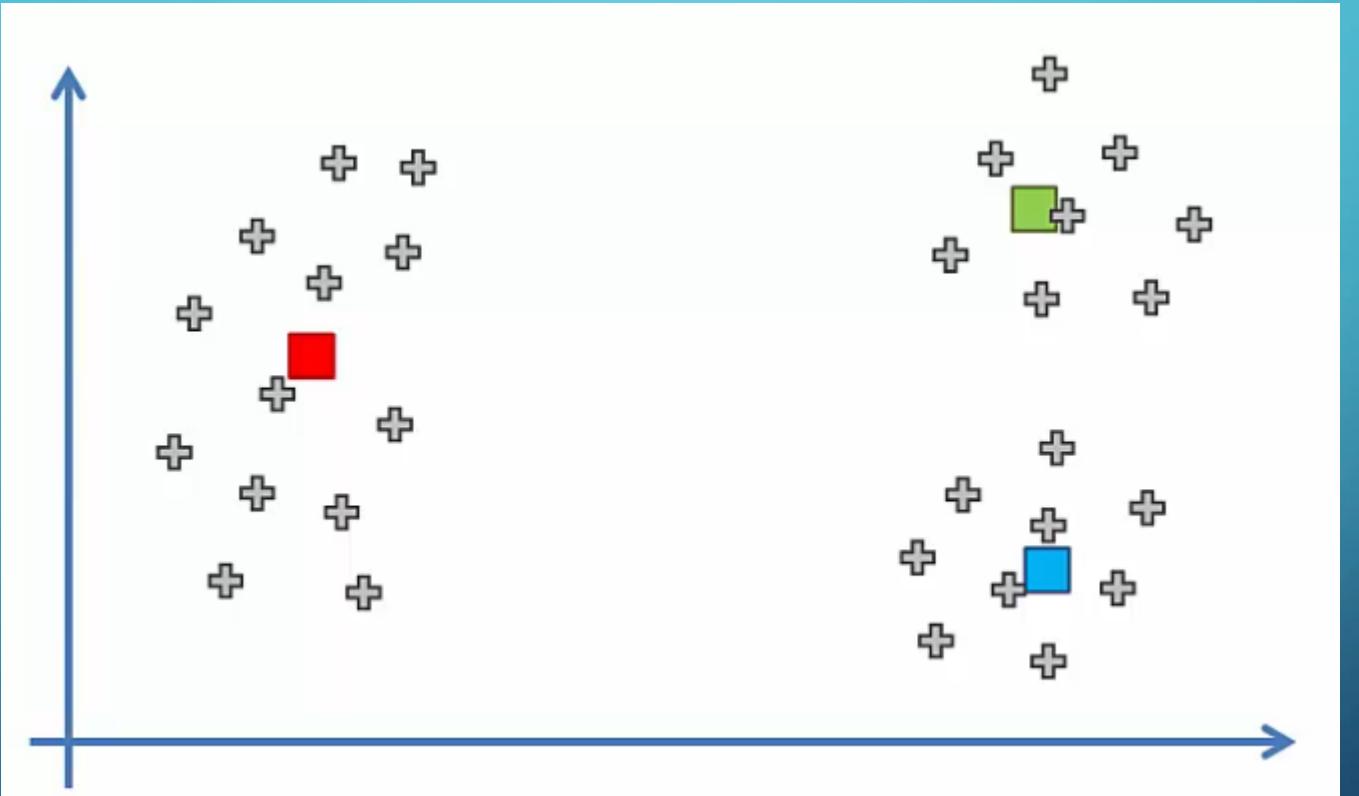
STEP 5: Reassign each data point to the new closest centroid.
If any reassignment took place, go to STEP 4, otherwise go to FIN.

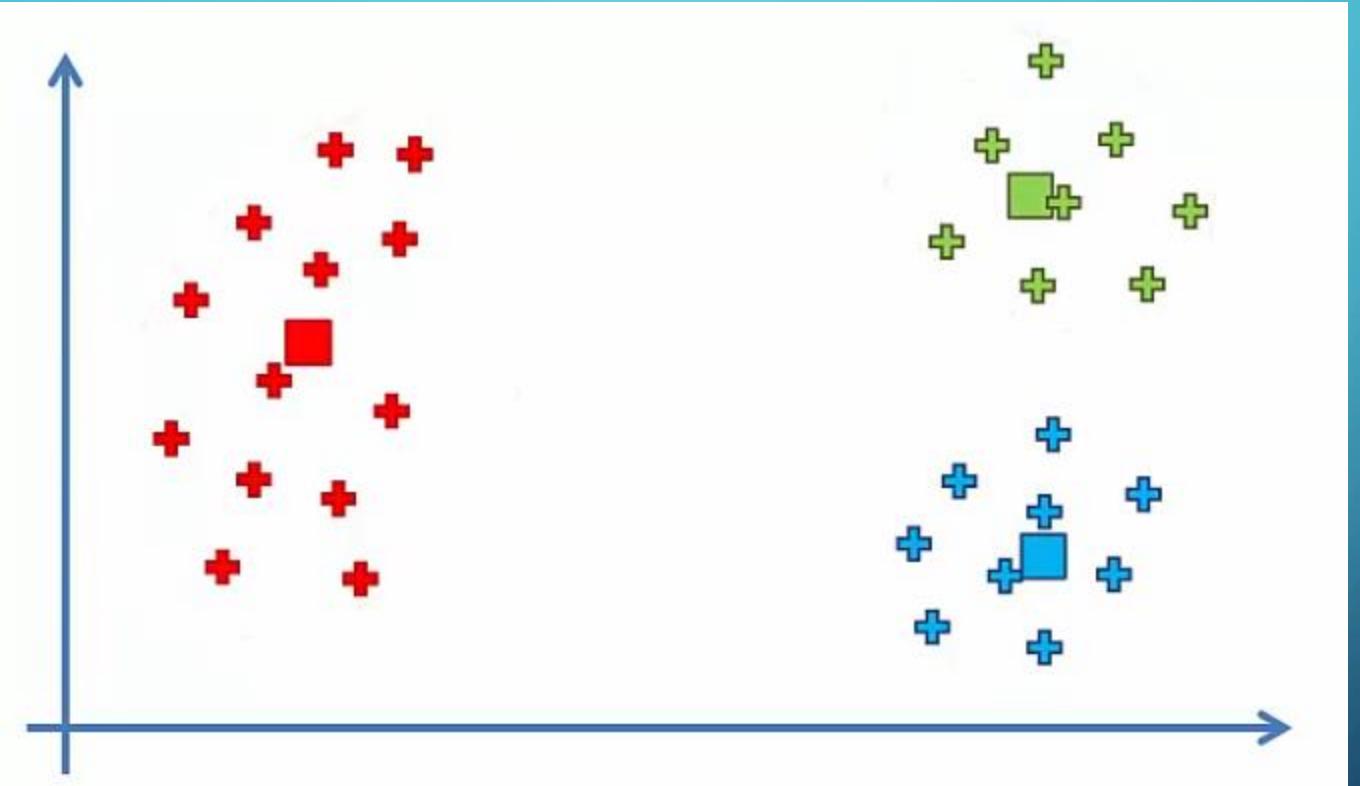


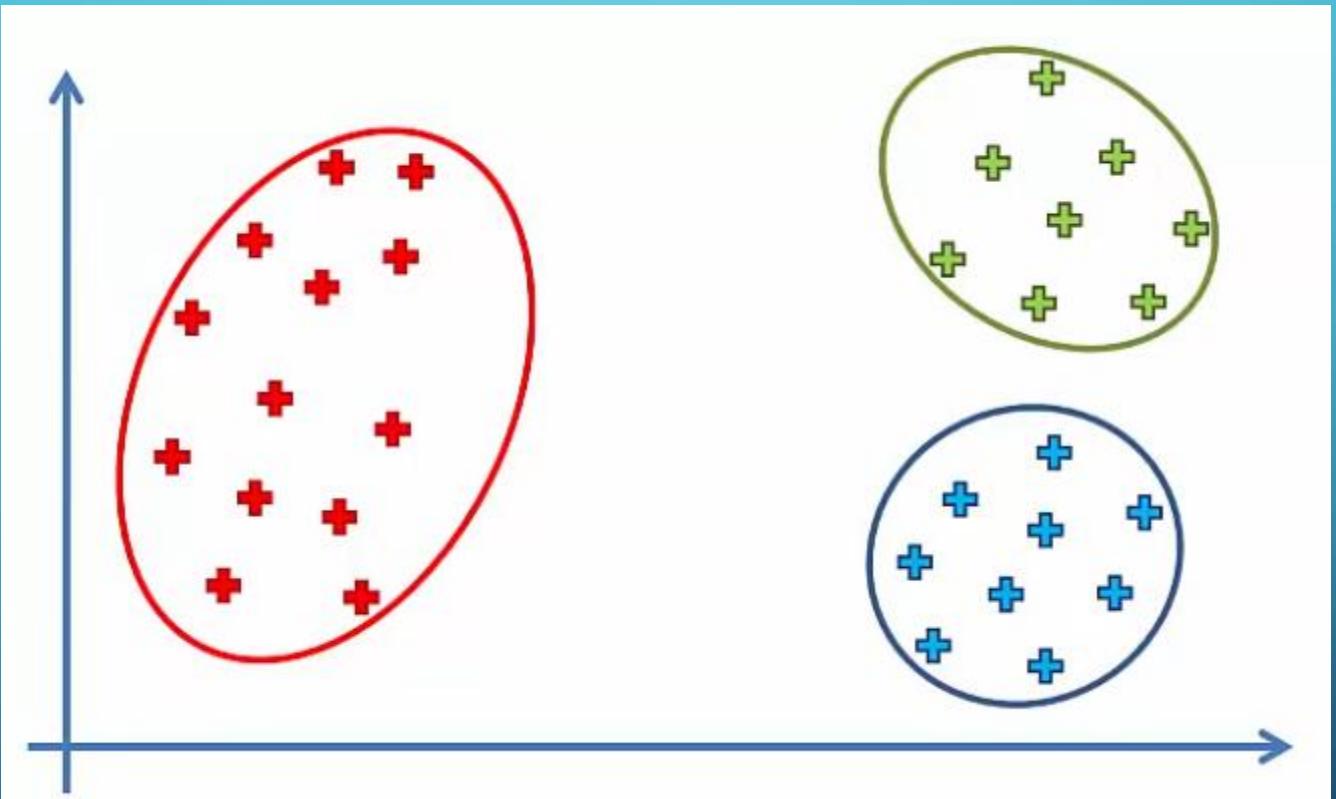
FIN: Your Model Is Ready



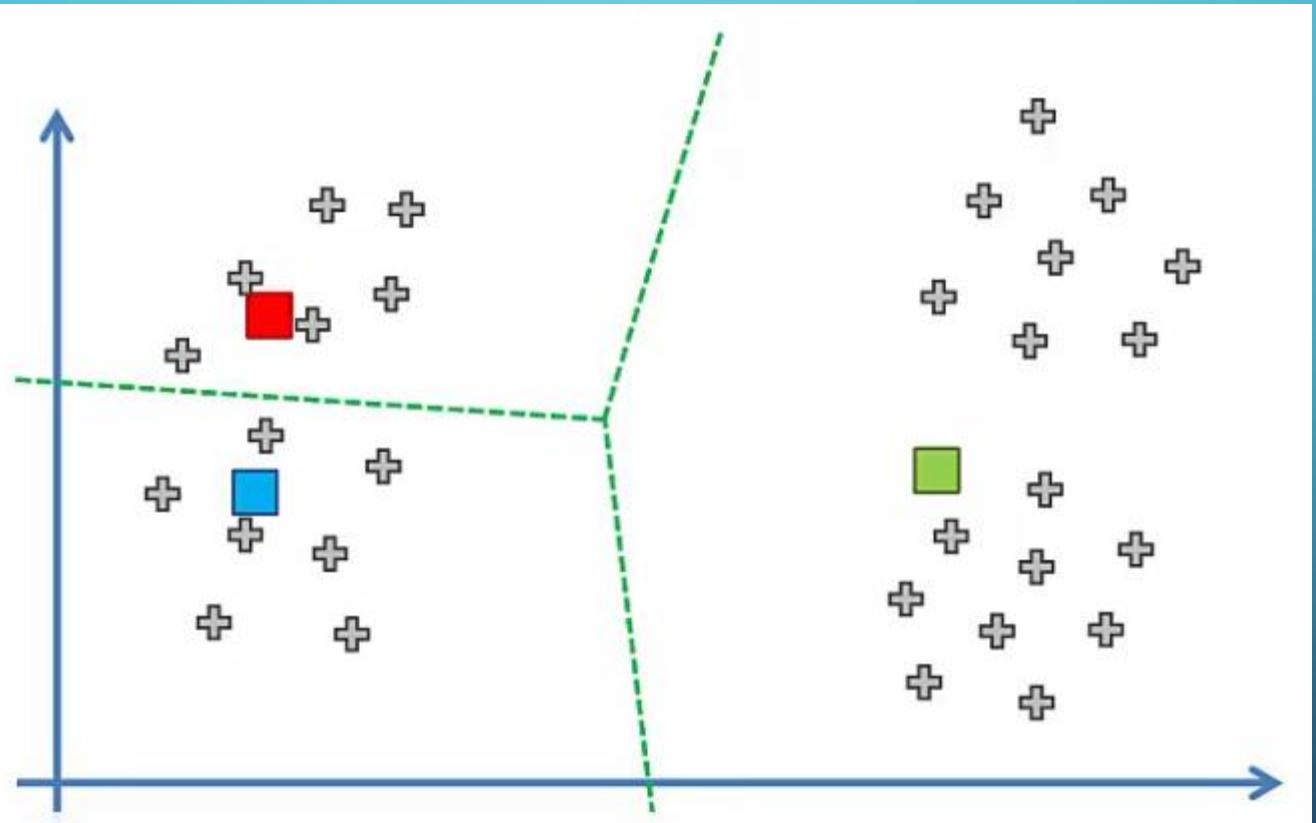
RANDOM INITIALIZATION TRAP

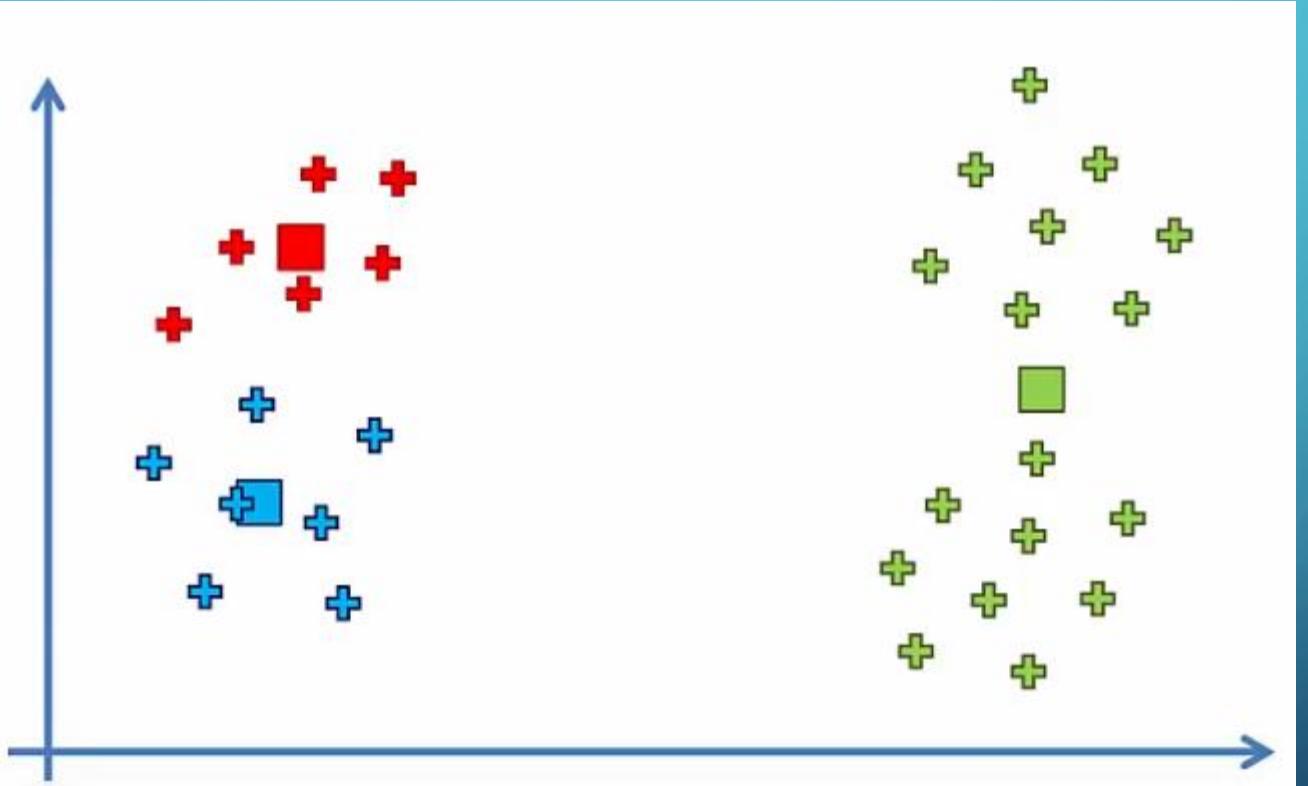


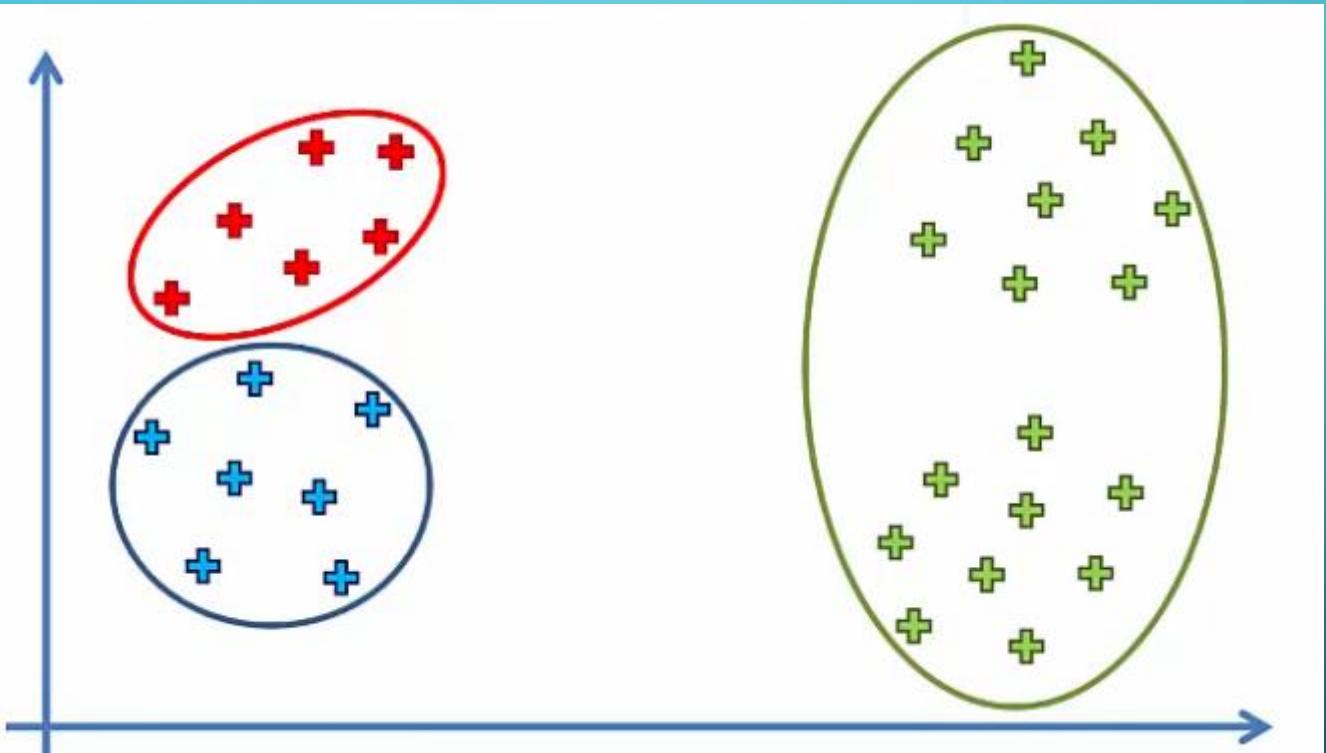


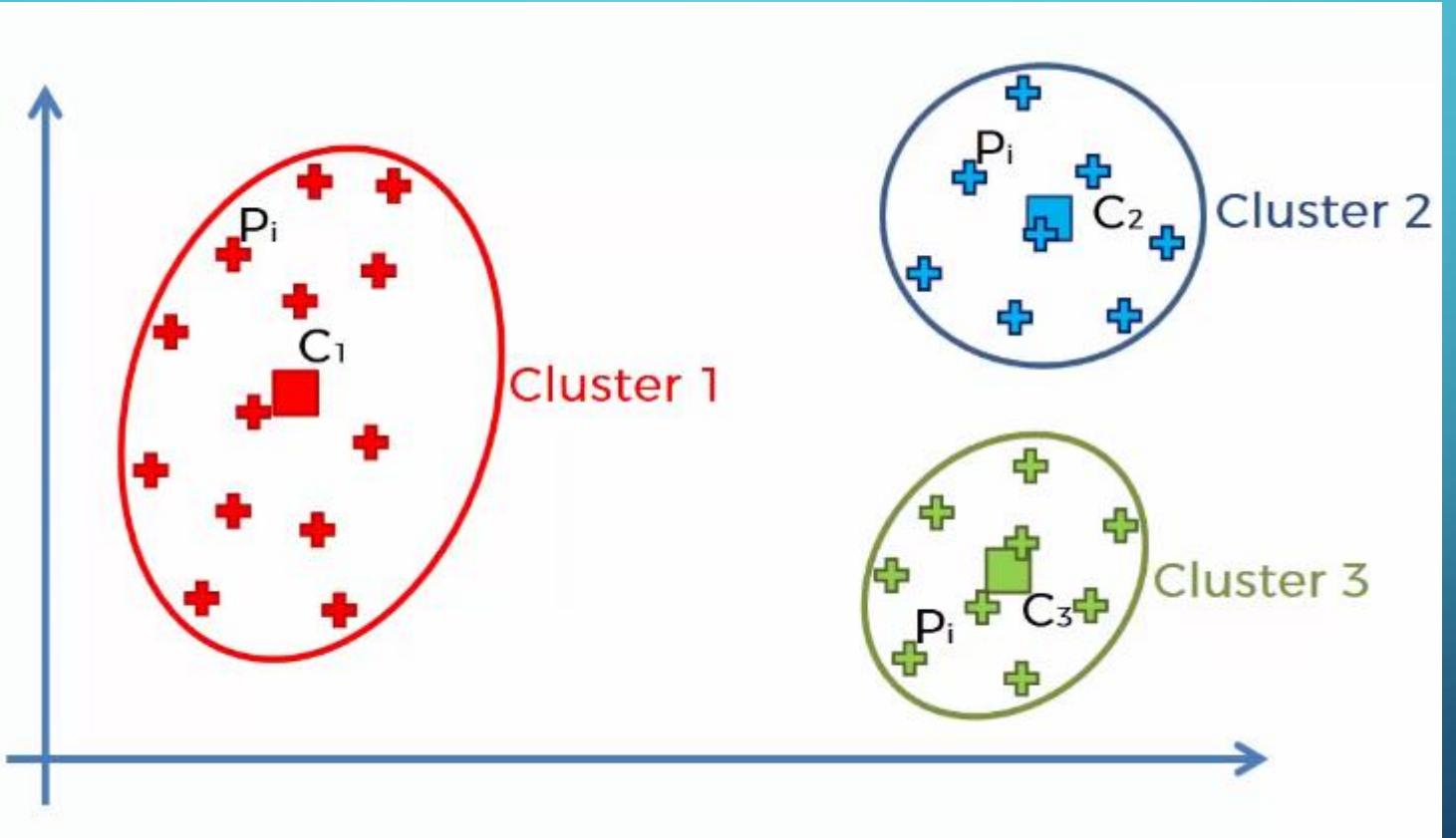






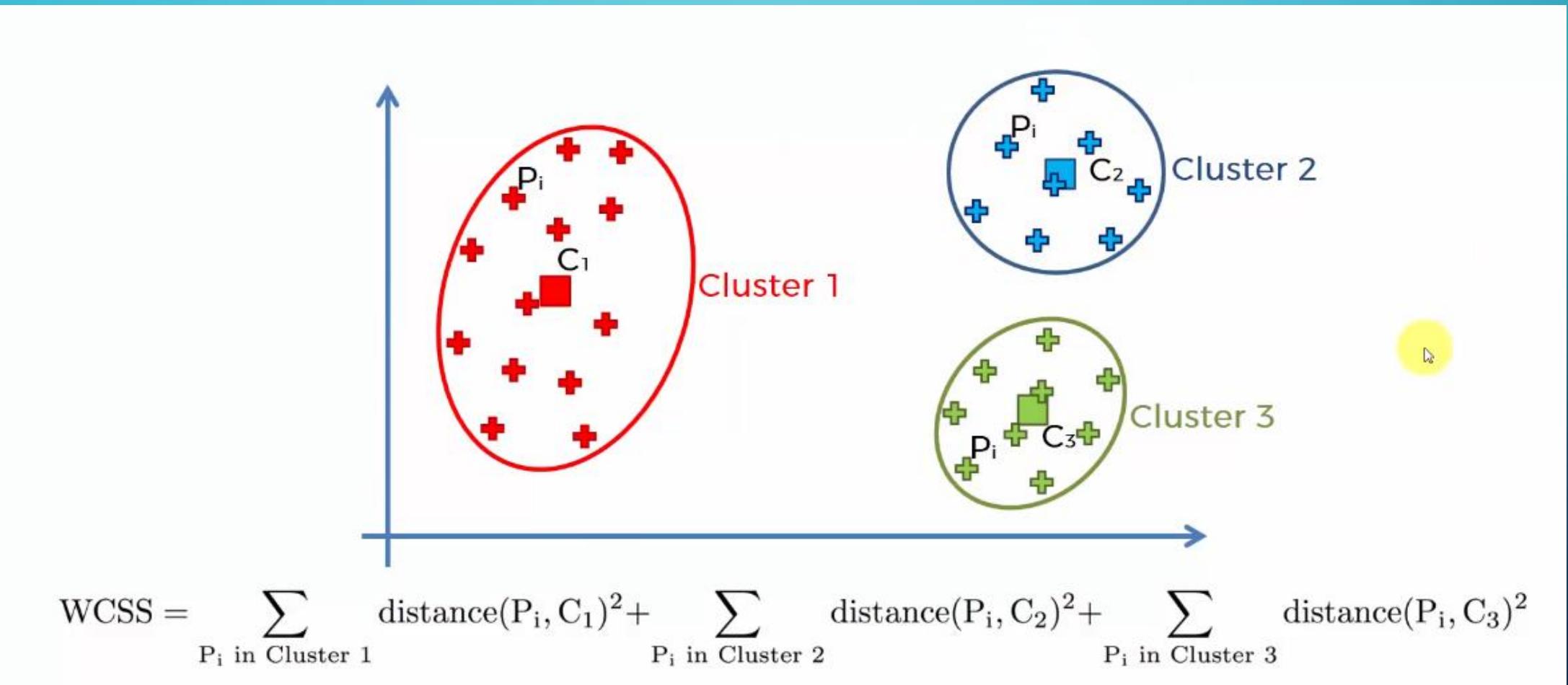




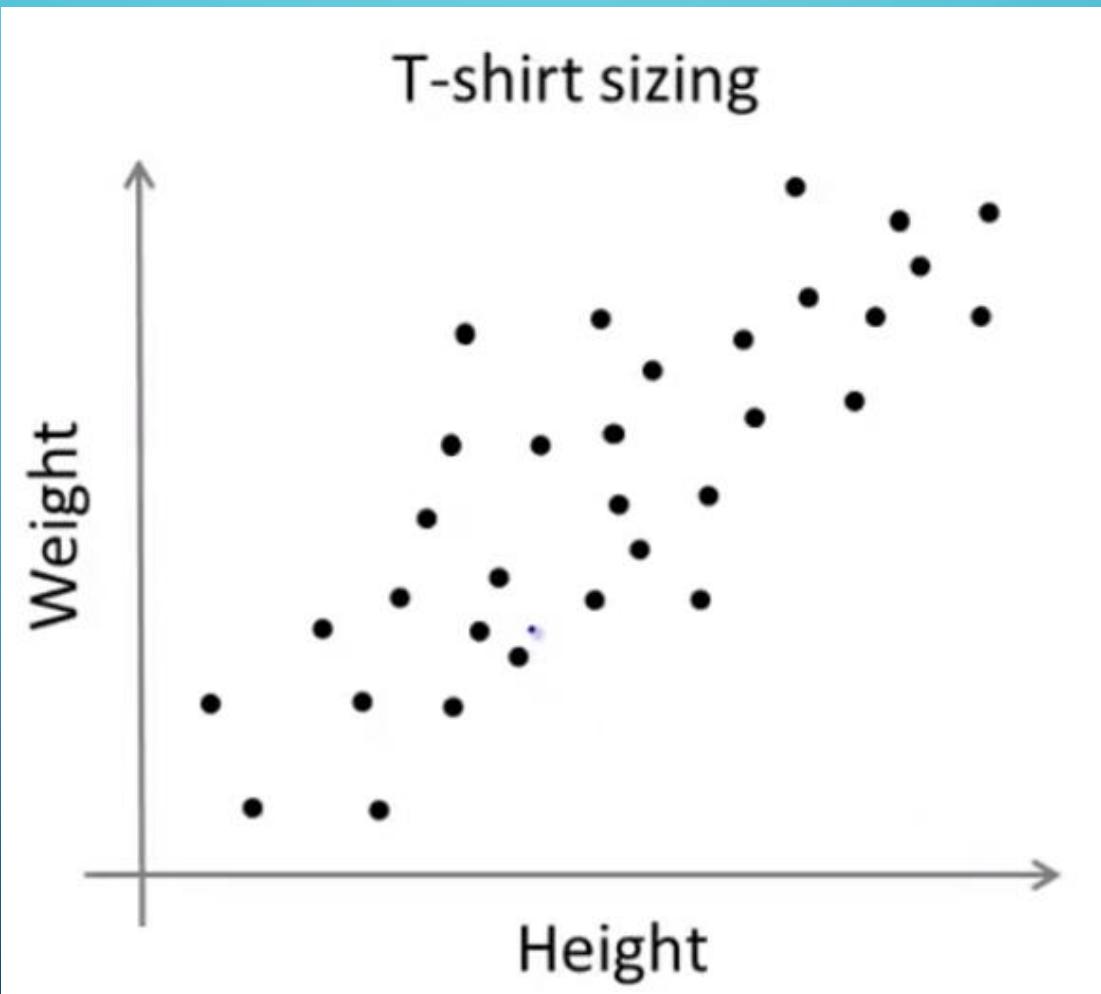


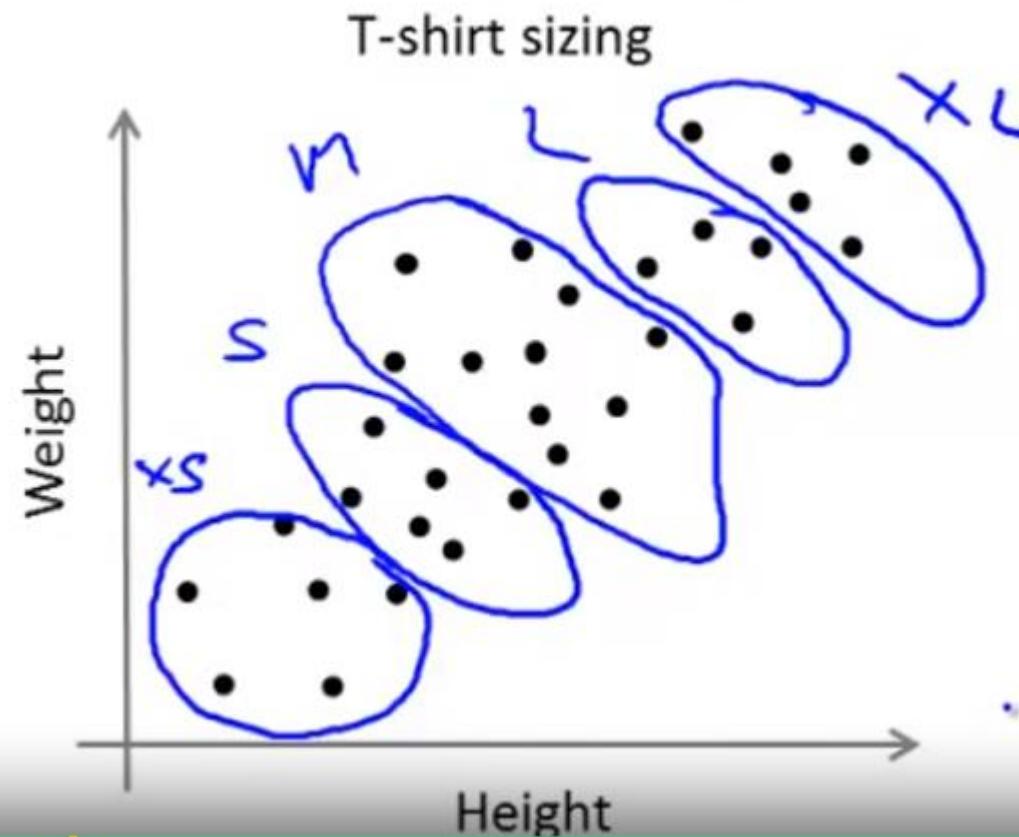
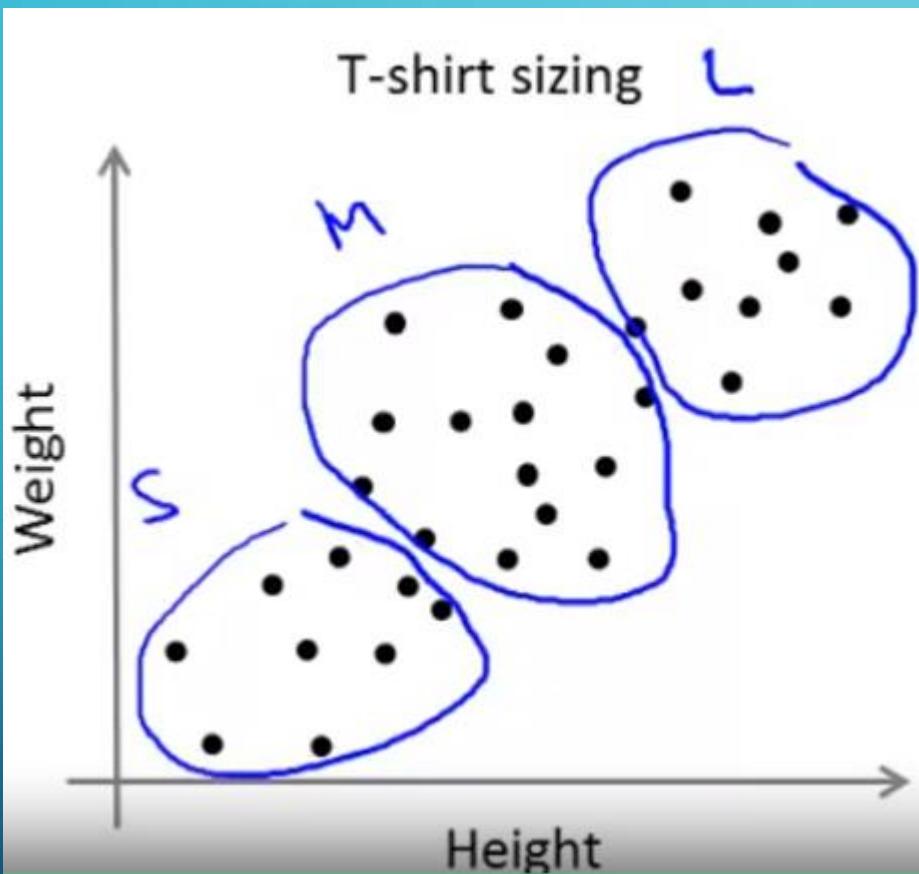
CHOOSING THE RIGHT CLUSTER

$$\text{WCSS} = \sum_{P_i \text{ in Cluster 1}} \text{distance}(P_i, C_1)^2 + \sum_{P_i \text{ in Cluster 2}} \text{distance}(P_i, C_2)^2 + \sum_{P_i \text{ in Cluster 3}} \text{distance}(P_i, C_3)^2$$



HOW TO CHOOSE NUMBER OF CLUSTERS







THANK YOU

AGENDA FOR TOMORROW:

- Doubt session
- Practical sessions
- Using Github as repository