



# AI 330: Machine Learning

## Fall 2023

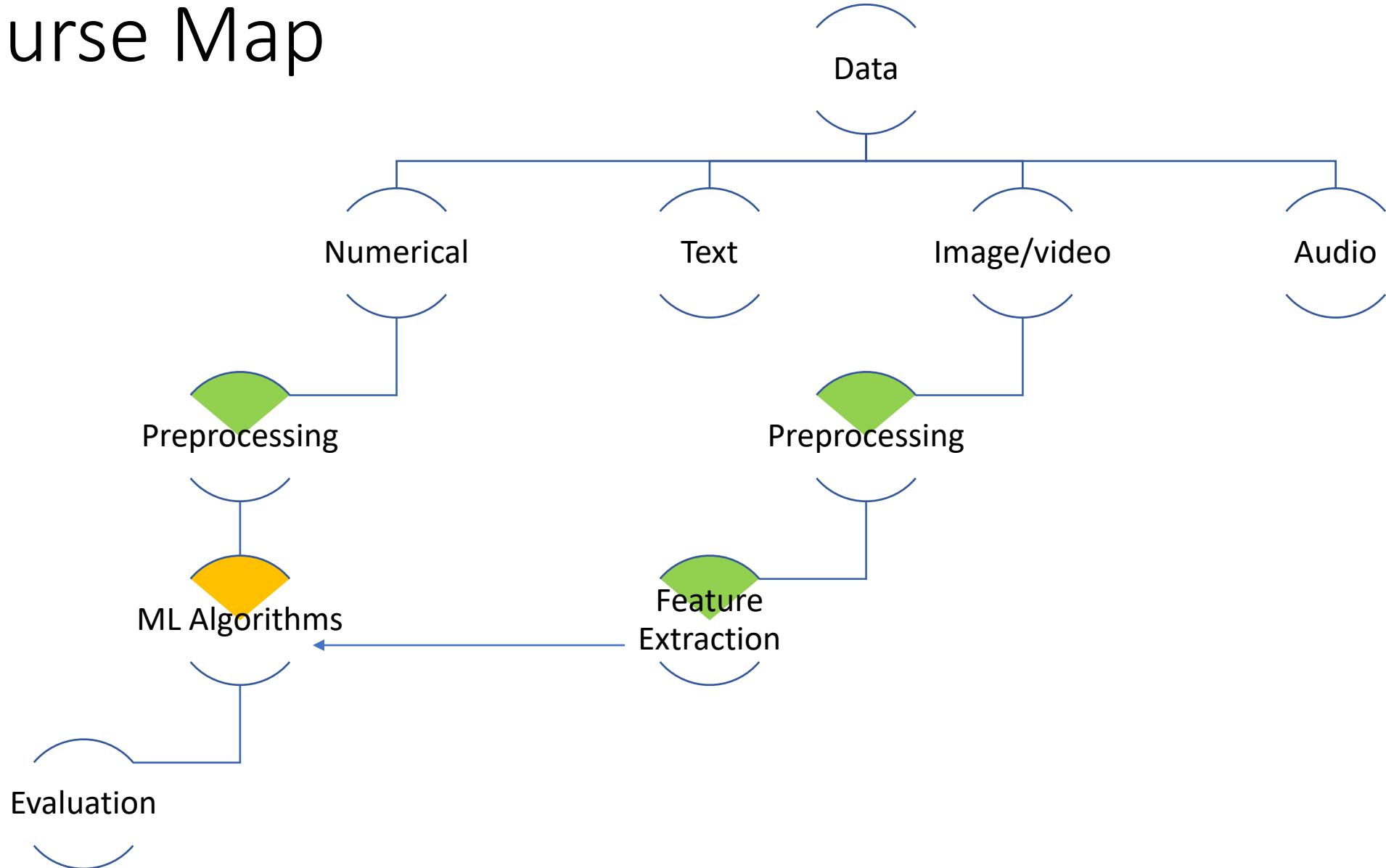
**Dr. Wessam EL-Behaidy**

Associate Professor, Computer Science Department,  
Faculty of Computers and Artificial Intelligence,  
Helwan University.

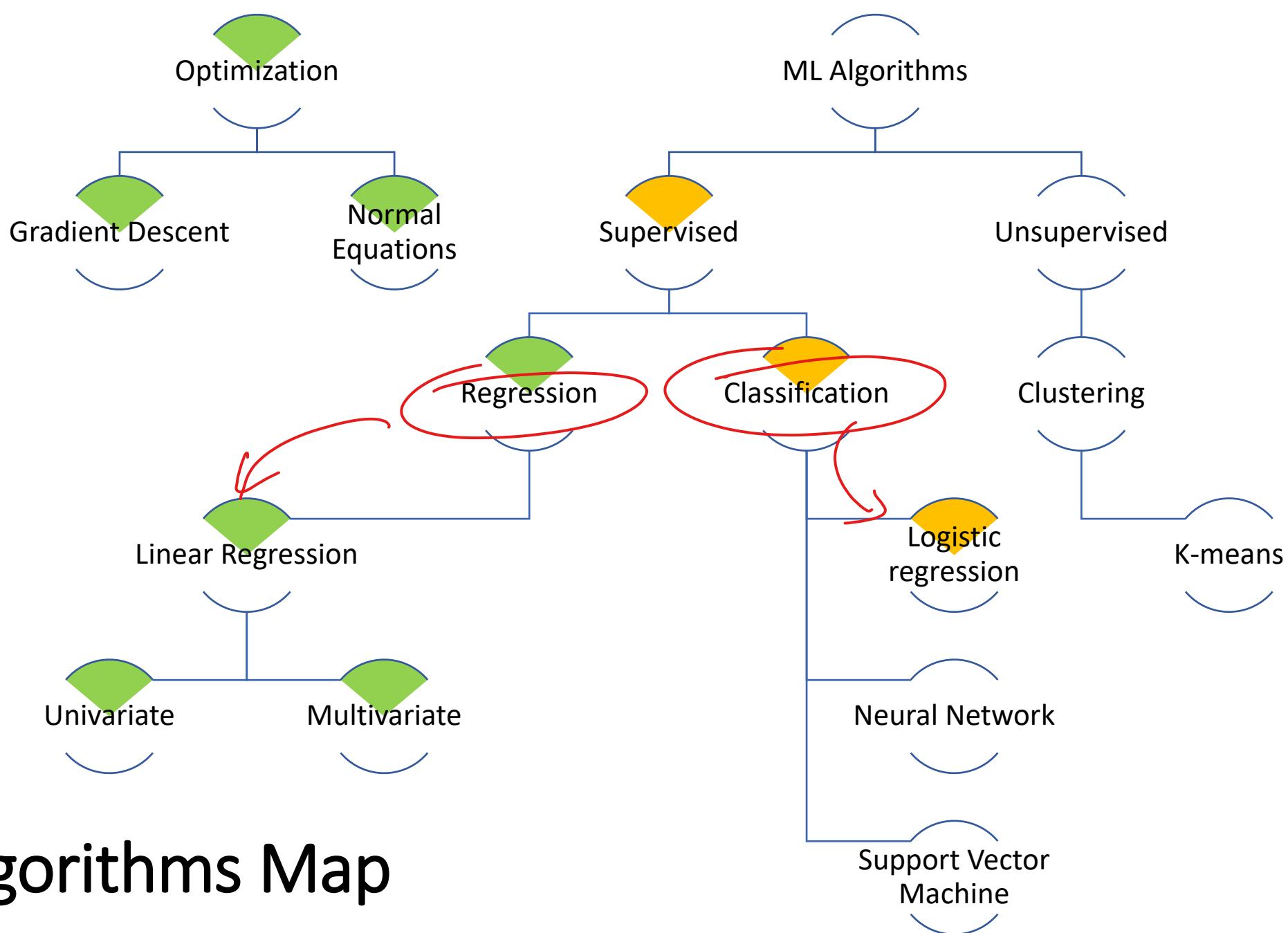
**Dr. Ensaif Hussein**

Associate Professor, Computer Science Department,  
Faculty of Computers and Artificial Intelligence,  
Helwan University.

# Course Map



# ML Algorithms Map



# Lecture 6

---

# Logistic Regression

&

# Regularization

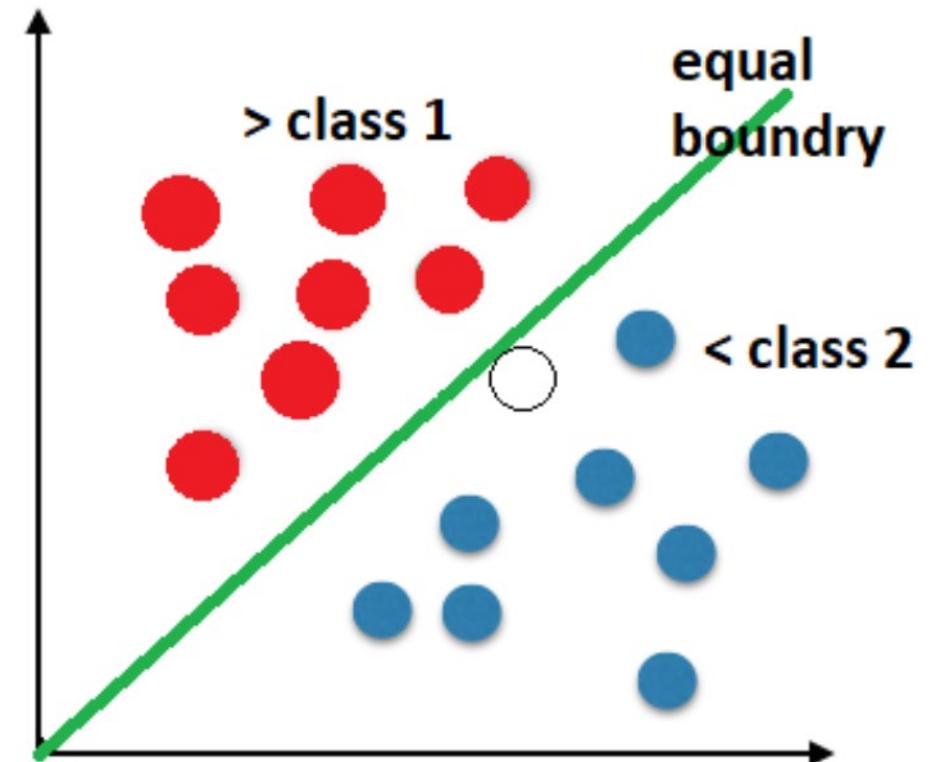
Slides of:

<https://www.coursera.org/learn/machine-learning> at Stanford University (Prof. Andrew Ng)

# Classification Problem: Logistic Regression

# Classification

- Discrete outcomes.
- Binary  $y \in \{0,1\}$ , 0 negative , 1 positive (normal / abnormal)
- Multi-class: a telescope that identifies whether an object in the night sky is a galaxy, star, or planet.  
 $y \in \{0,1,2,3\}$



- **Logistic regression** is a statistical model used for **binary classification** tasks, where the goal is to predict whether an outcome belongs to one of two possible classes (e.g., yes/no, 1/0). It's called "logistic" because it utilizes the logistic function to model the probability of an event occurring.
- In **logistic regression**, the input features are linearly combined with weights, and the result is transformed using the **logistic (sigmoid)** function, which maps the linear combination to a value between 0 and 1. This value represents the probability of the event being in the positive class.

- Mathematically, the logistic regression model can be expressed as:

$$\bullet P(Y=1|X) = 1 / (1 + e^{-(b_0 + b_1X_1 + b_2X_2 + \dots + b_nX_n)})$$

- Here:

$\downarrow$  X feature  $\rightarrow$  class(1)  $\rightarrow$  Y  $\uparrow$

- $P(Y=1|X)$  is the probability of the event Y being in class 1 given the input features X.
- $b_0, b_1, b_2, \dots, b_n$  are the model parameters (coefficients) to be learned from the training data.
- $X_1, X_2, \dots, X_n$  are the input features.

linear if  $b \neq 0$  is  
regression

the threshold is a value that is used to classify the predicted probabilities into two classes. Logistic regression predicts the probability that an instance belongs to a particular class, and the threshold is the point at which you decide which class the instance should be assigned to based on its predicted probability

- The model is trained by optimizing these parameters using techniques like maximum likelihood estimation. Once trained, it can be used to make predictions by classifying instances based on the calculated probabilities. Typically, a threshold (e.g., 0.5) is chosen, and if the predicted probability is above the threshold, the instance is classified as the positive class; otherwise, it's classified as the negative class.

above the threshold

under the threshold

$x=0$

$\text{hypothesis}$   
 $0.5 > \text{prob}$

- Logistic regression is widely used in various fields, such as healthcare, finance, and marketing, for tasks like spam detection, customer churn prediction, and medical diagnosis.

# Hypothesis Representation

- Logistic regression model

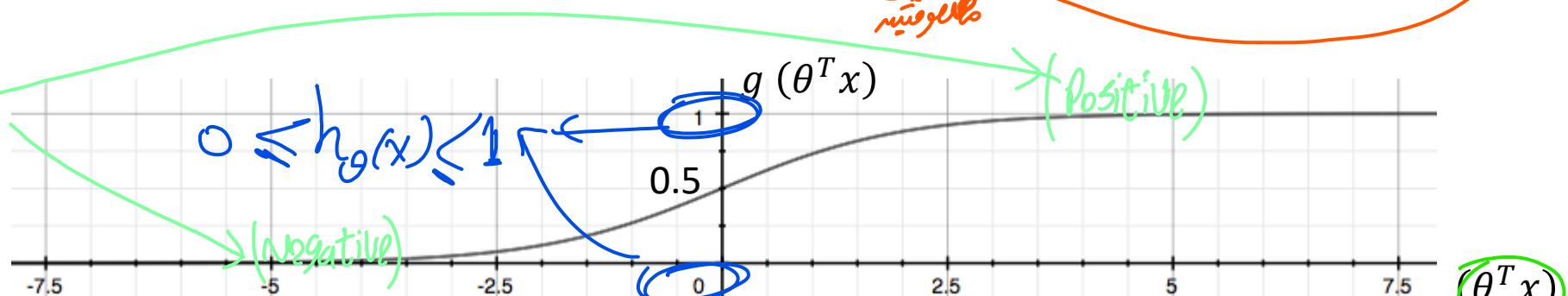
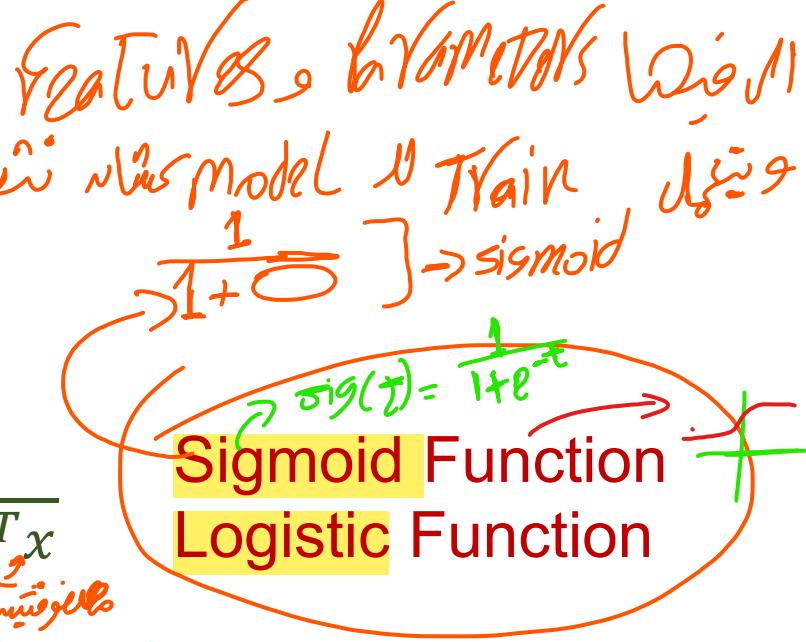
هذا المفهوم ينبع من  
الـ 1 والـ 0 والـ 0.5

وذلك لأنها ملحوظة  
وهي كلاسيكية

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

$$0 \leq h_{\theta}(x) \leq 1$$

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$



- We want our classifier to output values between 0 and 1

- When using linear regression we did  $h_{\theta}(x) = (\theta^T x)$

- For classification hypothesis representation we do  $h_{\theta}(x) = g((\theta^T x))$

# Interpretation of hypothesis Output

- $h_{\theta}(x)$  will give us the **probability** that our output is 1.

**For example,**

(1)  $h_{\theta}(x) = 0.7$  is 70%  $\rightarrow$  1

$h_{\theta}(x) = 0.7$  gives us a **probability** of 70% that our output is 1.

- Our probability that our prediction is 0 is just the complement of our probability that it is 1

**For example,**

if probability that it is 1 is 70%, then

the probability that it is 0 is 30%

0  $\rightarrow$  output 1 is 30%

# Interpretation of hypothesis Output

- $h_\theta(x)$  will give us the **probability** that our output is 1.

For example,

*features*  $\hookrightarrow \mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumourSize} \end{bmatrix}$ ,  $h_\theta(\mathbf{x})=0.7 \dots y = 1$

*about 70%  $x_0=1$*

→ 70% chance of a tumor being malignant.

- $h_\theta(x) = P(y = 1|x; \theta) = 1 - P(y = 0|x; \theta)$

Probability that  $y = 1$ , given  $x$ , parameterized by  $\theta$

# Decision Boundary

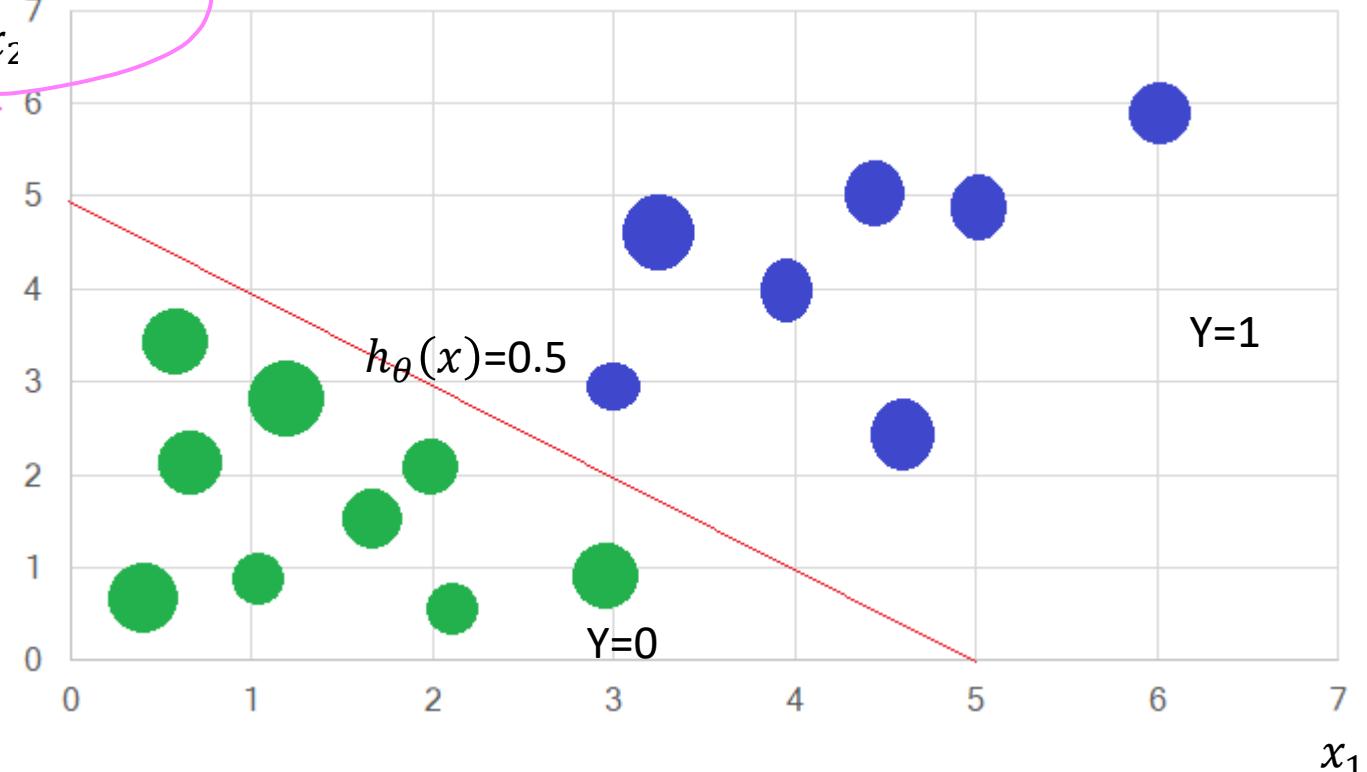
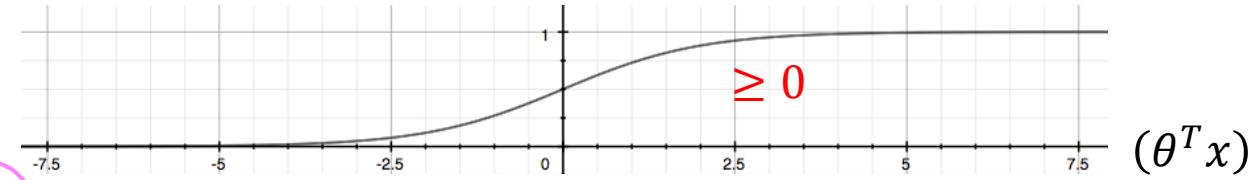
$$\begin{aligned}
 h_{\theta}(x) &= g(\theta^T x) \\
 &= g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) \\
 &= g(-5 + 1x_1 + 1x_2)
 \end{aligned}$$

$$\theta = \begin{bmatrix} -5 \\ 1 \\ 1 \end{bmatrix}$$

$\begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$

Predict  $y = 1$  if  $-5 + x_1 + x_2 \geq 0$

$x_1 + x_2 \geq 5$



# Non-linear decision boundaries

مکار (مکار) مکاری  
مکار (مکار) مکاری

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

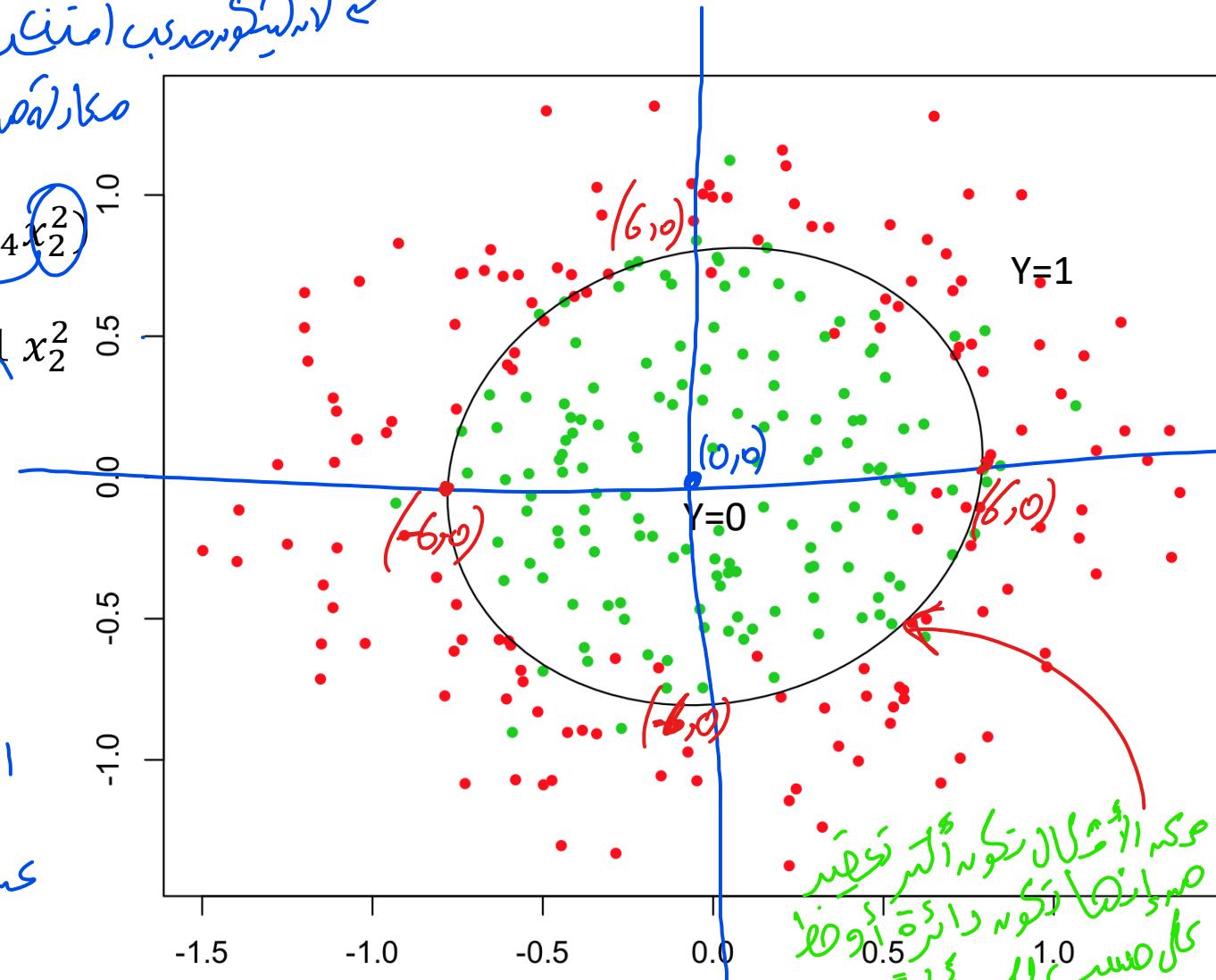
$\theta_0 = -0.6$  (bias)  
 $\theta_1, \theta_2 = 0$   
 $\theta_3, \theta_4 = 1$

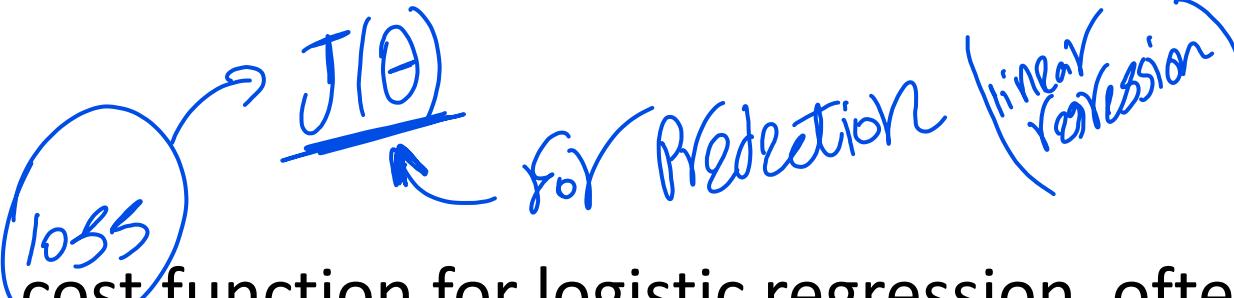
$$= -0.6 + 0 x_1 + 0 x_2 + 1 x_1^2 + 1 x_2^2$$

*the same 2 features*

$$y = 1, if -0.6 + x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 0.6$$





Not good for classification

- The cost function for logistic regression, often referred to as the "log loss" or "cross-entropy loss," measures how well a logistic regression model is performing in terms of its ability to predict the probability of an event belonging to a particular class. The goal is to minimize this cost function during the training of the model.

$y=0 \text{ or } y=1$  : not good for classification

- he cost function is defined as follows:

- $\text{Cost}(b_0, b_1, b_2, \dots, b_n) = -\frac{1}{m} \sum [y * \log(P) + (1 - y) * \log(1 - P)]$

- Here:

such as "Cost" is the cost function.

$\theta_0, \theta_1, \theta_2$  "b<sub>0</sub>, b<sub>1</sub>, b<sub>2</sub>, ..., b<sub>n</sub>" are the model parameters (coefficients).

- "m" is the number of training examples.

- " $\Sigma$ " represents the summation over all training examples.

correct class "y" is the actual class label (0 or 1) for a given example.

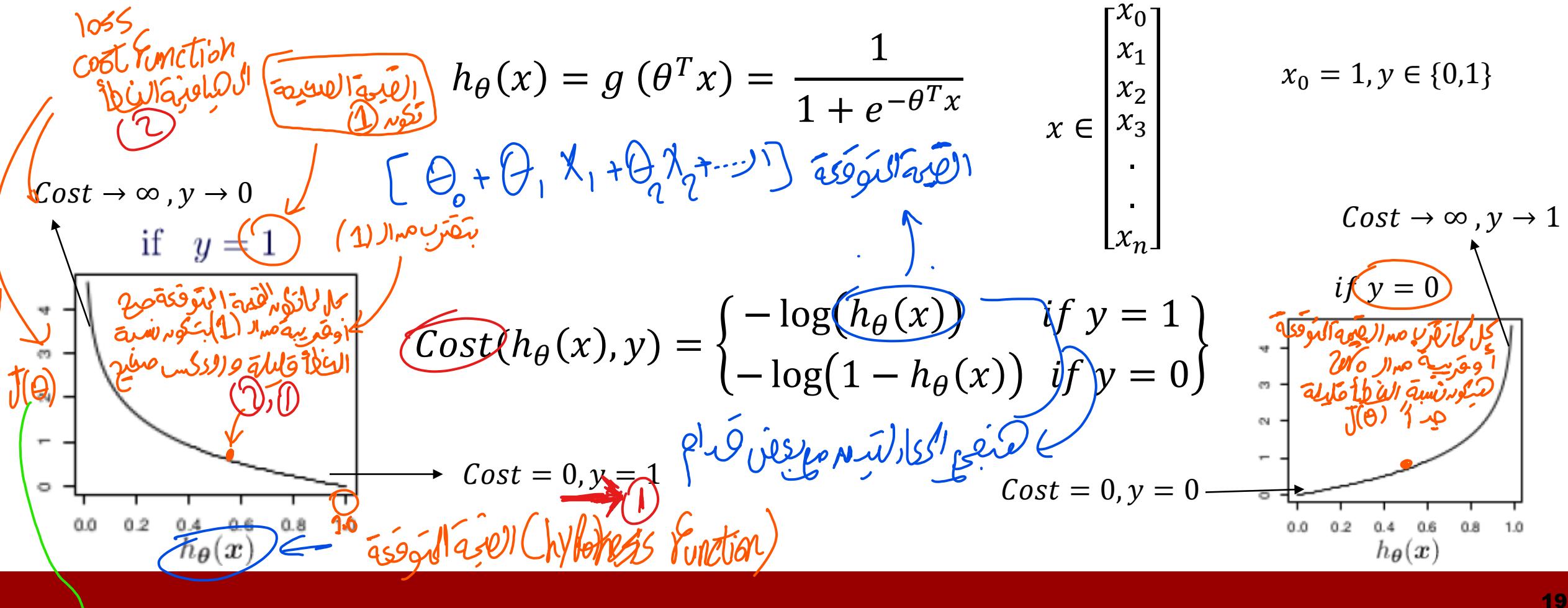
- "P" is the predicted probability that the example belongs to class 1, as calculated by the logistic regression model.

logistic

- The cost function penalizes the model more when its predicted probability is far from the actual class label. If the model's prediction is close to the true class label, the **log loss** is small; if it's far off, the log loss becomes significantly larger.  
*cost function* *مُخاول لـ الـ لـ لـ*
  - During the training process, the model's parameters ( $b_0, b_1, b_2, \dots, b_n$ ) are adjusted to minimize this cost function. Techniques like **gradient descent** or other optimization methods are used to find the values of the parameters that lead to the lowest cost.
  - In logistic regression, the ultimate goal is to find the parameter values that result in the **lowest cost**, indicating that the model's predicted probabilities are as close as possible to the true class labels for the training data. This ensures that the model makes accurate predictions on new, unseen data.

# How to choose parameter $\theta$ - Cost Function

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



$$J(\theta) = y - h_{\theta}(x)$$

$$y = 0.8 - 0.2 = 0.6$$

(0.8, 0.2) are the coordinates of the point

## Simplified Cost Function

$$Cost(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases} \equiv Cost(h_{\theta}(x), y) = -y \log h_{\theta}(x) - (1 - y) \log(1 - h_{\theta}(x))$$

when  $y=1$

$$J(h_{\theta}(x), y) = -\log h_{\theta}(x) - \text{zero}$$

when  $y$  is equal to 1, then the second term will be zero and will not affect the result.

If  $y$  is equal to 0, then the first term will be zero and will not affect the result.

$$\text{when } y=0 \quad J(h_{\theta}(x), y) = -\log(1 - h_{\theta}(x))$$

- Cost Function of Logistic regression using training set:

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

نقطة على خط

# Gradient Descent for Logistic Regression

loss fun

- Cost Function of Logistic regression using training set:

*loss fun*

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

- Using Gradient Descent (GD) to  $\min_{\theta} J(\theta)$ :

# Repeat

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

gratuito möglicher

# Feature Pipeline

~~Repeat until convergence :~~

$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$ , simultaneously updates  $\theta$  for every  $j = 0, 1, 2, \dots, n$

$$\frac{\partial}{\partial \theta_j} J(\theta) = -y^{(i)} x_j^{(i)}$$

where,  $h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$

# Multiclass Classification

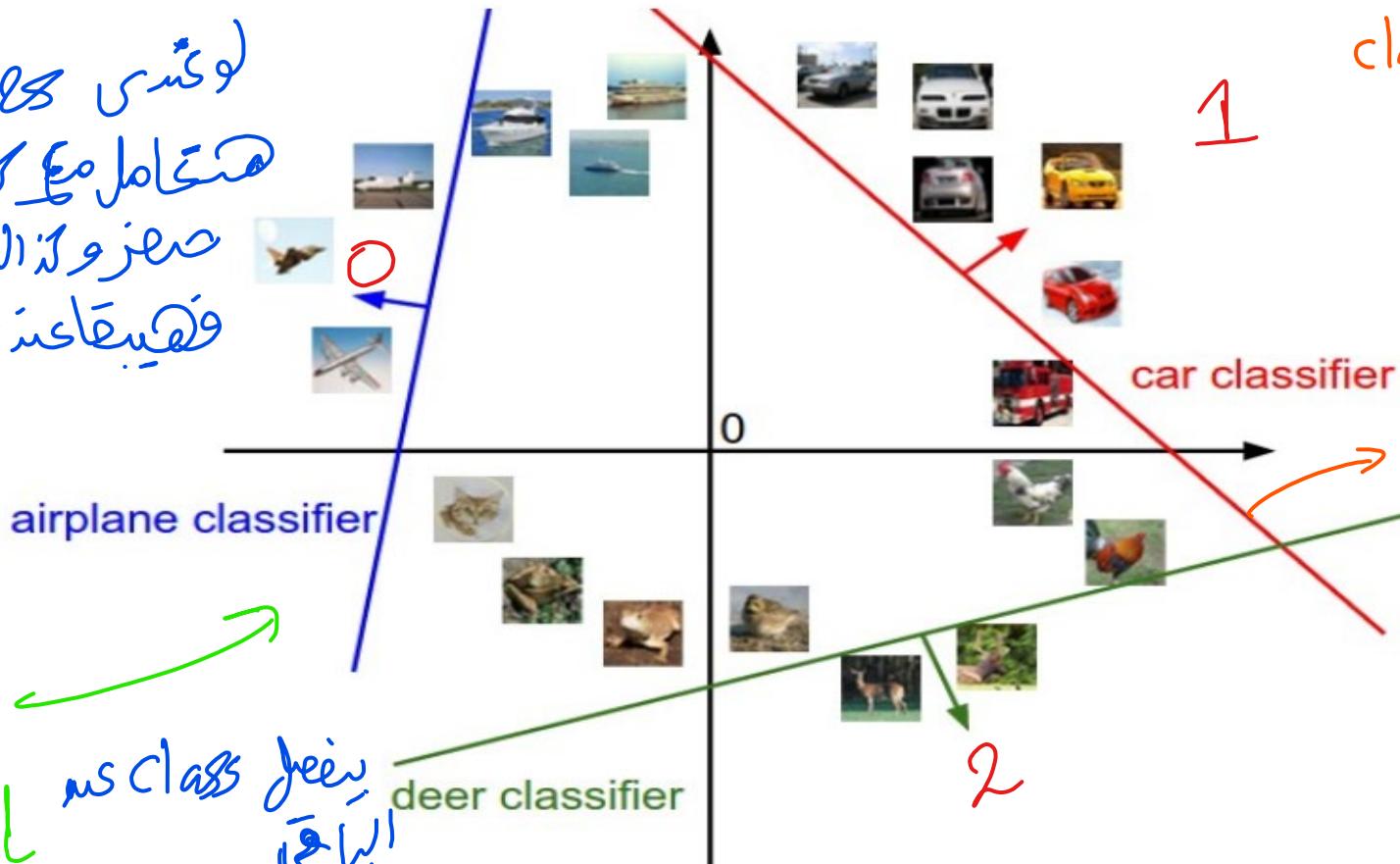
- Train a logistic regression classifier  $h_{\theta}(x)$  for each class

....., ۱, ۳, ۲ بس مگ ۱ و ۰ میں کمترین

3 classes (وئیسی کل کلاسیفایر) کا محتوا کل کلاسیفایر کے علی ۱ و ۰ باقی  
class کا محتوا کل کلاسیفایر کے علی ۱ و ۰ باقی  
فہیقانی کے لئے ۳ بھرپور مادیں

OVA

one Verses all (versus class deer) (۱) ایسی کا نامیہ و ابھی کافی ہے نامیہ تائپیہ



class 1 کا item (کوئی کوئی کا item) کا item کا item  
کا item کا item کا item  
کا item کا item کا item  
کا item کا item کا item

کل ۳ مولیکارہ کے  
line equation  
کل ۳ مولیکارہ کے  
line equation  
کل ۳ مولیکارہ کے  
line equation

# Multiclass Classification (one-vs-all)

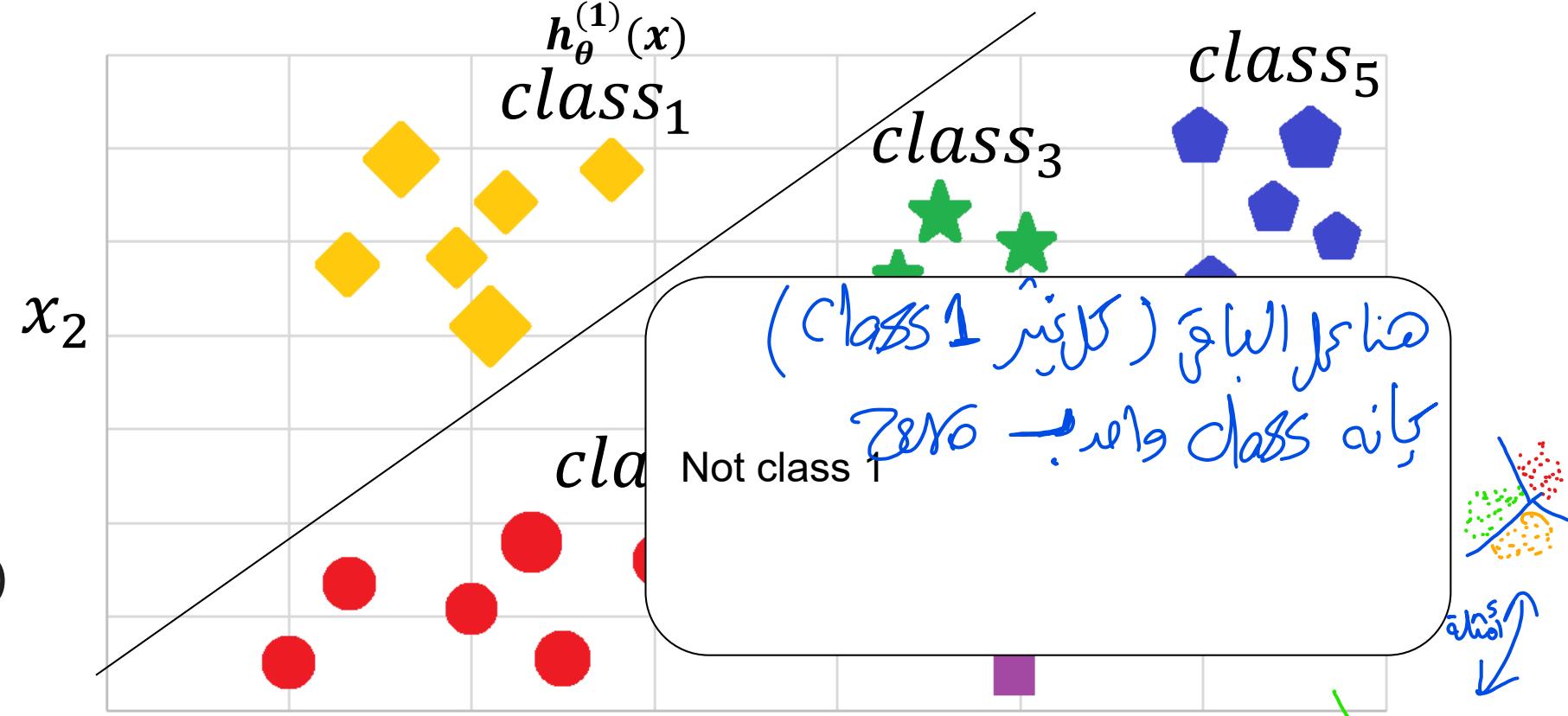
OVA

$$h_{\theta}^{(i)} = P(y = i|x; \theta)$$

$$i = 1, 2, 3 \dots$$

Pick the class  $i$  that maximize

$$\text{prediction} = \max_i (h_{\theta}^{(i)})$$



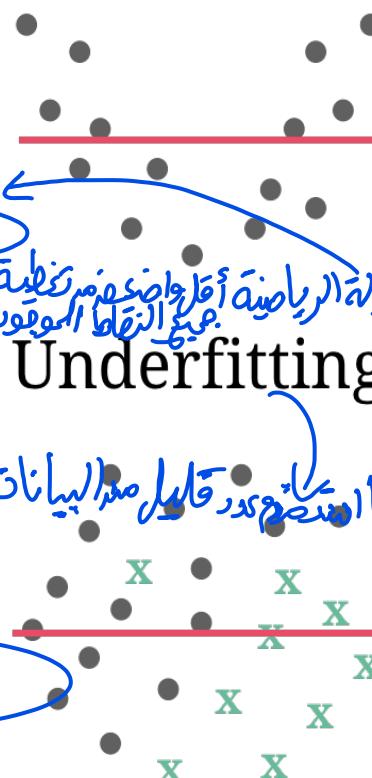
Suppose you have a multi-class classification problem with  $k$  classes (so  $y \in \{1, 2, \dots, k\}$ ). Using the one-vs.-all method, how many different logistic regression classifiers will you end up training?  $K$

OVA, all vs all  
class wise logistic regression classifier  
Technique

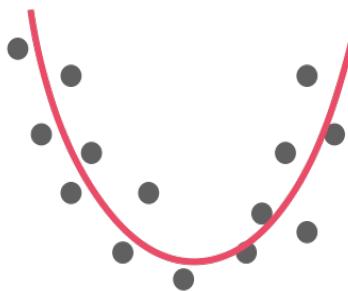
# The Problem of Overfitting

With Regression, classification, ...  
one classes بار محدود

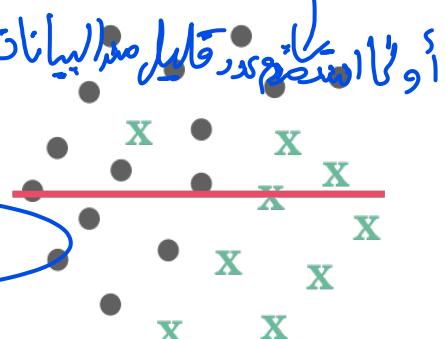
- **Underfitting**, or **high bias**, is when the form of our hypothesis function  $h$  maps poorly to the trend of the data. It is usually caused by a function that is too simple or uses too few features.
- At the other extreme, **overfitting**, or **high variance**, is caused by a hypothesis function that fits the available data but **does not generalize** well to predict new data. It is usually caused by a complicated function that creates a lot of unnecessary curves and angles unrelated to the data.



Regression



Desired



Classification



Overfitting

# Addressing overfitting

- There are two main options to address the issue of overfitting:

## 1) Reduce the number of features:

- Manually select which features to keep.
- Use a model selection algorithm.

پیوچ کردن فیچرز کم کردن

## 2) Regularization

- Keep all the features, but reduce the magnitude of parameters  $\theta_j$ .
- Regularization works well when we have a lot of slightly useful features.

برای ارزیابی برآوردهای آزمایشی Testing و آموزشی Training داده‌set می‌گیریم  
80% 20%

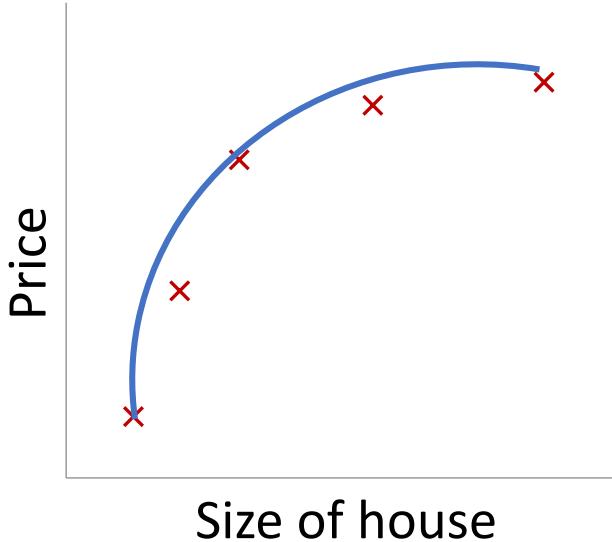
Overfitting ایجاد شود، Testing نتایج آزمایشی Training نیست

# Regularization (To avoid overfitting)

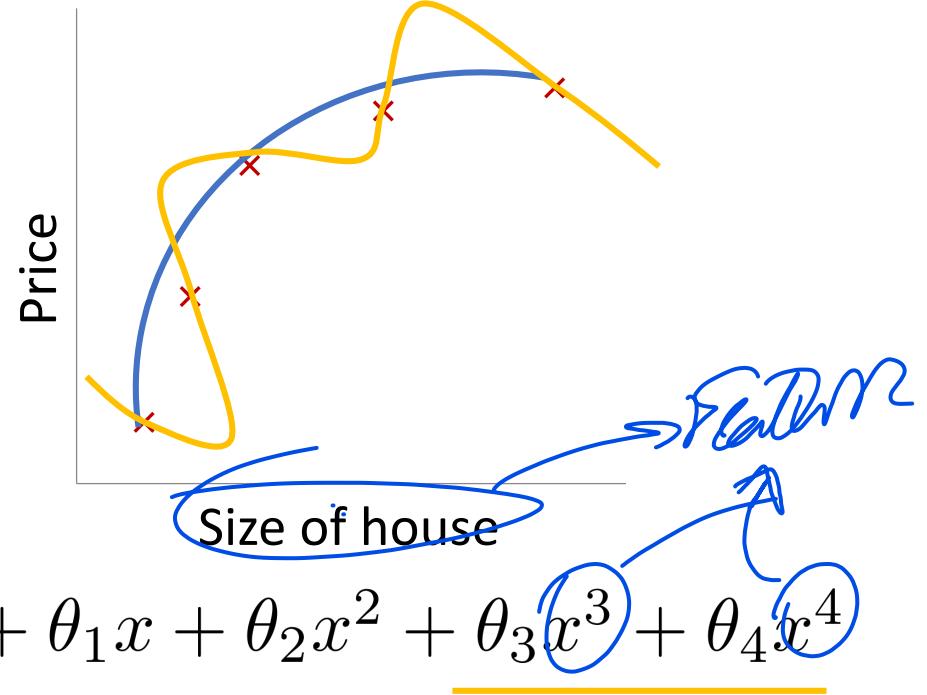
Regularization for Linear Regression

ومن المهم أن نكون حذرين من خطأ المبالغة في التعلم (Overfitting) في التعلم الآلي  
حيث يمكن أن يؤدي ذلك إلى تعلم محدود ومتضخم (Underfitting) في التعلم الآلي

# Regularization Intuition



$$\theta_0 + \theta_1 x + \theta_2 x^2$$



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 \textcolor{blue}{x^3} + \theta_4 \textcolor{blue}{x^4}$$

- 1

  - Suppose we penalize and make  $\theta_3, \theta_4$  **really small**.
  - Small values for parameters  $\theta_0, \theta_1, \dots, \theta_n$ :

- simpler hypothesis

- less prone to overfitting

میزہ-غرض میں ایکاں اور overfitting

لما جب رفاقت نايمان  
نار في Training تكنولوجیه  
صيغة لما افراخ و افراد

١. ← جاول اهل الـ  $\theta_1$  و  $\theta_2$   $\rightarrow$   $\theta_1 = \theta_2 = 0$   
پیمانه  $\theta_1 = \theta_2 = \pi$   $\rightarrow$   $\theta_1 = \theta_2 = \pi$   
و  $\theta_1 = \theta_2 = \frac{\pi}{2}$   $\rightarrow$   $\theta_1 = \theta_2 = \frac{\pi}{2}$   
و  $\theta_1 = \theta_2 = \frac{3\pi}{2}$   $\rightarrow$   $\theta_1 = \theta_2 = \frac{3\pi}{2}$

loss function  $J(\theta)$  مموجة في  
ذلك

# Regularization for linear Regression

②

- Simpler hypothesis  $\rightarrow$  small values of  $\theta_1, \theta_2, \dots, \theta_n$ .

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

cost fun factor from (1)

$\theta_0$  is no cost fun

choose small  $\theta$  to min  $J(\theta)$

$J(\theta)$  increases

- The  $\lambda$ , or lambda, is the **regularization parameter**.  $0 \leq \lambda \leq$  large number It determines how much the costs of our theta parameters are inflated.

- Example: we have 2 sets of parameters  $\theta = [1.35 \ 3.5]$  and  $\theta = [45.2 \ 75.6]$

- regularization factor
- If  $\lambda$  is chosen to be 0, the cost function act as usual with no penalty on  $\theta$   
 $\rightarrow$  choose the large ones  $\theta = [45.2 \ 75.6]$

- If  $\lambda$  is chosen to be large, small values of  $\theta$  are chosen instead of large one.  
 $\rightarrow$  choose the large ones  $\theta = [1.35 \ 3.5]$

اختيار مقدار  $\theta$  لتجنب التفتق  
تفتق overfitting يدل على انتهاج

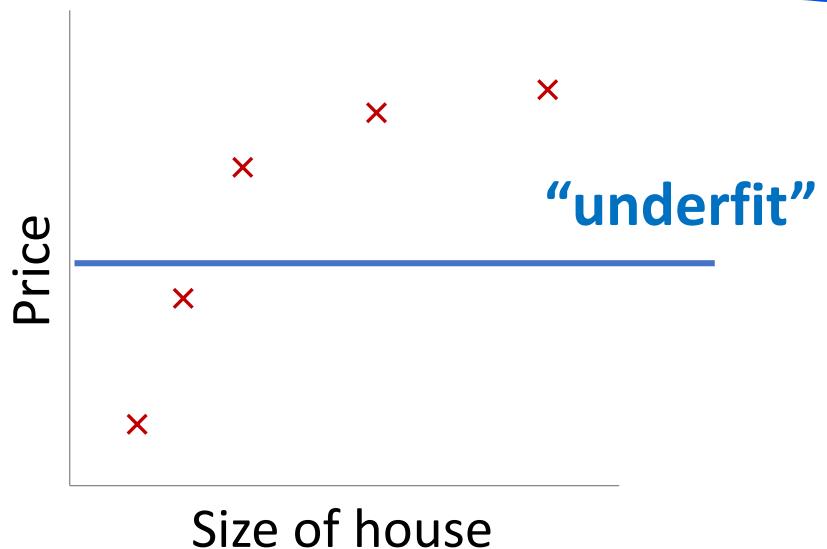
# Regularization for linear Regression

$J(\theta)$  Major

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

- What if  $\lambda$  is set to an extremely large value (perhaps too large for our problem, say  $\lambda = 10^{10}$ )?  
→ it may smooth out the function too much and cause *underfitting*.

Why?



$$\theta_1 \approx 0, \theta_2 \approx 0, \theta_3 \approx 0, \theta_4 \approx 0$$

$$\rightarrow h_{\theta}(x) = \theta_0$$

$$\theta_0 + \cancel{\theta_1}x + \cancel{\theta_2}x^2 + \cancel{\theta_3}x^3 + \cancel{\theta_4}x^4$$

# Regularization for linear Regression

- Gradient descent

Repeat { مکمل ترین و محدود نهادن } که در اینجا از اینجا که در اینجا از اینجا

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \quad j \notin \{0\}$$
$$\theta_j := \theta_j - \alpha \left[ \left( \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right) + \frac{\lambda}{m} \theta_j \right] \quad j \in \{1, 2, \dots, n\}$$

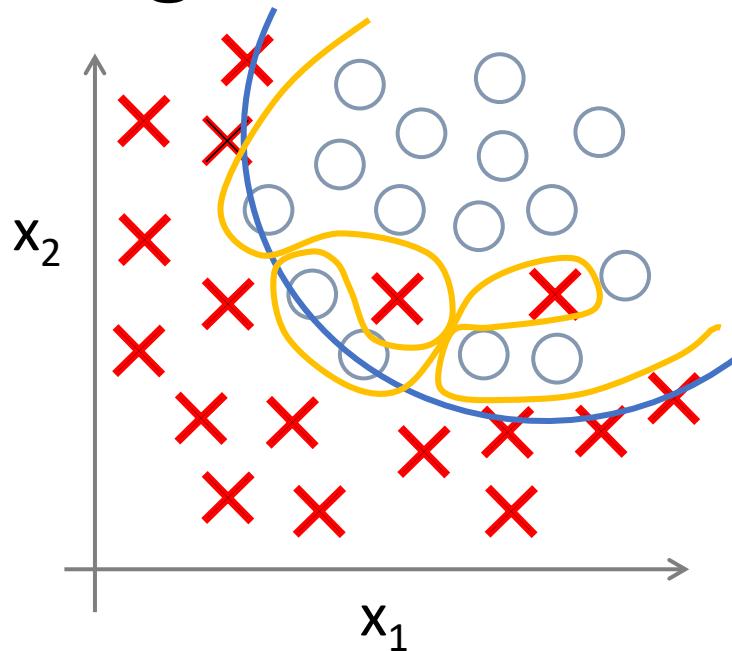
}

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

# Regularization (To avoid overfitting)

Regularization for **Logistic Regression**

# Regularization for 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:

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

**Small values of parameters**

# Regularization for Logistic Regression

- Gradient descent

Repeat {

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

$$\theta_j := \theta_j - \alpha \left[ \left( \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right) + \frac{\lambda}{m} \theta_j \right] \quad j \in \{1, 2, \dots, n\}$$

}

$\theta_0$  fiktiv

fiktiv

$j \notin \{0\}$

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

where,  $h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$

# Thanks