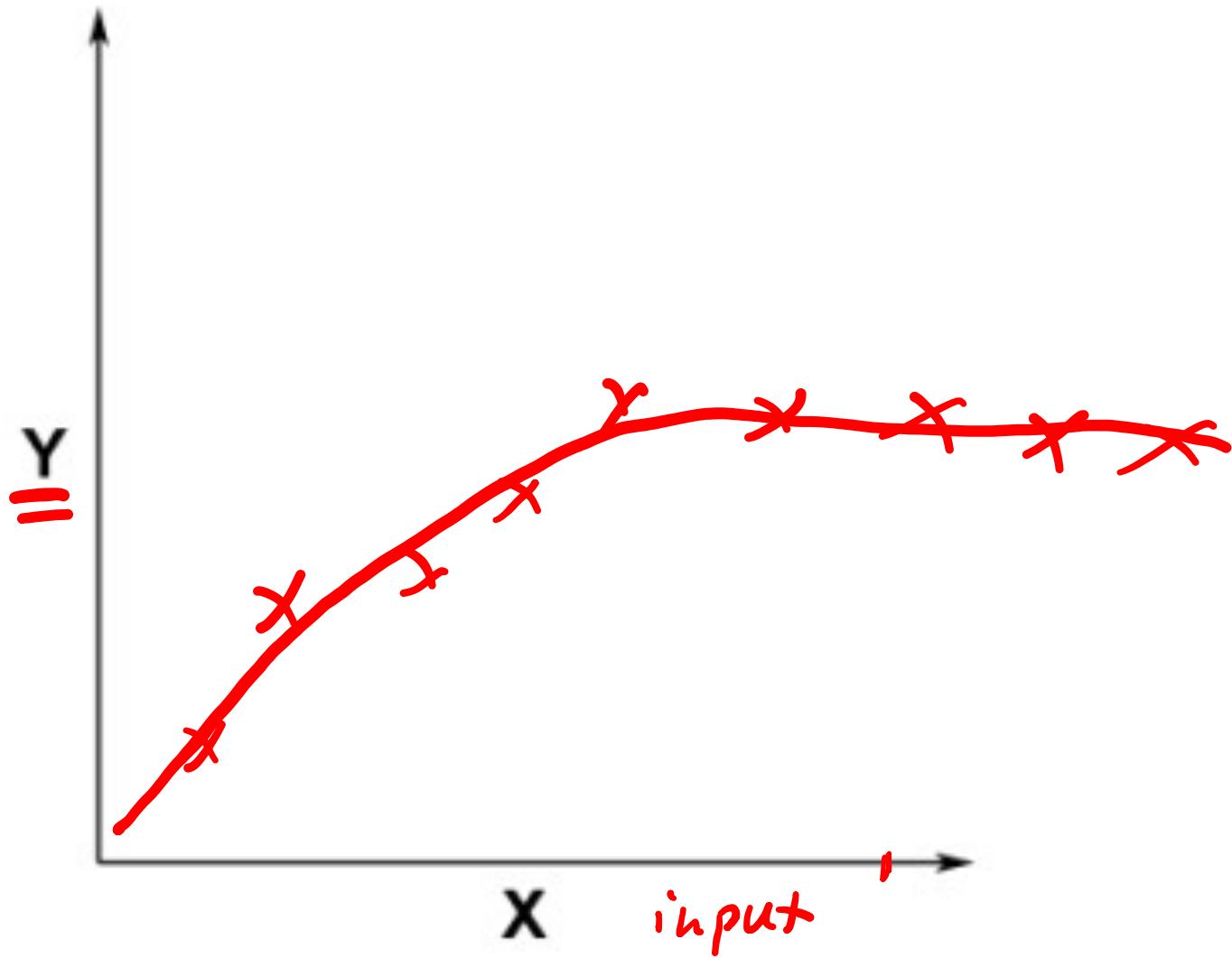


深度学习

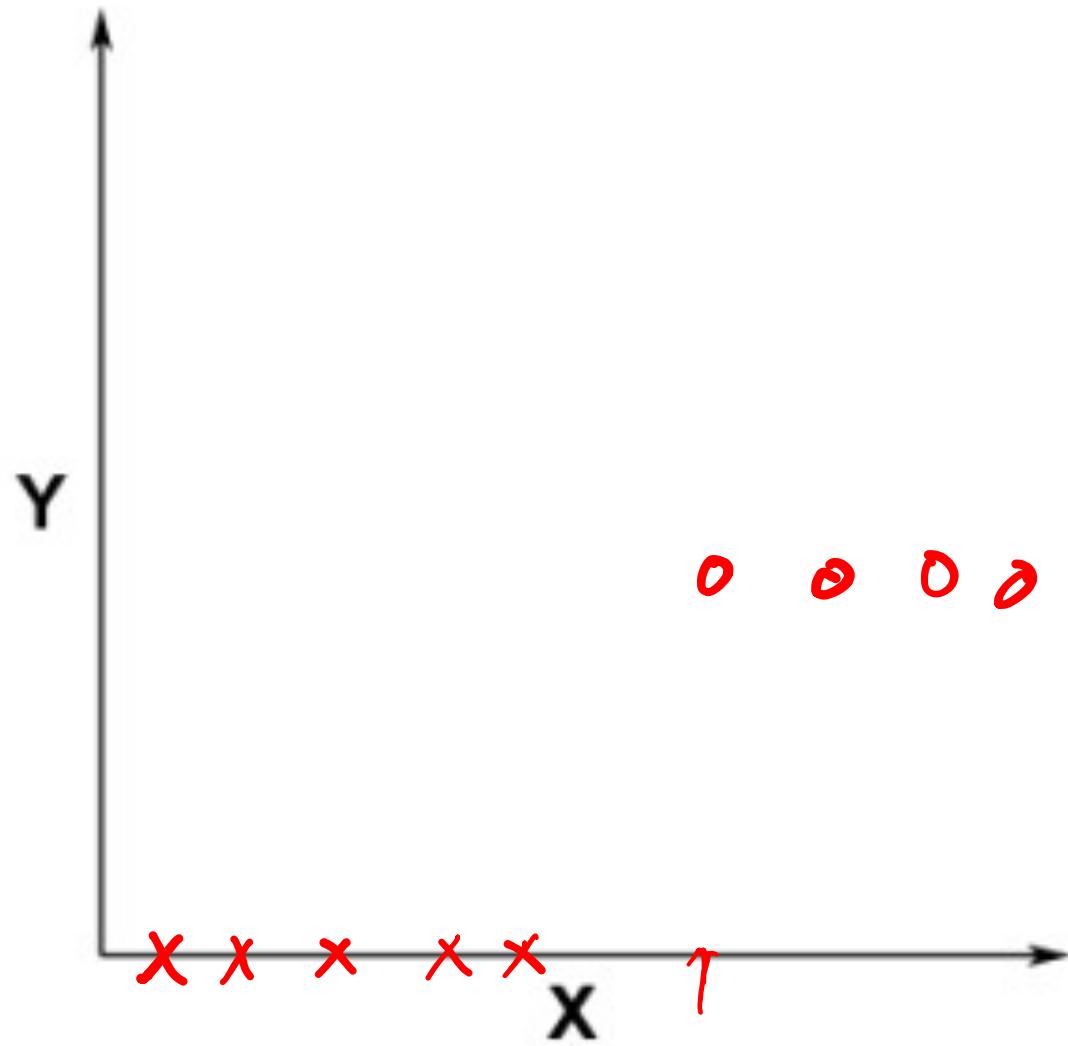
第二讲

王胤

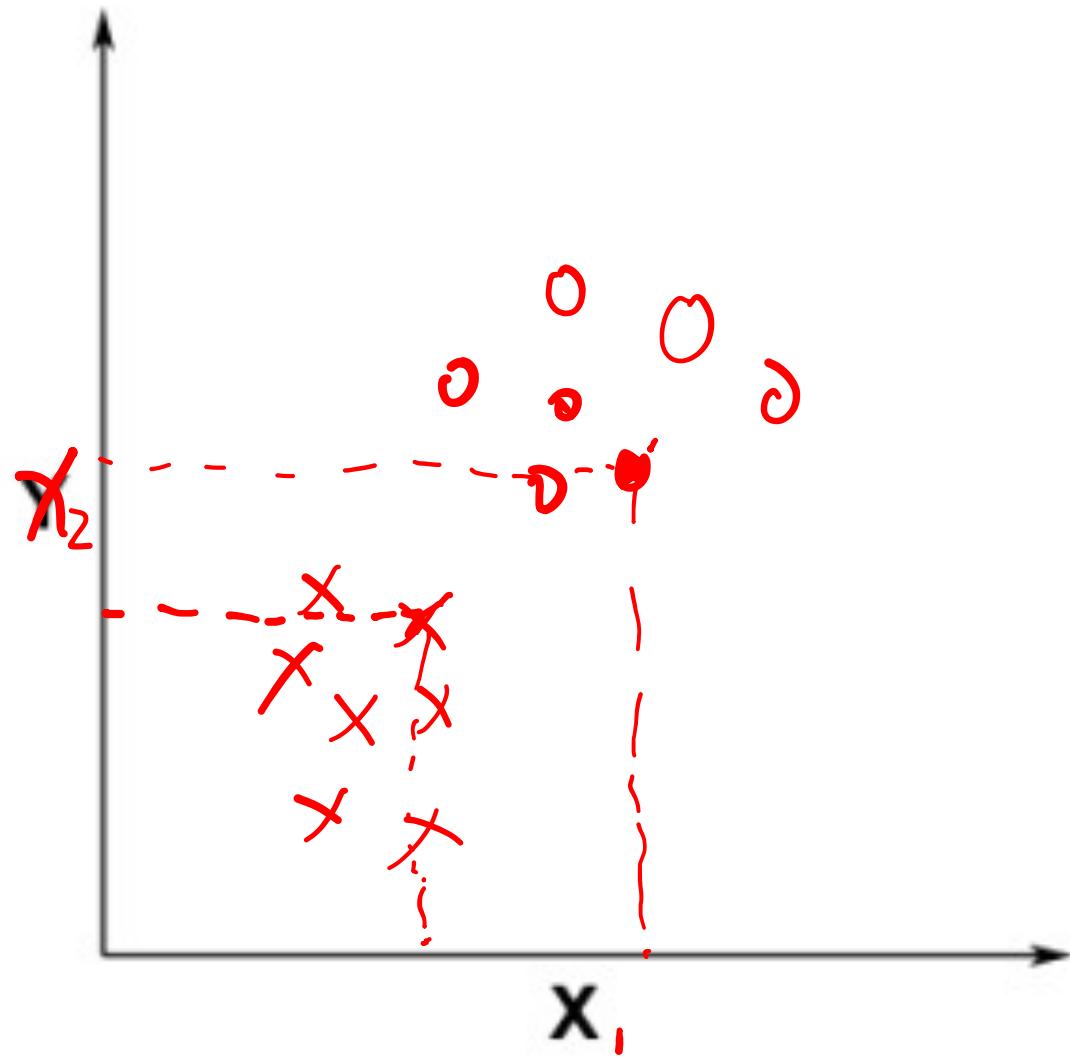
Regression



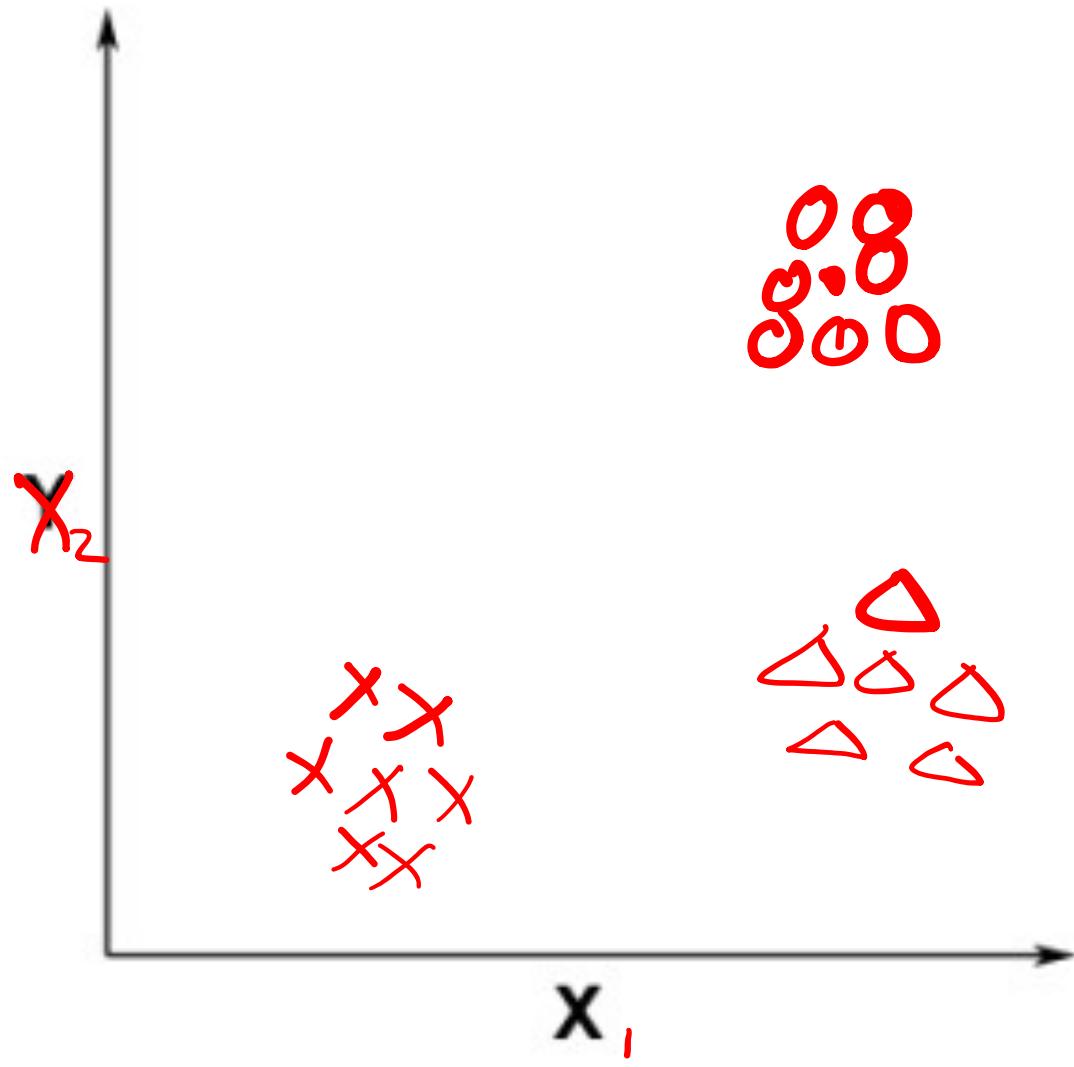
Classification



Classification



Unsupervised



Linear Regression

- Training set (this is your data set)
- Notation (*used throughout the course*)
 - m = number of training examples
 - x's = input variables / features
 - t's = output variable "target" variables
 - (x, t) - single training example
 - (x^i, t^i) - specific example (i^{th} training example)

Size in feet ² (x)	Price (\$) in 1000's (t)
2104	460
1416	232
1534	315
852	178
...	...

$m = 47$

Hypothesis, Cost Function

$$\underline{y_{\omega}(x) = \omega_0 + \omega_1 x}$$

$$C(\omega) = \frac{1}{2m} \sum (y_{\omega}(x^{(i)}) - t^{(i)})^2$$

$$\min_{\omega} C(\omega)$$

Assume $\omega_0 = 0$

$$y_\omega(x) = \omega_1 x$$

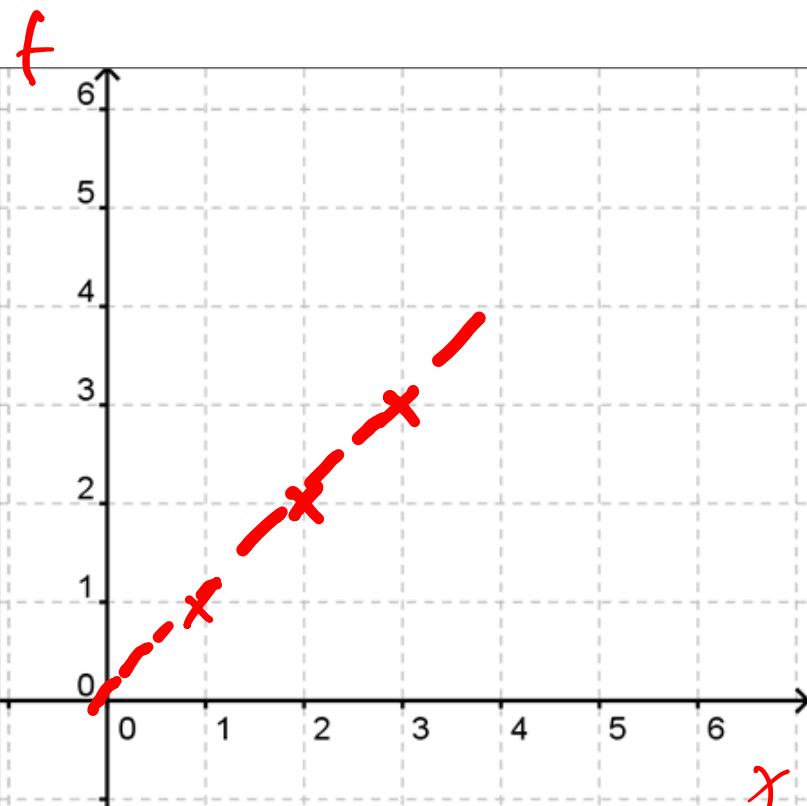
$$C(\omega) = \frac{1}{2m} \sum_{i=1}^m (\omega_1 x^{(i)} - t^{(i)})^2$$

$$\underline{C(\omega_1)} = \frac{1}{2m} \sum_{i=1}^m (y_\omega(x^{(i)}) - t^{(i)})^2$$

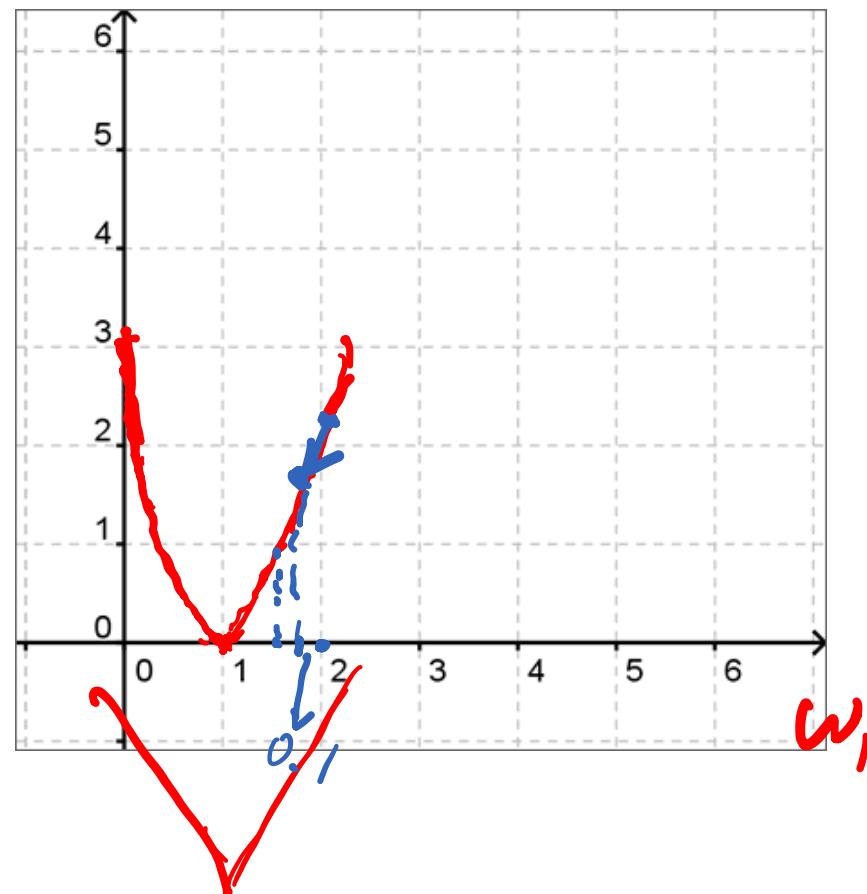
minimize $c(\omega_1)$

$y_\omega(x)$ vs. $C(\omega_1)$

$$y_\omega(x) = \omega, x$$

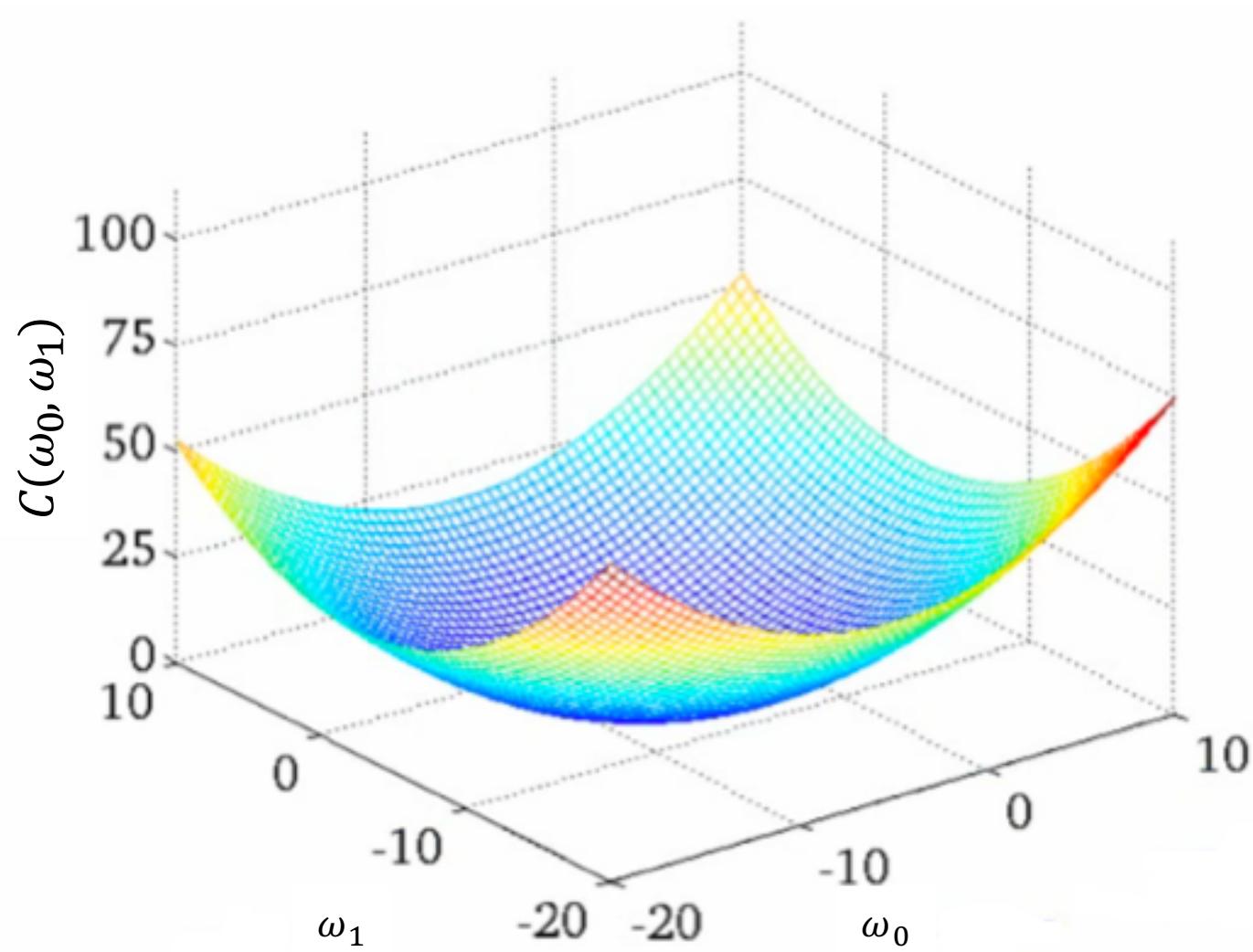


$C(\omega_1)$



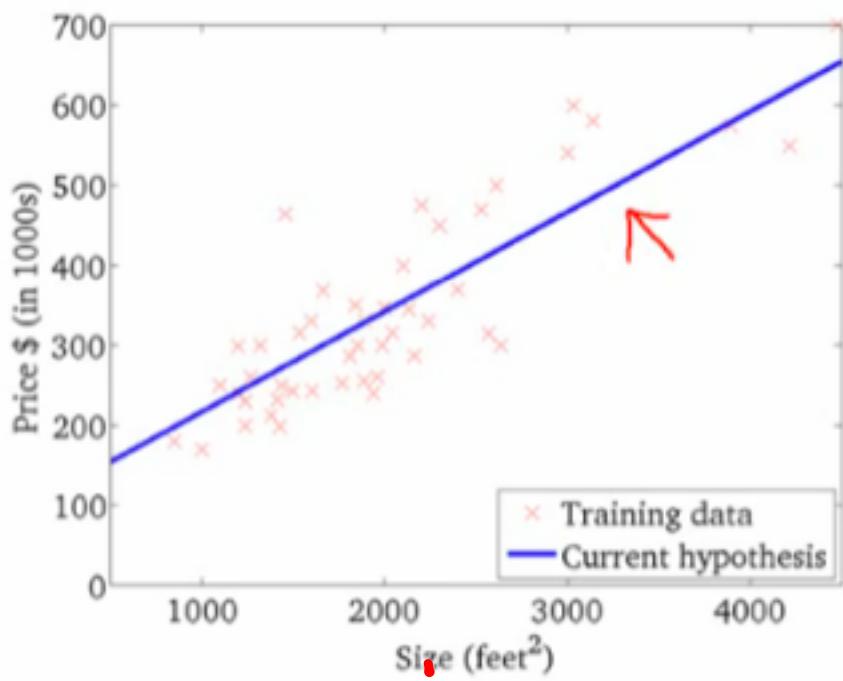
$C(\omega_0, \omega_1)$

$$y = w_0 + w_1 x$$

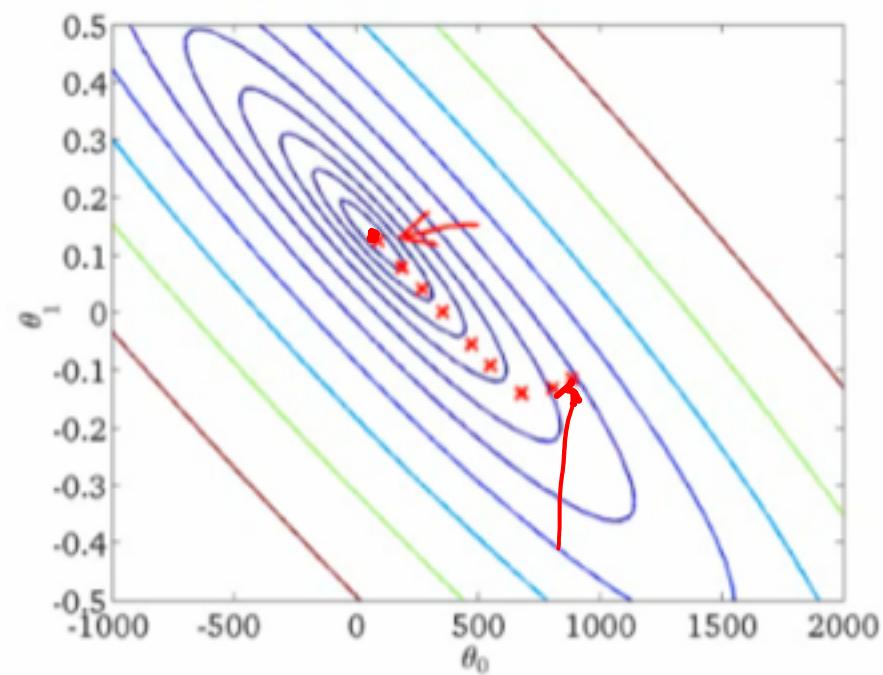


$y_\omega(x)$ vs. $C(\omega_0, \omega_1)$

$y_\omega(x)$
(for fixed ω_0, ω_1 , this is a function of x)



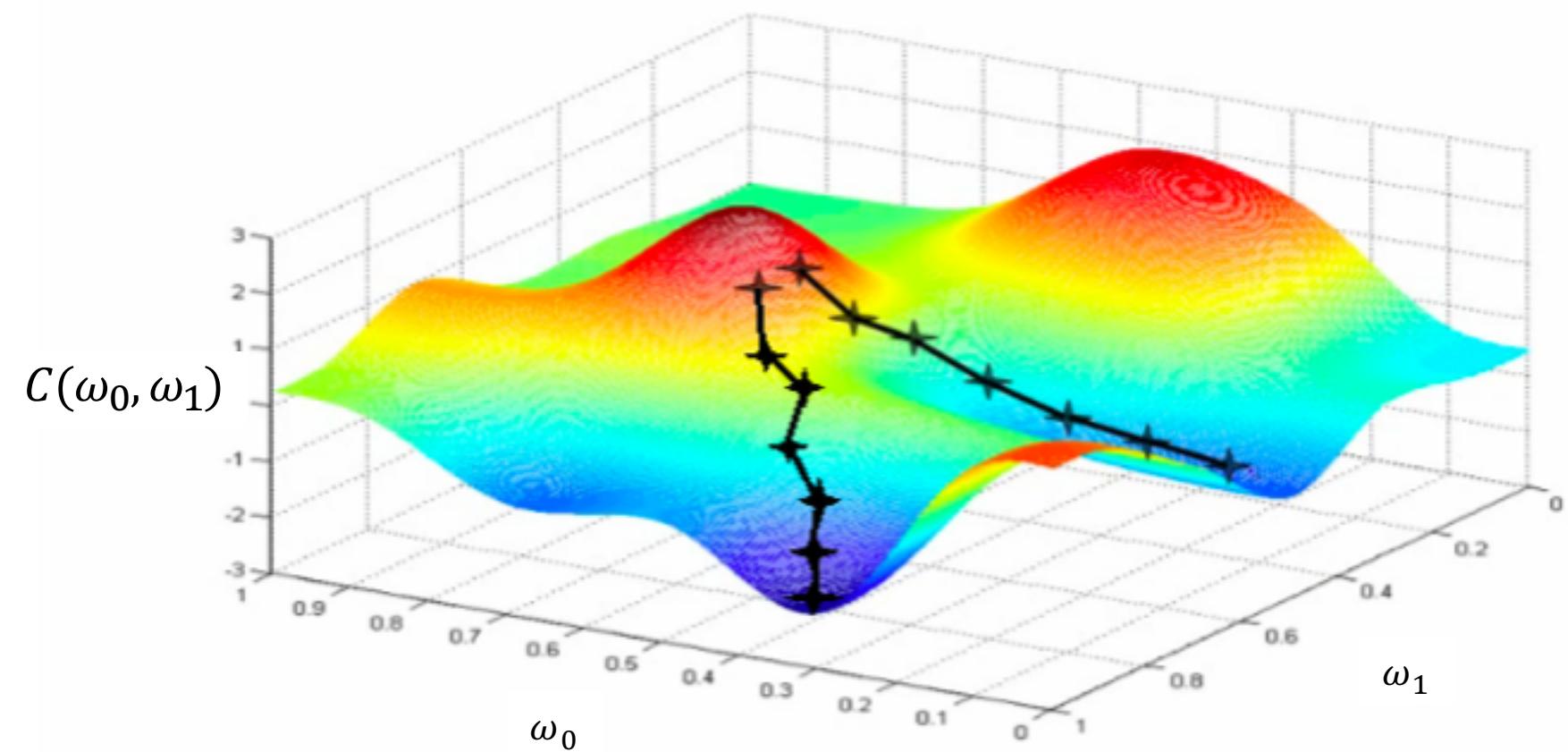
$C(\omega_0, \omega_1)$
(function of the parameters ω_0, ω_1)



Gradient Descent Algorithm

- Problem
 - We have $C(\omega_0, \omega_1)$
 - We want to get $\min C(\omega_0, \omega_1)$
- Gradient descent applies to more general functions
 - $C(\omega_0, \omega_1, \dots \omega_n)$
 - $\min C(\omega_0, \omega_1, \dots \omega_n)$
- Start with initial guesses
 - Start at 0,0 (or any other value)
 - Keeping changing ω_0 and ω_1 a little bit to try and reduce $C(\omega_0, \omega_1)$
- Each time you change the parameters, you select the gradient which reduces $C(\omega_0, \omega_1)$ the most possible
- Repeat until you converge to a local minimum

Local Minimum



A More Formal Definition

- Do the following until convergence

$$\omega_j := \omega_j - \underbrace{\left[\alpha \frac{\partial}{\partial \omega_j} C(\omega_0, \omega_1) \right]}_{\text{DW}}$$

(for j = 0 and j = 1)

- There is a subtlety about how this gradient descent algorithm is implemented

$$\text{temp0} := \omega_0 - \alpha \frac{\partial}{\partial \omega_0} C(\omega_0, \omega_1)$$

$$\text{temp1} := \omega_1 - \alpha \frac{\partial}{\partial \omega_1} C(\omega_0, \omega_1)$$

$$\omega_0 := \text{temp0}$$

$$\omega_1 := \text{temp1}$$

$$\omega_j = \omega_j - \alpha \Delta \omega_j$$

↑ ↑

Linear Regression Case

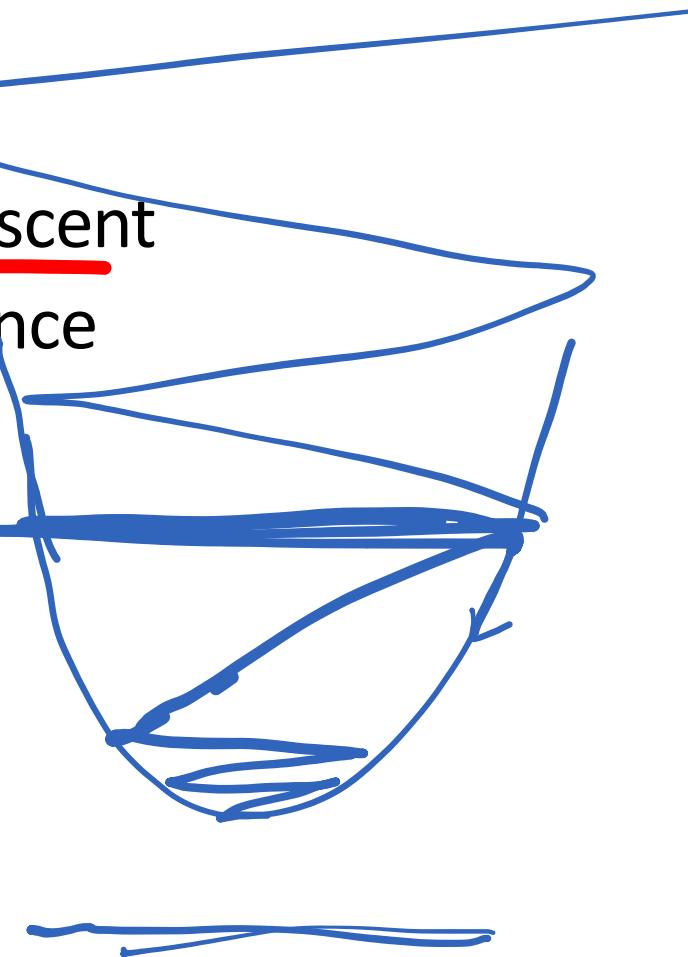
$$\begin{aligned}\frac{\partial}{\partial \omega_j} C(\omega_0, \omega_1) &= \frac{\partial}{\partial \omega_j} \frac{1}{2m} \sum \left(\underline{\omega_0 + \omega_1 x^i} - \underline{t^i} \right)^2 \\ &= \frac{1}{m} \sum \frac{\partial}{\partial \omega_j} (\underline{\omega_0 + \omega_1 x^i}) \cdot (\underline{y_{(x)} - t})\end{aligned}$$

$$\cancel{\omega_0} \rightarrow \omega_0 - \cancel{2} \sum (y_{(x)} - t)$$

$$\omega_i \rightarrow \omega_i - \cancel{2} \left(\sum_j (y_{(x)} - t) \cdot x_i^j \right)$$

Practical Issues

- Global vs. local minimum
- Single vs. batch gradient descent
- Learning rate and convergence



Multiple Variables

- Multiple variables = multiple features
- In original version we had
 - X = house size, use this to predict
 - t = house price
- If in a new scheme we have more variables (such as number of bedrooms, number floors, age of the home)
 - x_1, x_2, x_3, x_4 are the four features
 - x_1 - size (feet squared)
 - x_2 - Number of bedrooms
 - x_3 - Number of floors
 - x_4 - Age of home (years)
 - t is the output variable (price)

More Notation

• n

- number of features ($n = 4$)

• m

- number of examples (i.e. number of rows in a table)

• x^i

- vector of the input for an example (so a vector of the four parameters for the i^{th} input example)
- i is an index into the training set
- So
 - x is an n -dimensional feature vector
 - x^3 is, for example, the 3rd house, and contains the four features associated with that house

• x_j^i

- The value of feature j in the i^{th} training example
- So
 - x_2^3 is, for example, the number of bedrooms in the third house

Hypothesis

- Previously our hypothesis took the form;
 - $y_{\omega}(x) = \omega_0 + \omega_1 x$
 - Here we have two parameters (theta 1 and theta 2) determined by our cost function
 - One variable x
- Now we have multiple features
 - $y_{\omega}(x) = \omega_0 + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + \omega_4 x_4$
- For example
 - $y_{\omega}(x) = 80 + 0.1x_1 + 0.01x_2 + 3x_3 - 2x_4$
 - An example of a hypothesis which is trying to predict the price of a house
 - Parameters are still determined through a cost function

Vector Form

- For convenience of notation, $x_0 = 1$
- Considering this, hypothesis can be written

$$y_{\omega}(x) = \omega_0 + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + \omega_4 x_4$$

- More conveniently

$$y_{\omega}(\underline{x}) = \underline{\omega^T x}$$

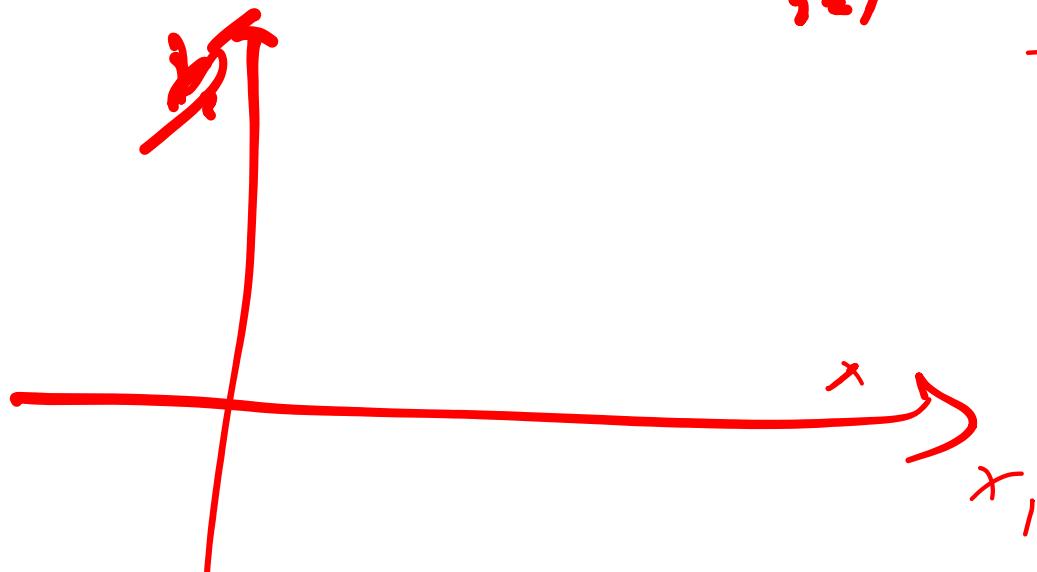
$$Y = \omega^T X$$

Gradient Descent For Multiple Variables

- Our cost function is

$$C(\omega_0, \omega_1, \dots, \omega_n) = \frac{1}{2m} \sum_{i=1}^m (y_\omega(x^{(i)}) - t^{(i)})^2$$

$$\omega_j = \omega_j - \alpha \frac{1}{m} \sum_{i=1}^m (y_\omega(x^{(i)}) - t^{(i)}) \cdot x_j$$



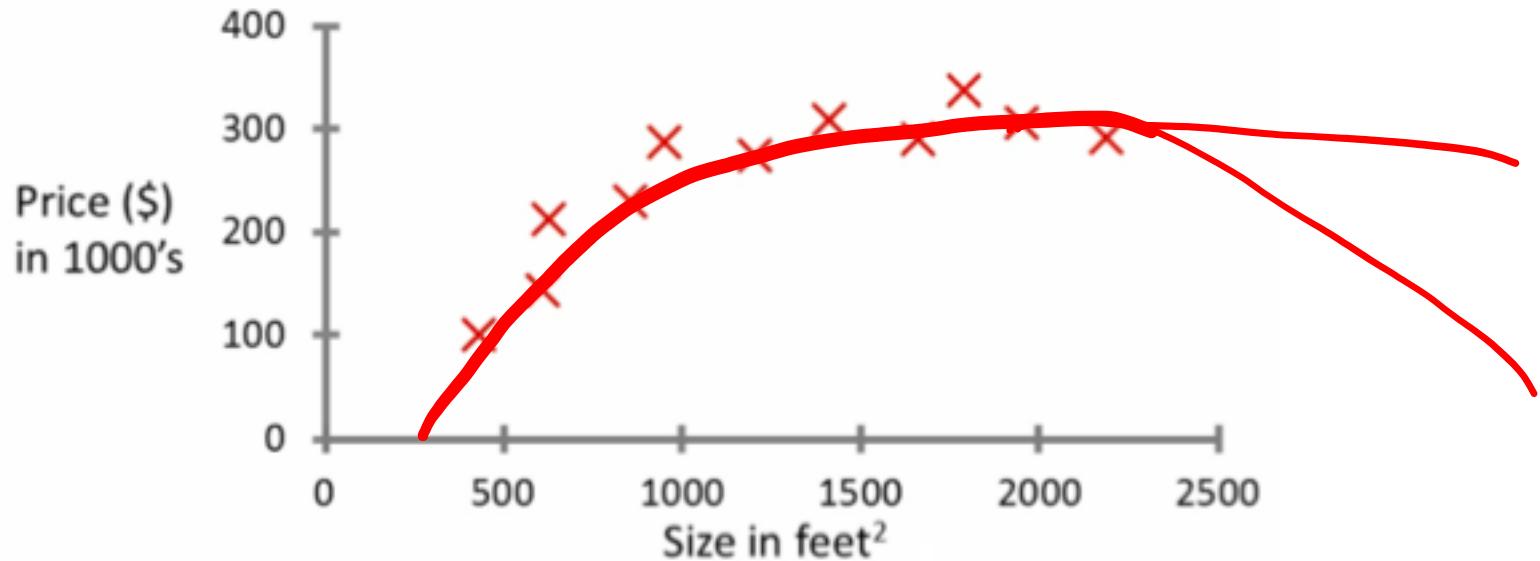
Practical Issues

- Feature scaling
- Mean normalization
- Learning rate



Polynomial Regression

Housing price prediction.



$$\begin{matrix} x_1 \\ x_1^2 \\ \vdots \\ x_2 \end{matrix}$$

Normal Equation (Least Square)

$$\min C(\omega_0, \omega_i)$$

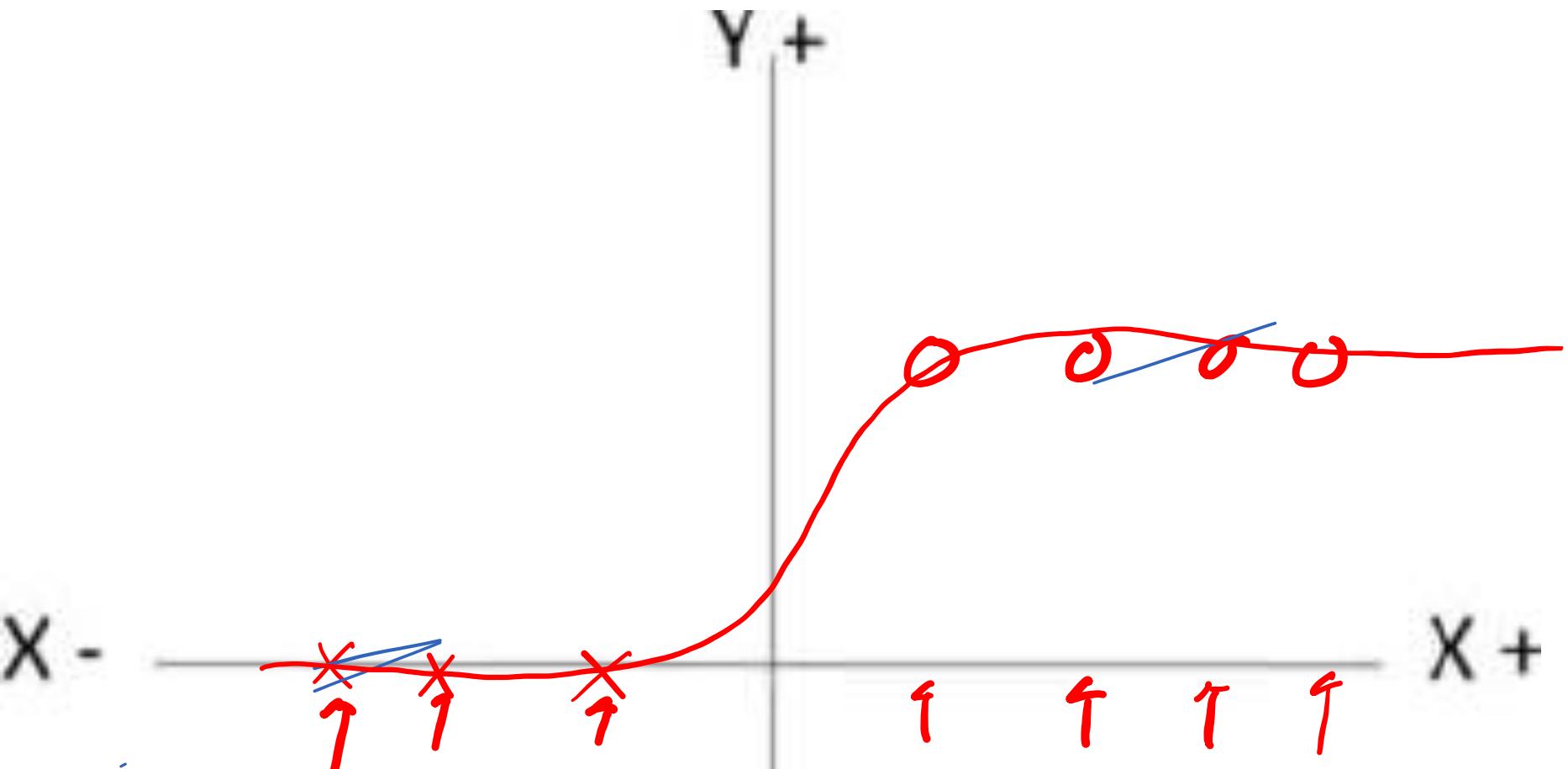
$$\sum_{i=0}^m$$

$O(m)$

$$\omega = \underbrace{(X^T X)^{-1} X^T t}_{O(n^3)}$$

Logistic Regression

Sigmoid function



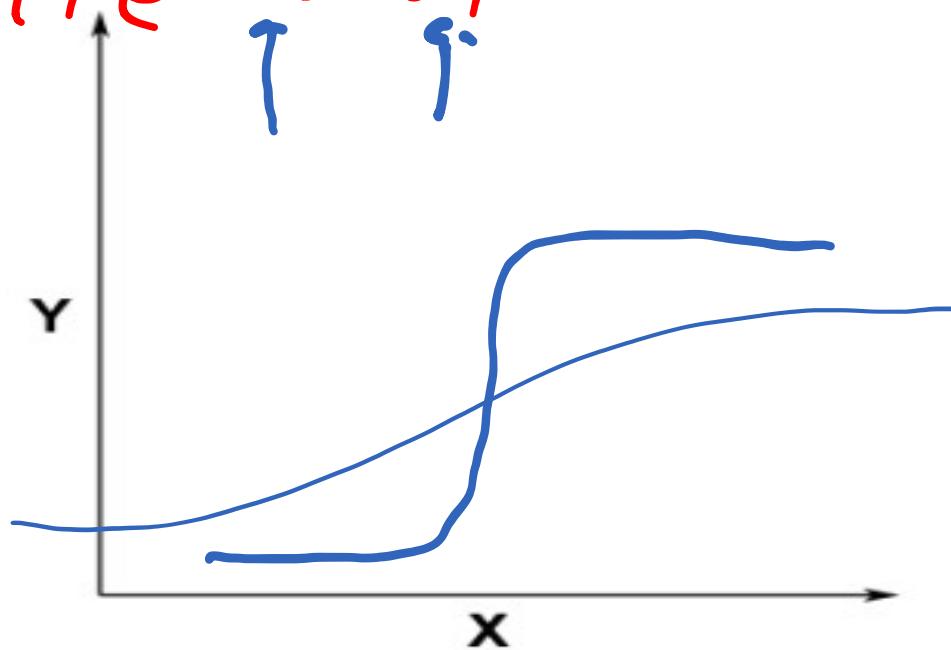
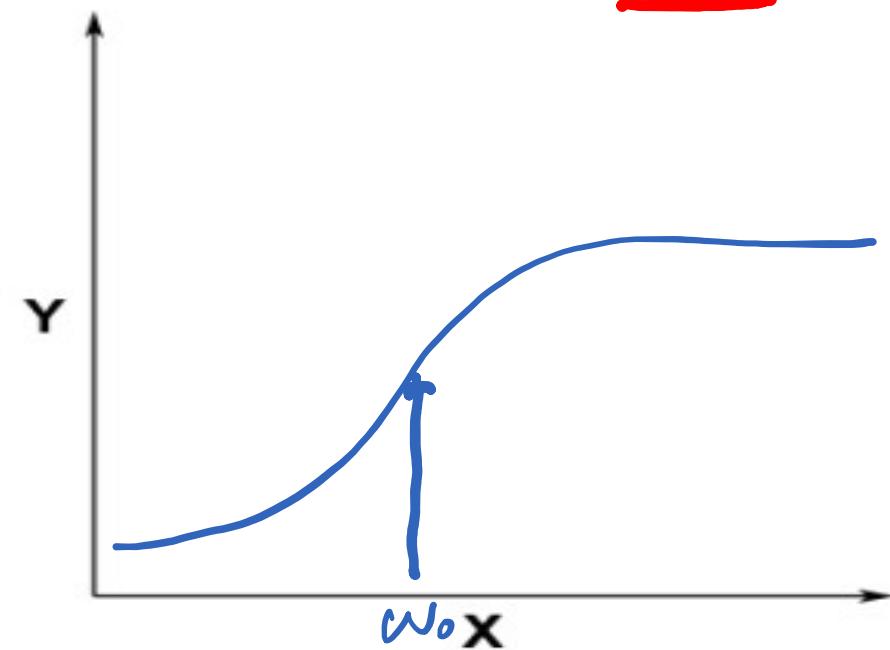
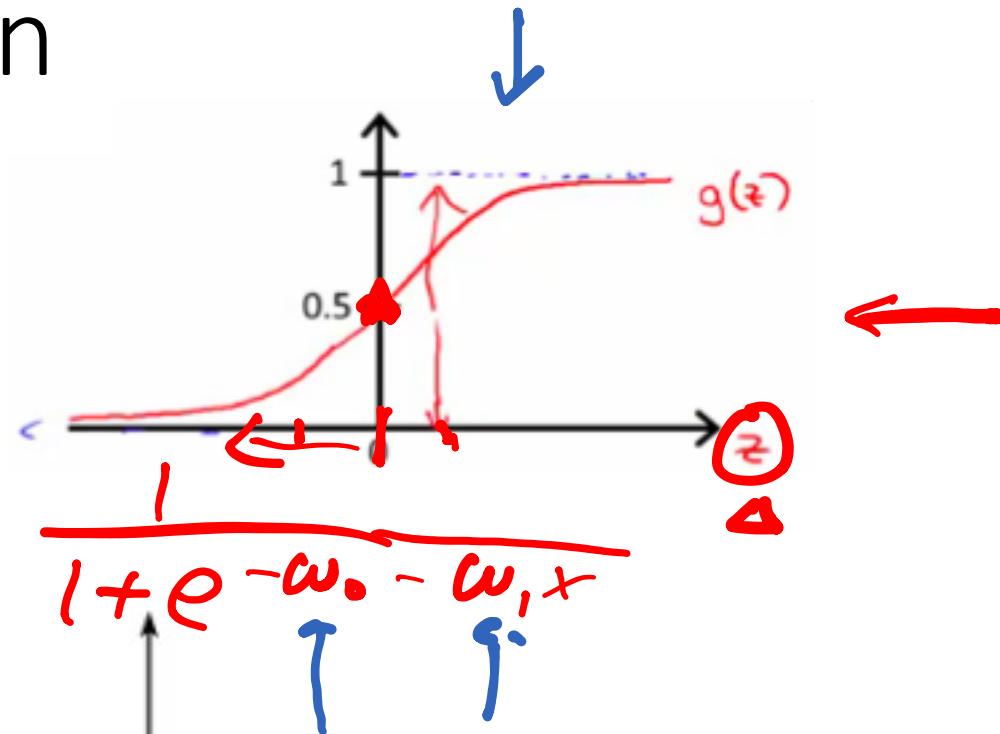
Logistic Function

$0 \sim \infty$

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

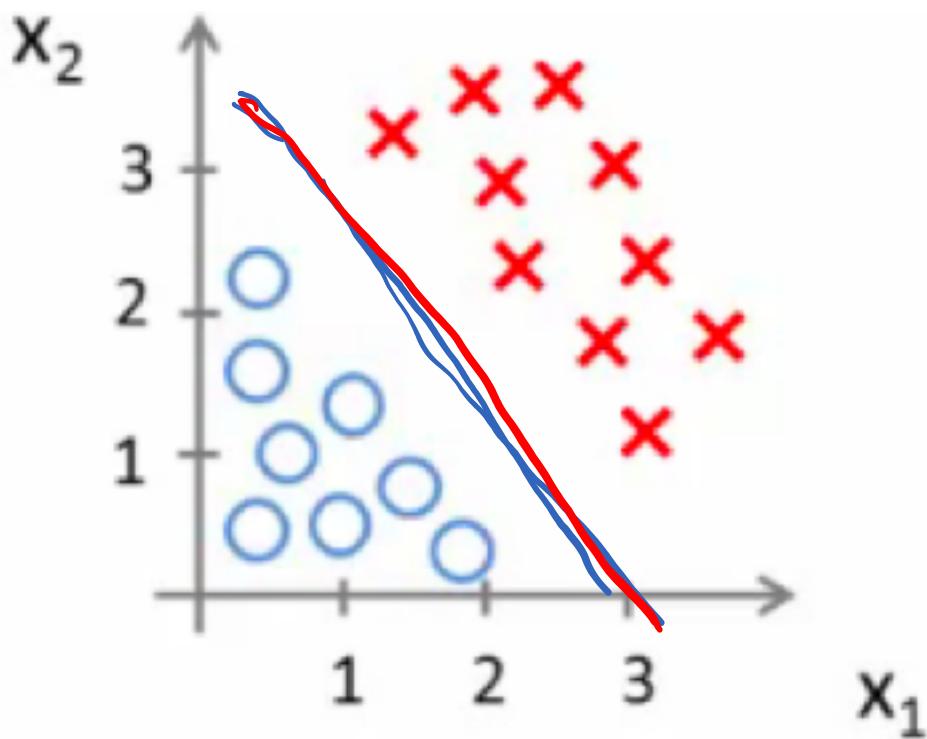


$$y_{\omega}(x) = \frac{1}{1 + e^{-\omega^T x}} = \frac{1}{1 + e^{-\omega_0 - \omega_1 x}}$$



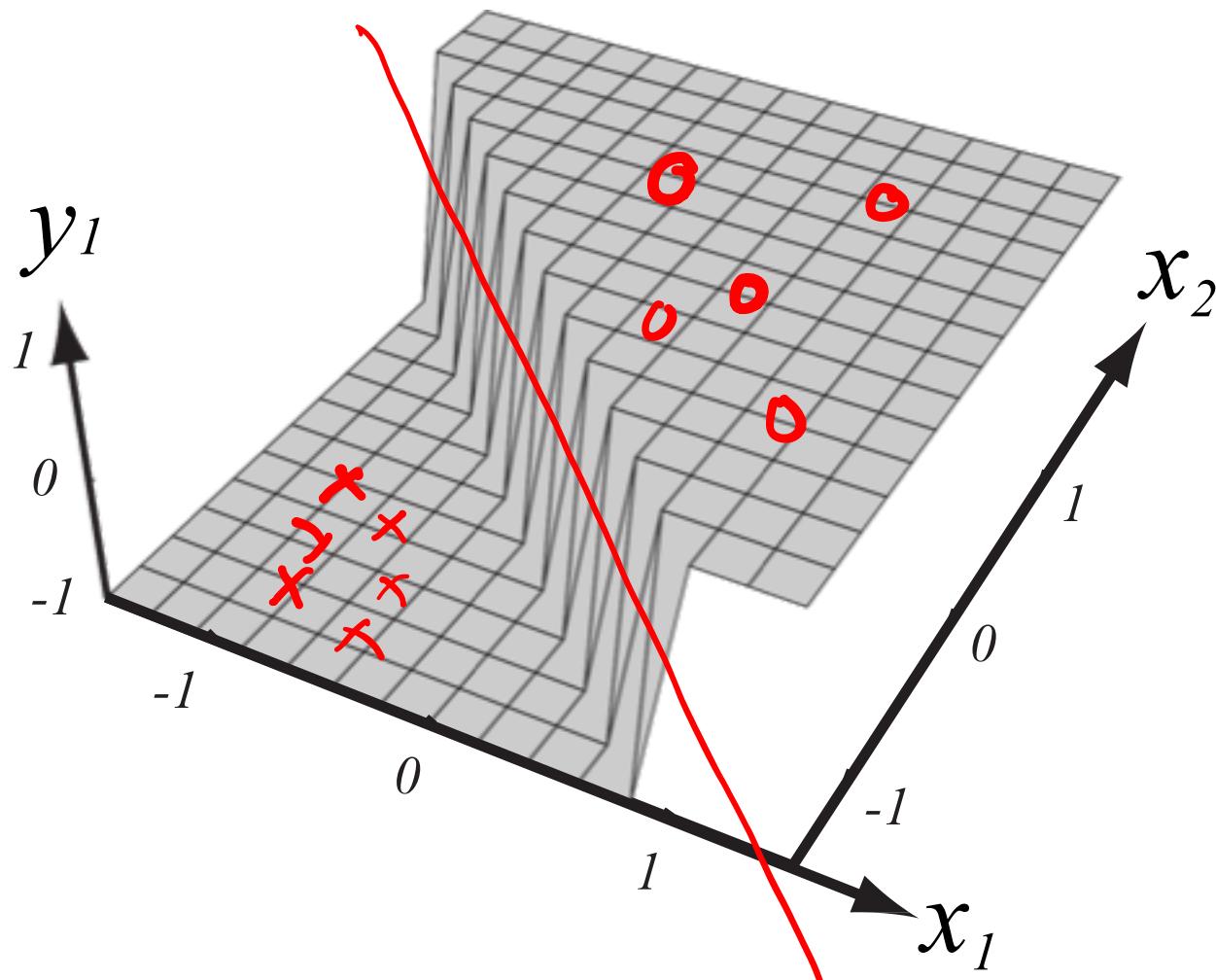
2-Dimension

$$\frac{1}{1 + e^{-w_0 + w_1 x_1 + w_2 x_2}}$$

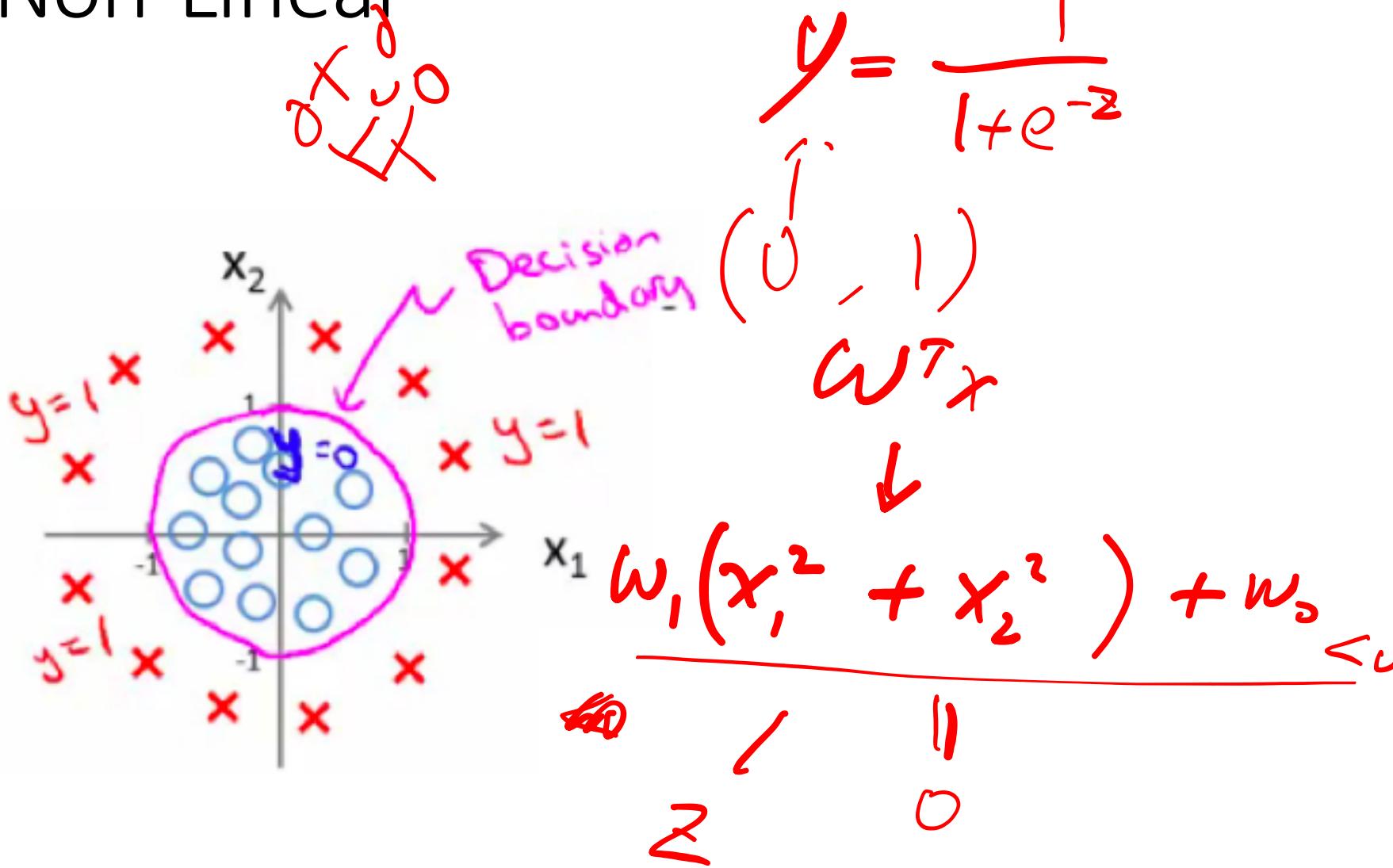


2-Dimension

$\geq w^T x$



Non-Linear



Cost Function

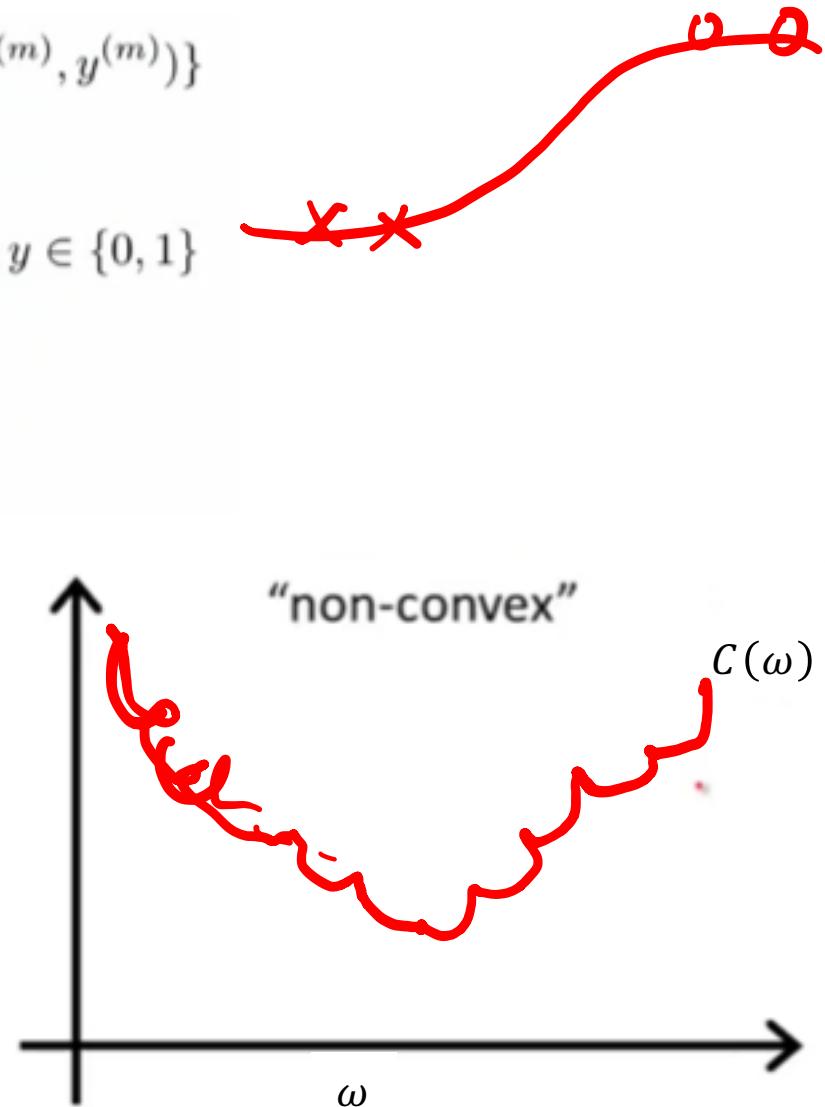
Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

m examples

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{bmatrix} \quad x_0 = 1, y \in \{0, 1\}$$

$$y_{\omega}(x) = \frac{1}{1 + e^{-\omega^T X}}$$

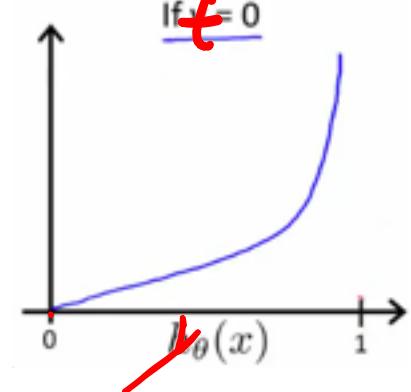
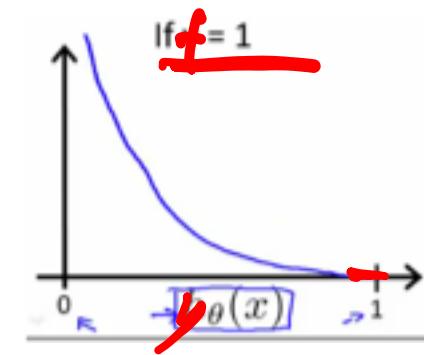
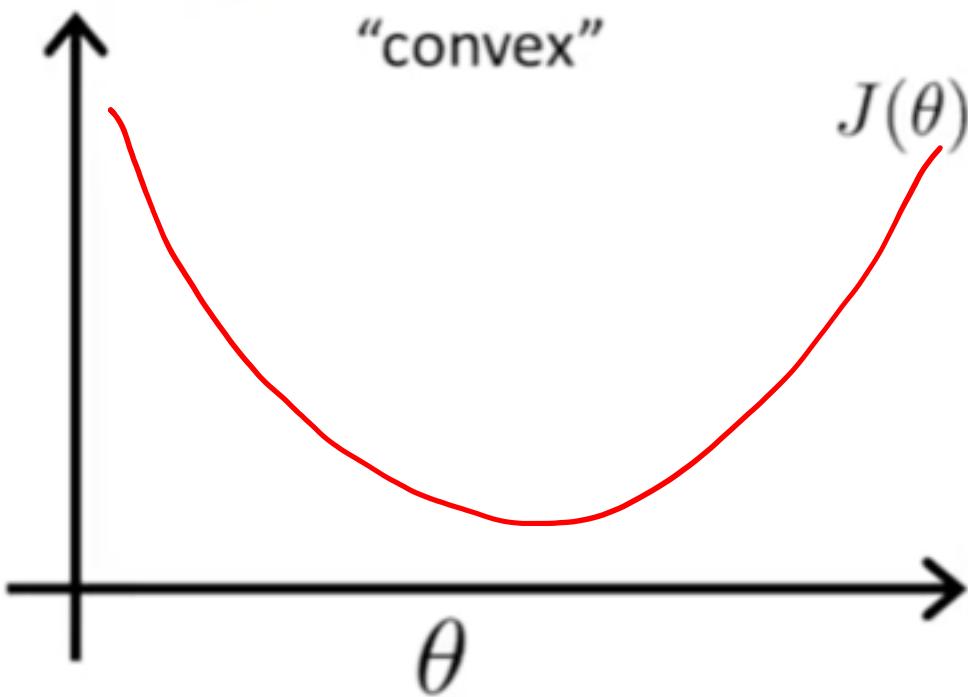
$$C(\omega) = \frac{1}{2m} \sum_{i=1}^m (y_{\omega}(x^{(i)}) - t^{(i)})^2$$



A Convex Cost Function

$$C(\omega) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(y_\omega(x), t)$$

$$\text{Cost}(y_\omega(x), t) = \begin{cases} -\log(1 - y_\omega(x)) & \text{if } t = 0 \\ -\log(y_\omega(x)) & \text{if } t = 1 \end{cases}$$



A Simpler Form

$$C(\omega) = \frac{1}{m} \sum_{i=1}^m Cost(y_\omega(x), t)$$

Cross entropy

$$Cost(y_\omega(x), t) = \begin{cases} -\log(1 - y_\omega(x)) & \text{if } t = 0 \\ -\log(y_\omega(x)) & \text{if } t = 1 \end{cases}$$

$$\begin{array}{ll} t_0 = 1, & t_1 = 0 \\ t_1 = 1, & t_0 = 0 \end{array}$$

$$C(\omega) = -\frac{1}{m} \sum_{i=1}^m t^{(i)} \log(y_\omega(x^{(i)})) + (1 - t^{(i)}) \log(1 - y_\omega(x^{(i)}))$$

- Why this cost function?

$$C(\omega) = -\frac{1}{m} \sum_i t^{(i)} \log y^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^L t_k^{(i)} \cdot \log y_k^{(i)}$$

$$y_\omega(x) = y_\omega(x) \quad y_0(x) = 1 - y_\omega(x)$$

Gradient Descent

$$C(\omega) = -\frac{1}{m} \sum_{i=1}^m t^{(i)} \log(y_\omega(x^{(i)})) + (1 - t^{(i)}) \log(1 - y_\omega(x^{(i)}))$$

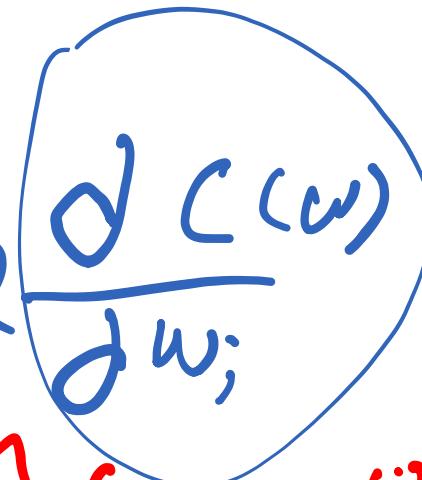
$$y_\omega(x) = \frac{1}{1 + e^{-\omega^T x}}$$

$$\omega_j = \omega_j - \alpha$$

$$\omega_j = \omega_j - \alpha \frac{1}{m} \sum_{i=1}^m (y_\omega(x^{(i)}) - t^{(i)}) \cdot x^{(i)}$$

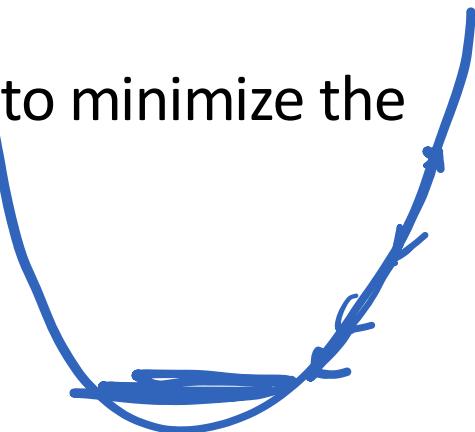
$$\omega_x \downarrow$$

$$\frac{1}{1 + e^{-\omega_x}}$$

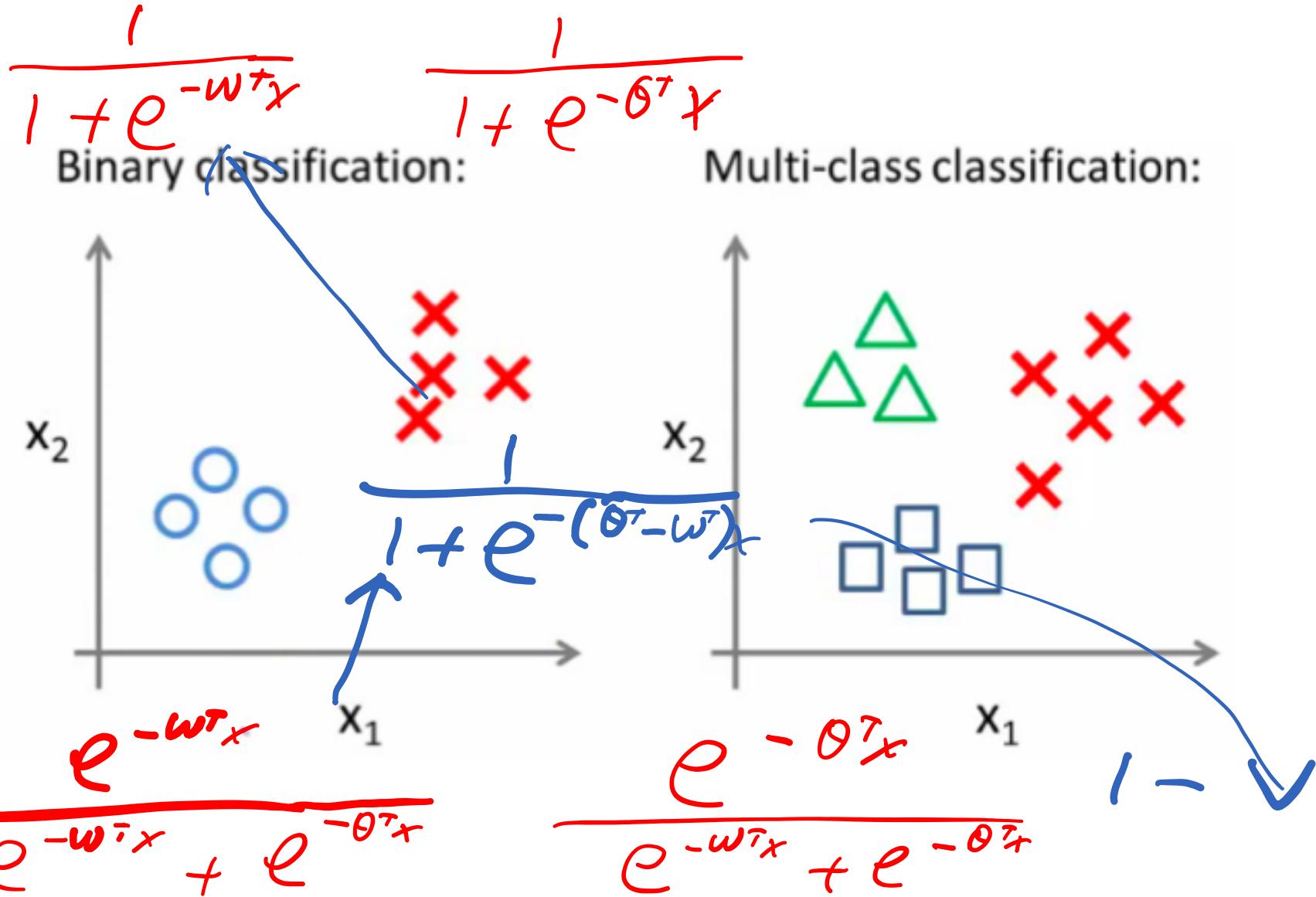


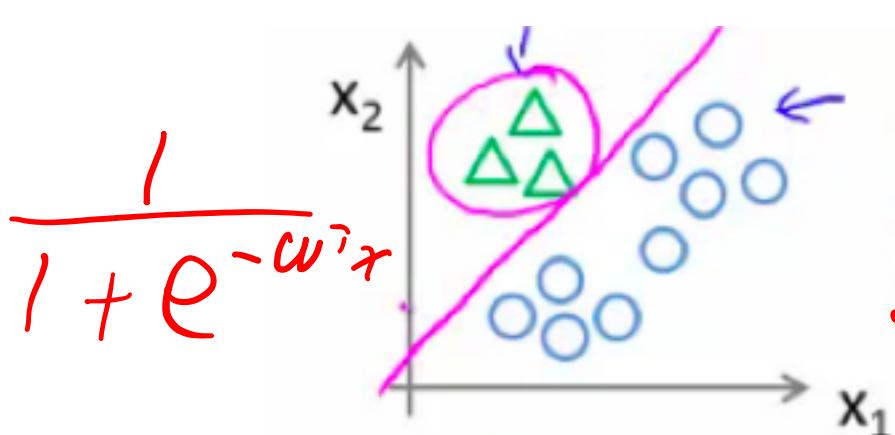
Advanced Optimization

- So, we must;
 - Supply code to compute $C(\omega)$ and the derivatives
 - Then plug these values into gradient descent
- Alternatively, instead of gradient descent to minimize the cost function we could use
 - **Conjugate gradient**
 - **BFGS** (Broyden-Fletcher-Goldfarb-Shanno)
 - **L-BFGS** (Limited memory - BFGS)
- These are *very* complicated algorithms
 - **Advantages**
 - No need to manually pick alpha (learning rate)
 - Have a clever inner loop (line search algorithm) which tries a bunch of alpha values and picks a good one
 - Often faster than gradient descent
 - Do more than just pick a good learning rate
 - Can be used successfully without understanding their complexity
 - **Disadvantages**
 - Could make debugging more difficult
 - Should not be implemented themselves
 - Different libraries may use different implementations - may hit performance



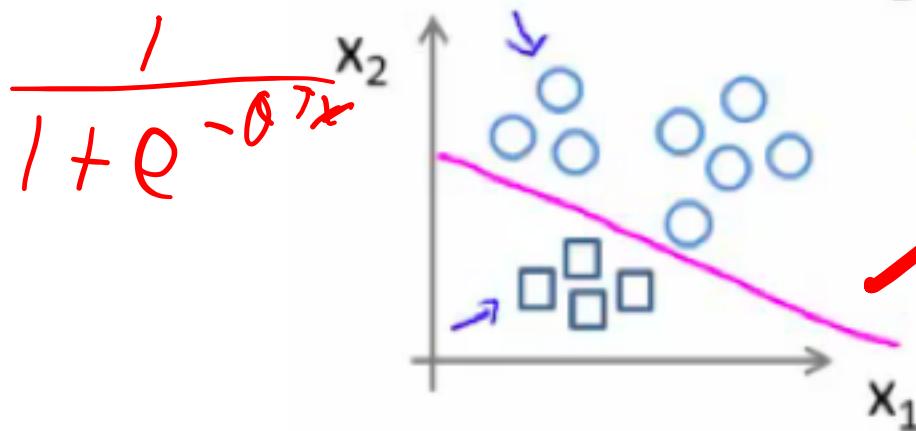
Multiclass Classification



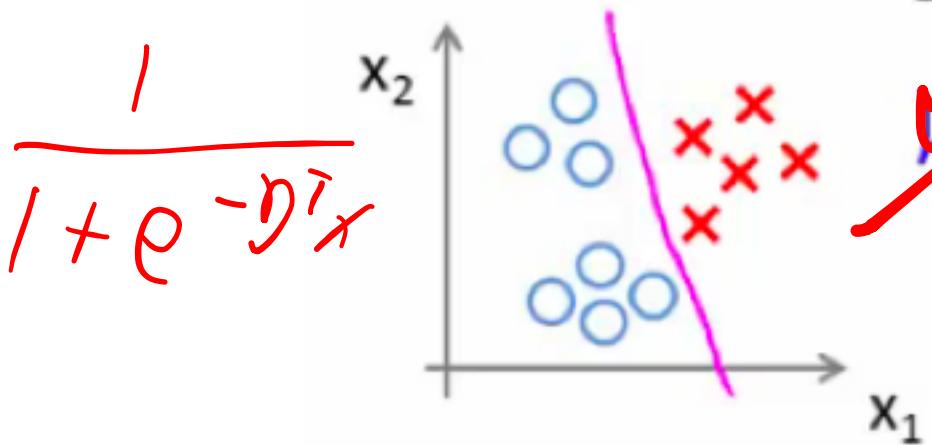


$$\frac{e^{\omega^T x}}{e^{\omega^T x} + e^{-\theta^T x} + e^{-\eta^T x}}$$

softmax



$$\frac{e^{\theta^T x}}{e^{\omega^T x} + e^{-\theta^T x} + e^{-\eta^T x}}$$



$$\frac{e^{-\gamma^T x}}{e^{\omega^T x} + e^{-\theta^T x} + e^{-\eta^T x}}$$

