



# Applied Machine Learning

**BITS Pilani**  
Pilani Campus

Dr. Harikrishnan N B  
Computer Science and Information Systems



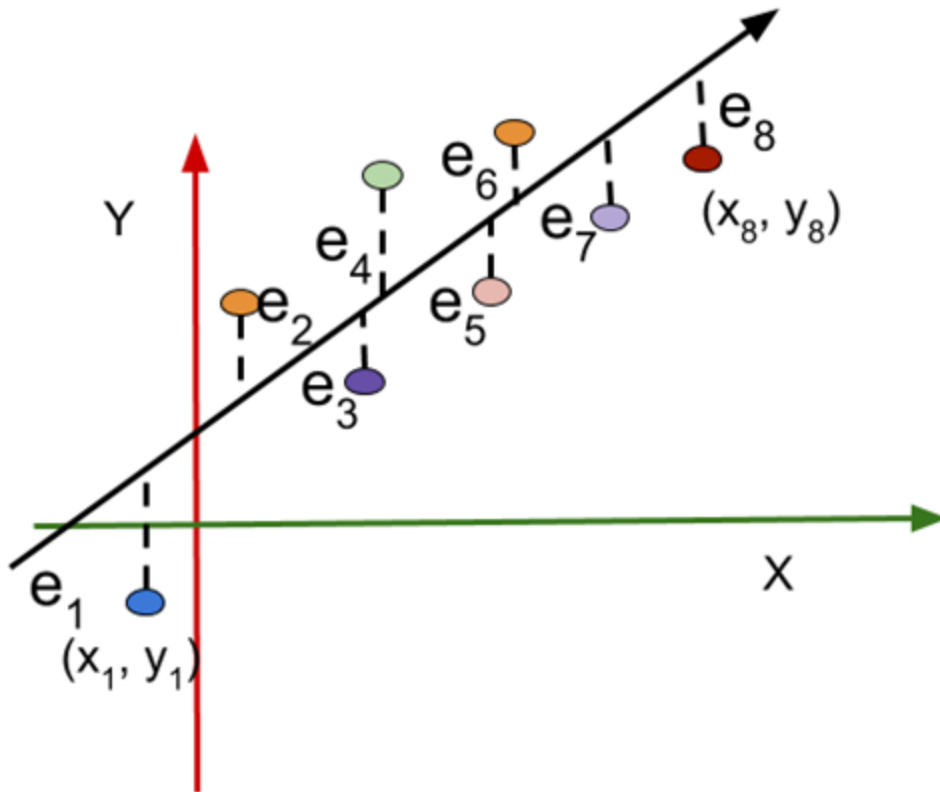
# **SE ZG568 / SS ZG568, Applied Machine Learning Lecture No. 11 [13 April 2025]**

---

# Linear Regression using Optimization

## Logistic Regression

# Linear Least Square Regression



$$y_1 = mx_1 + c + e_1$$

$$y_2 = mx_2 + c + e_2$$

$$y_3 = mx_3 + c + e_3$$

$$y_4 = mx_4 + c + e_4$$

$$y_5 = mx_5 + c + e_5$$

$$y_6 = mx_6 + c + e_6$$

$$y_7 = mx_7 + c + e_7$$

$$y_8 = mx_8 + c + e_8$$

# Quick Recap



Linear regression models the relationship between the input  $x$  and the output  $y$  as a linear combination of the features:

$$\underline{h_\theta(x)} = \underline{\theta_0} + \theta_1 \overset{!}{x_1} + \theta_2 \overset{!}{x_2} + \dots + \theta_n \overset{!}{x_n} \rightarrow \text{feature}$$

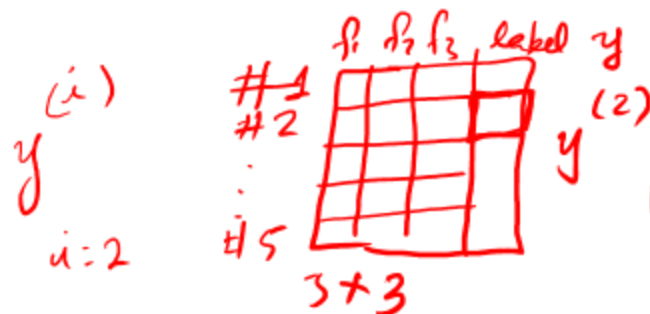
where:

- $x = [1, x_1, x_2, \dots, x_n]^T$  (including  $x_0 = 1$  for bias)
- $\theta = [\theta_0, \theta_1, \dots, \theta_n]^T$  are the parameters (weights)
- $h_\theta(x)$  is the predicted output.

In vector form:

$$h_\theta(x) = \theta^T x$$

To measure how well our model predicts, we use the **Mean Squared Error (MSE)**:

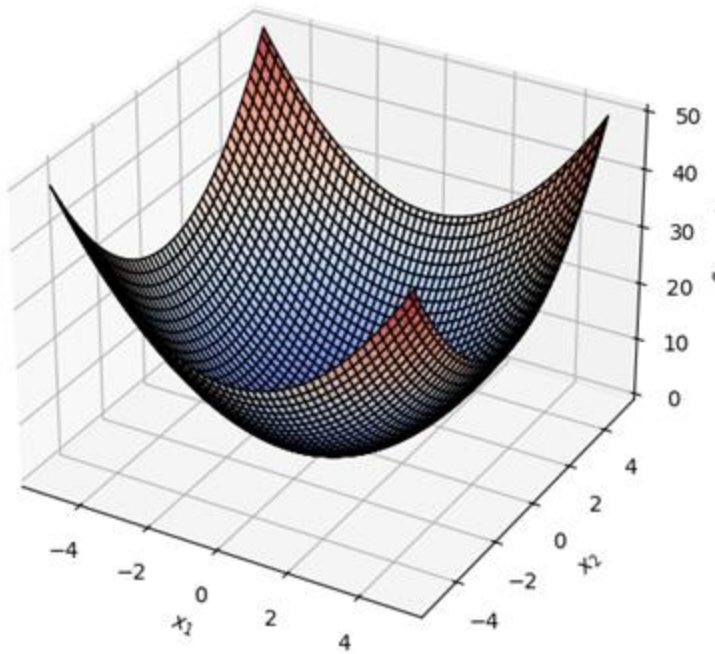


$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - \underline{y^{(i)}})^2$$

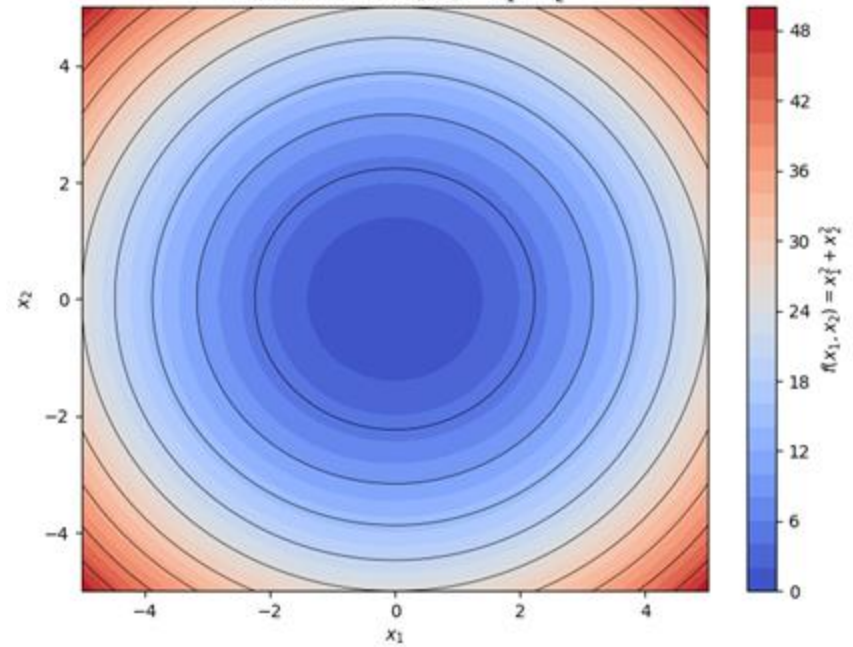
where:

- $m$  is the number of training examples
- $y^{(i)}$  is the actual output for the  $i$ -th training example
- $h_{\theta}(x^{(i)})$  is the predicted output.
- The factor  $\frac{1}{2}$  is used for mathematical convenience when differentiating.

3D Plot of  $f(x_1, x_2) = x_1^2 + x_2^2$



Contour Plot of  $f(x_1, x_2) = x_1^2 + x_2^2$



Gradient descent is used to minimize  $J(\theta)$ . The update rule is:

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

where:

- $\alpha$  is the **learning rate**, controlling step size
- $\frac{\partial J(\theta)}{\partial \theta_j}$  is the gradient (partial derivative)

Computing the gradient:

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

Thus, the update step becomes:

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

This process repeats until convergence (i.e., when  $J(\theta)$  stops changing significantly).



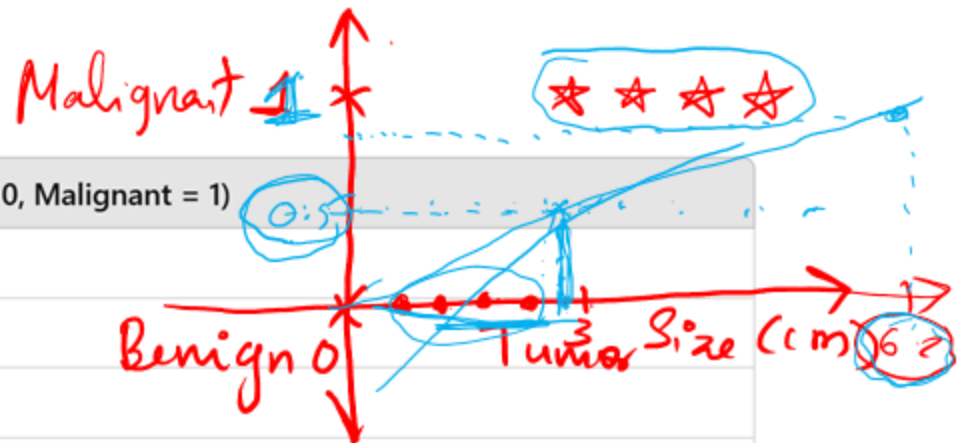
# Logistic Regression

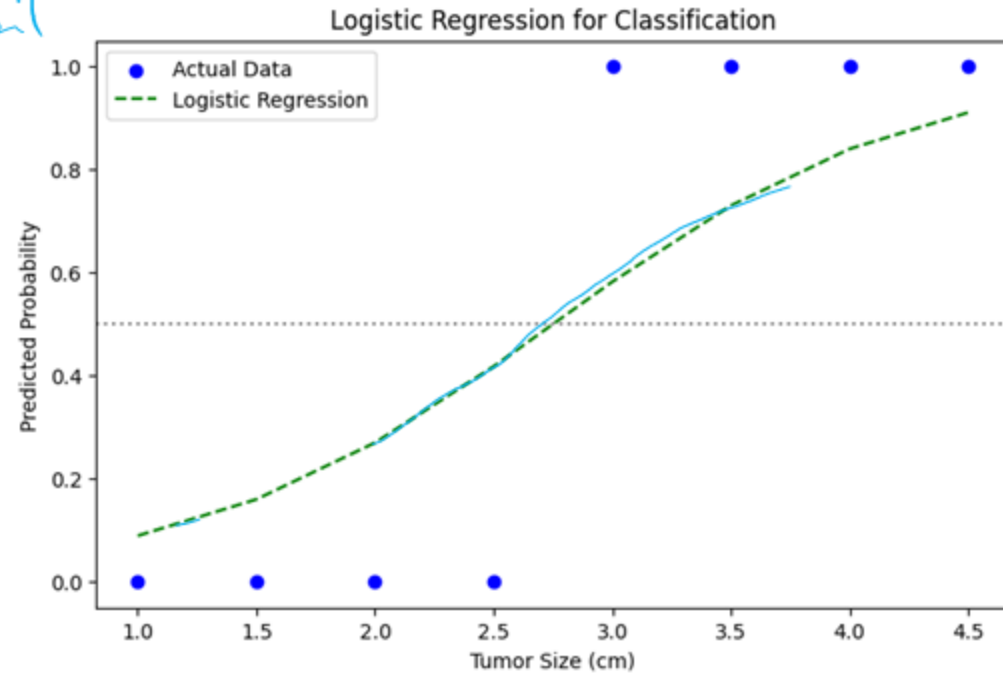
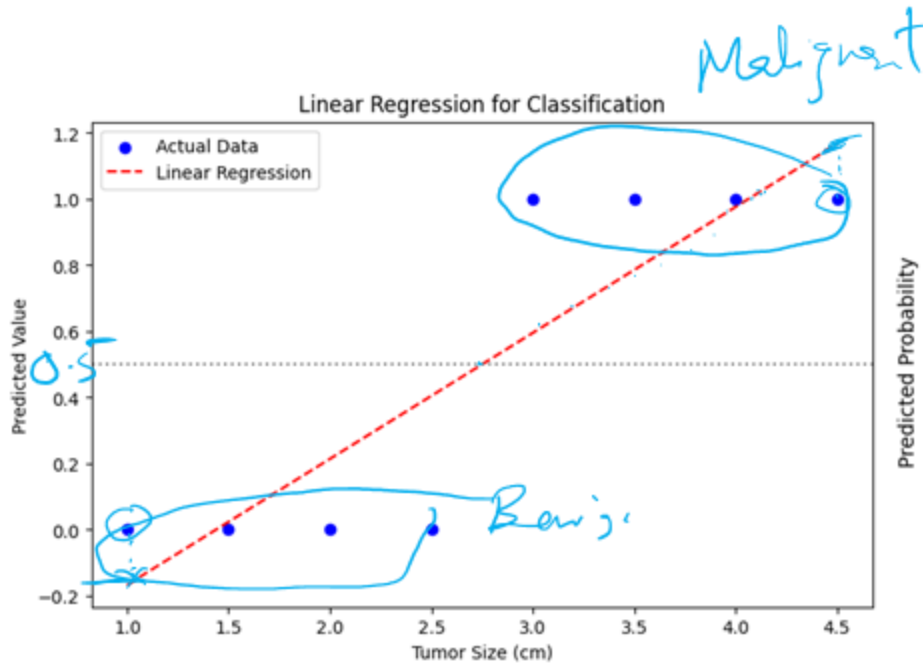
## Scenario

Doctors want to classify tumors as **benign (0)** or **malignant (1)** based on tumor size. Larger tumors are more likely to be malignant, but there's no hard cutoff—so we use logistic regression to model this probability.

## Hypothetical Data

| Tumor Size (cm) | Diagnosis (Benign = 0, Malignant = 1) |
|-----------------|---------------------------------------|
| 1.0 →           | 0 (Benign)                            |
| 1.5             | 0                                     |
| 2.0             | 0                                     |
| 2.5             | 0                                     |
| 3.0             | 1                                     |
| 3.5             | 1                                     |
| 4.0             | 1                                     |
| 4.5             | 1                                     |





# Effect of Outlier



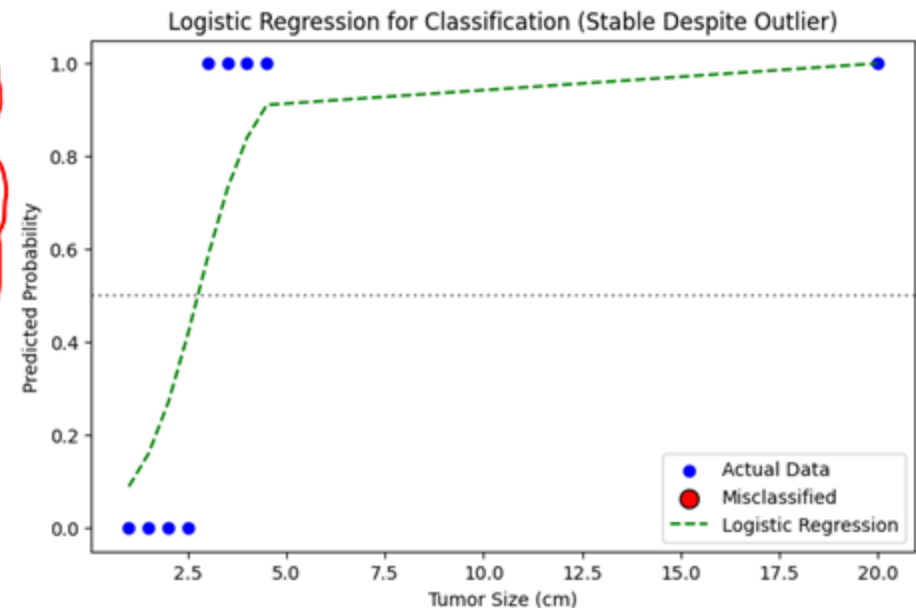
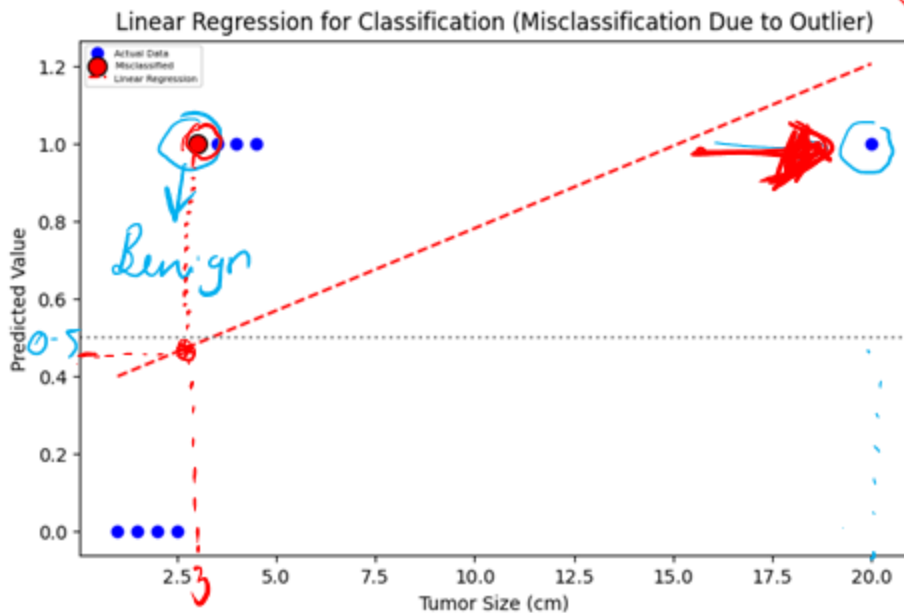
| Tumor Size (cm) | Diagnosis (0 = Benign, 1 = Malignant) |
|-----------------|---------------------------------------|
| 1.0             | 0                                     |
| 1.5             | 0                                     |
| 2.0             | 0                                     |
| 2.5             | 0                                     |
| 3.0             | 1                                     |
| 3.5             | 1                                     |
| 4.0             | 1                                     |
| 4.5             | 1                                     |
| 20.0            | 1 (Outlier)                           |

# Effect of Outlier



Regularized

$$\min_x \left( \|Ax - b\|_2^2 + \lambda \|Dx\|_2^2 \right)$$



Threshold = 0.5,  
Accuracy = 0.89

## 1. Effect on Linear Regression

- Linear regression tries to fit a straight line, minimizing squared errors.
- A large outlier (e.g., tumor size = 10 cm, malignant = 1) will pull the line upwards, making it steeper.
- This shifts the decision threshold (where the prediction crosses 0.5), possibly classifying smaller tumors incorrectly.

## 2. Effect on Logistic Regression

- Logistic regression models probabilities using the sigmoid function, which is bounded between 0 and 1.
- A large outlier won't distort the S-curve as much because it saturates towards 1.
- The decision boundary (where probability = 0.5) remains relatively stable.

# Logistic Regression



$$x = [x^{(1)} \quad x^{(2)}]$$

$$\theta^T x^{(1)} = \theta_0 + \theta_1 x_1^{(1)} + \theta_2 x_2^{(1)}$$
 Goal: learn  $\theta_0, \theta_1, \theta_2$

The hypothesis function in logistic regression is given by:

$$h_{\theta}(x) = g(\theta^T x)$$

function of  $\theta^T x$

$$h_{\theta}(x) = g(\theta^T x)$$

and hypothesis function of linear regression  

$$h_{\theta}(x) = w x + b$$

$$\begin{bmatrix} w & b \end{bmatrix} \begin{bmatrix} x \\ 1 \end{bmatrix}$$

$$\theta^T \begin{bmatrix} x \\ 1 \end{bmatrix}$$

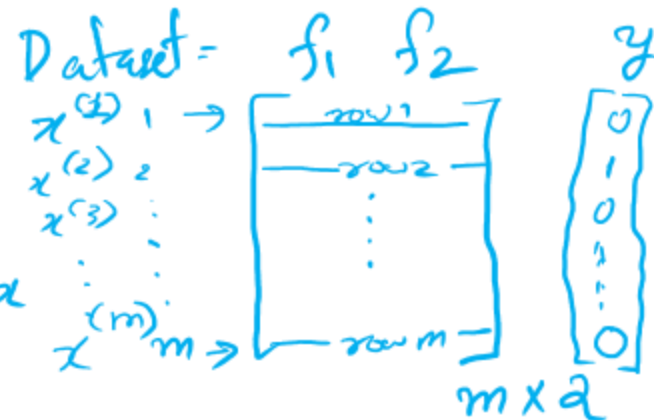
where:

- $\theta$  is the parameter vector (weights).
- $x$  is the feature vector (input data).
- $\theta^T x$  represents the linear combination of features and parameters.
- $g(z)$  is the sigmoid function, defined as:

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

Thus, expanding  $h_{\theta}(x)$ :

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



This function maps any real number to the range (0,1), making it suitable for probability estimation in classification problems.

# Notation

Data
 
$$\begin{matrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \\ \vdots \\ x^{(m)} \end{matrix} \begin{bmatrix} f_1 & f_2 \\ x_1^{(1)} & x_2^{(1)} \\ \vdots & \vdots \\ x_1^{(m)} & x_2^{(m)} \end{bmatrix}$$

$$\theta^T x^{(i)} \Rightarrow \theta_0 + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$

where  $i \in \{1, 2, \dots, m\}$

Goal: learn  $\theta_0, \theta_1, \theta_2$

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

Is  $\theta^T x^{(i)}$  is scalar or vector?

$$x=2 \quad z=2$$

$$k = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

$$x^{(i)} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

$$\theta^T x^{(i)} = \begin{bmatrix} \theta_0 & \theta_1 & \theta_2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \theta_0 + \theta_1 2 + \theta_2 3$$

scalar



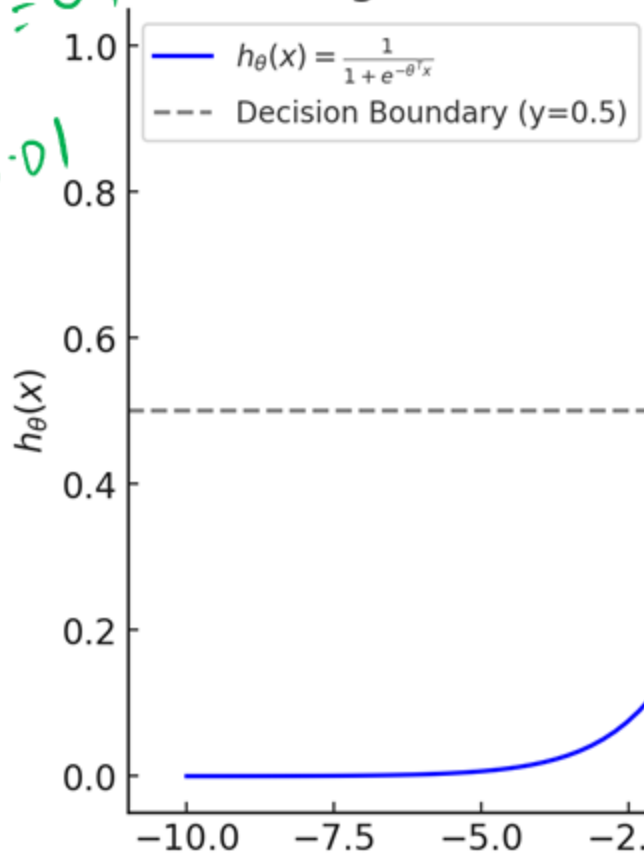
$\frac{1}{1 + e^{-\theta^T x}}$  vs  $\theta^T x$   
y-axis vs x-axis

(are I what happens when

$\theta^T x \rightarrow \infty$

### Sigmoid Function for Logistic Regression

$\frac{1}{b} = 0.1$   
 $\frac{1}{100} = 0.01$   
 $\frac{1}{1000}$



$1 \left( \frac{1}{1 + \frac{1}{e^{\infty}}} \right)$

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

Case II  
what happens  
 $\theta^T x \rightarrow -\infty$

Case III  
what happens  
when  $\theta^T x = 0$

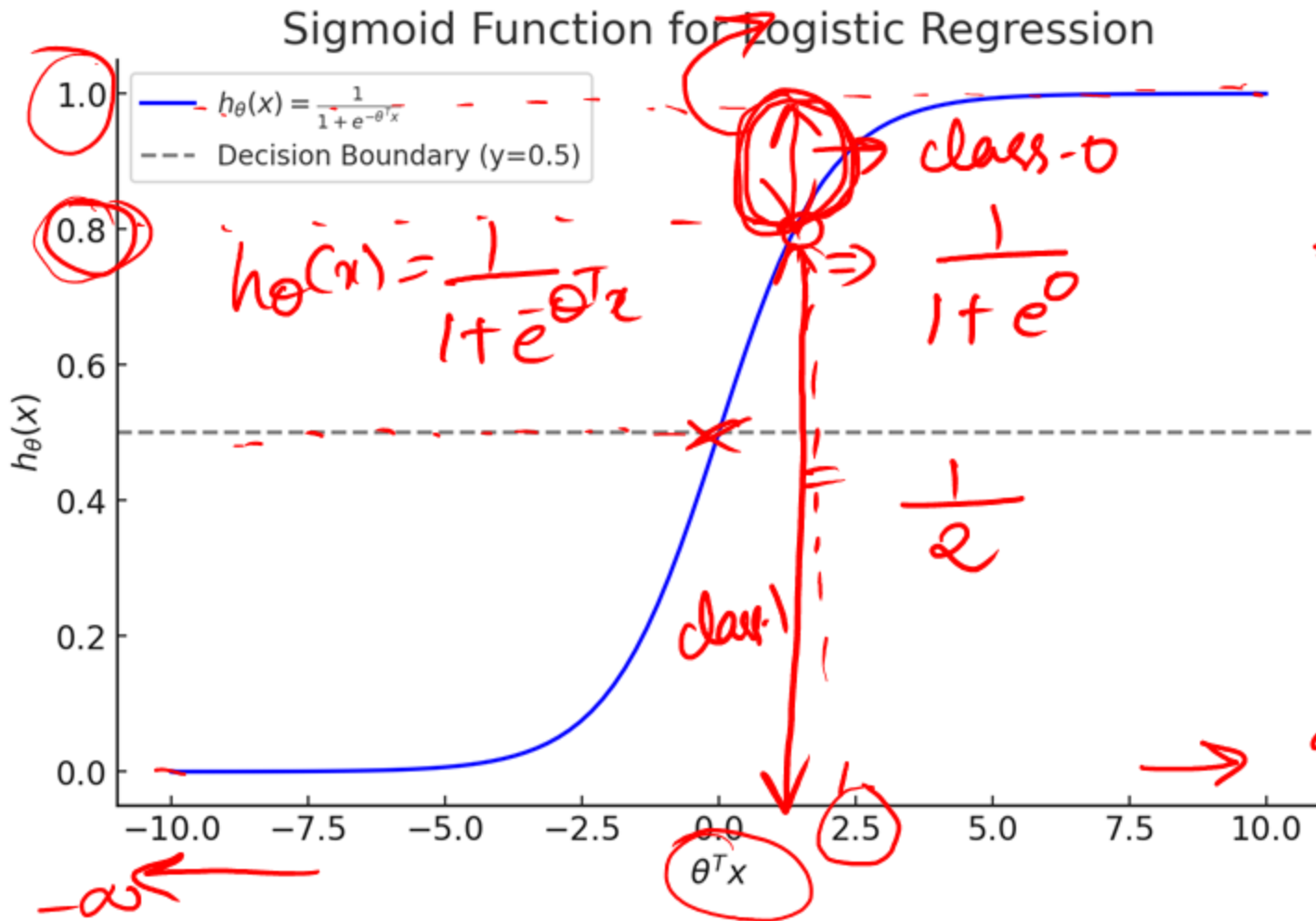


Case II  
 $\Theta^T x = 0$

Case II  $\theta^T x \rightarrow -\infty$



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



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

$$= \frac{1}{1 + e^w}$$

$$= \frac{1}{1+\infty} = 0$$

## Mathematical Formulation of Logistic Regression

We have already defined the hypothesis function as:

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

to Probabilities

Now, let's go through the remaining key mathematical components of logistic regression.

# Binary Classification Problem class labels are 0 or 1



$$P(y=y|\theta;x) = h_{\theta}(x)^y (1-h_{\theta}(x))^{1-y}$$

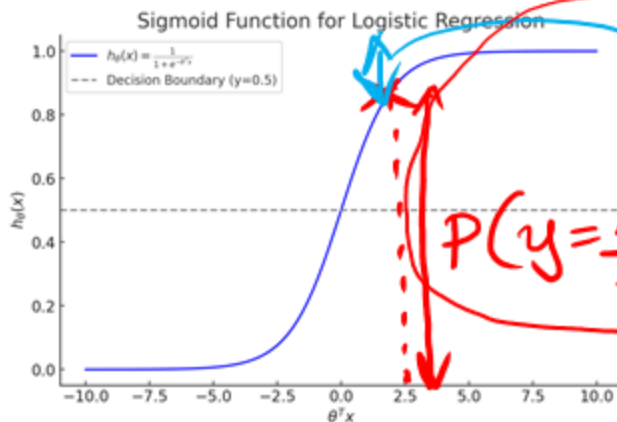
## 1. Probability Interpretation

$$P(y=1|\theta;x) = h_{\theta}(x)^1 (1-h_{\theta}(x))^0 \Rightarrow h_{\theta}(x)$$

Since  $h_{\theta}(x)$  represents the probability of  $y = 1$ , we can write:

$$P(y=1|x;\theta) = h_{\theta}(x)$$

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



$$P(y=0|\theta^T x) = 1 - h_{\theta}(x)$$

$$P(y=1|\theta^T x) = h_{\theta}(x)$$

$$P(y=0|\theta;x) = h_{\theta}(x)^0 (1-h_{\theta}(x))^{1-0} \Rightarrow 1 - h_{\theta}(x)$$

# Likelihood Function

independent

$$P(A \cap B \cap C) = P(A) \cdot P(B) \cdot P(C)$$

Prob of getting Head (P),  $P(T) = 1 - p$

innovate

achieve

lead

H T T H H T H T T H

$$P \cdot (1-p) \cdot (1-p) \cdot P \cdot P \cdot P \cdot P \cdot (1-p) \cdot P$$

$$= P^5 \cdot (1-p)^5$$

## 1. Likelihood Function

Given a dataset  $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$ , we assume that each training example follows a Bernoulli distribution:

$$P(y|x; \theta) = h_{\theta}(x)^y (1 - h_{\theta}(x))^{(1-y)}$$

For all  $m$  training samples, assuming independence, the likelihood function is the product of the probabilities for all data points:

$\theta_0, \theta_1, \theta_2$

$$L(\theta) = \prod_{i=1}^m h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{(1-y^{(i)})}$$

product

where:

- $h_{\theta}(x^{(i)}) = \frac{1}{1 + e^{-\theta^T x^{(i)}}}$  is the predicted probability of  $y = 1$ .
- If  $y^{(i)} = 1$ , the term  $(1 - h_{\theta}(x^{(i)}))^{(1-y^{(i)})}$  vanishes.
- If  $y^{(i)} = 0$ , the term  $h_{\theta}(x^{(i)})^{y^{(i)}}$  vanishes.

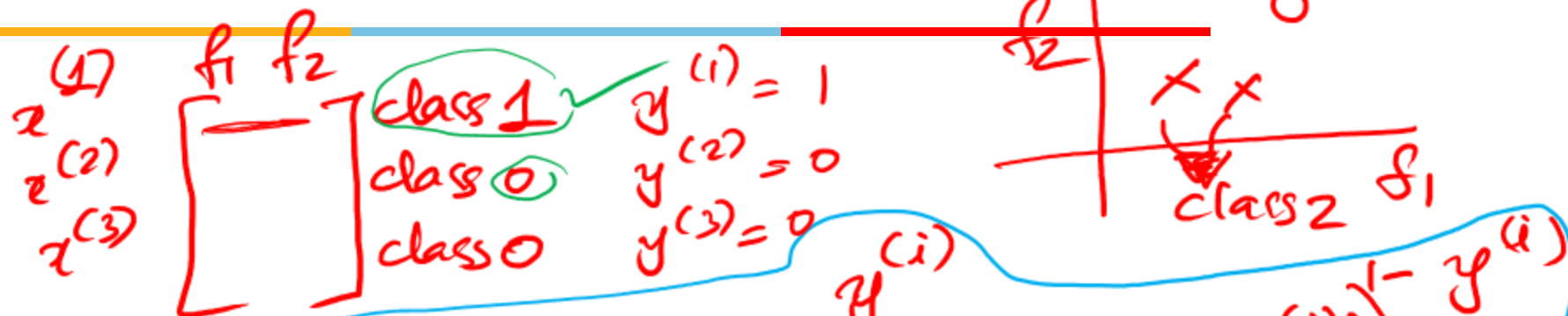
$y^{(i)} = 1$   
 $h_{\theta}(x^{(i)}) = 0.01$

Data  $x_1, x_2$

$x^{(1)}$   
 $x^{(2)}$   
 $x^{(3)}$   
 $\vdots$

Training Data

Take some random  $\theta_0, \theta_1, \theta_2$

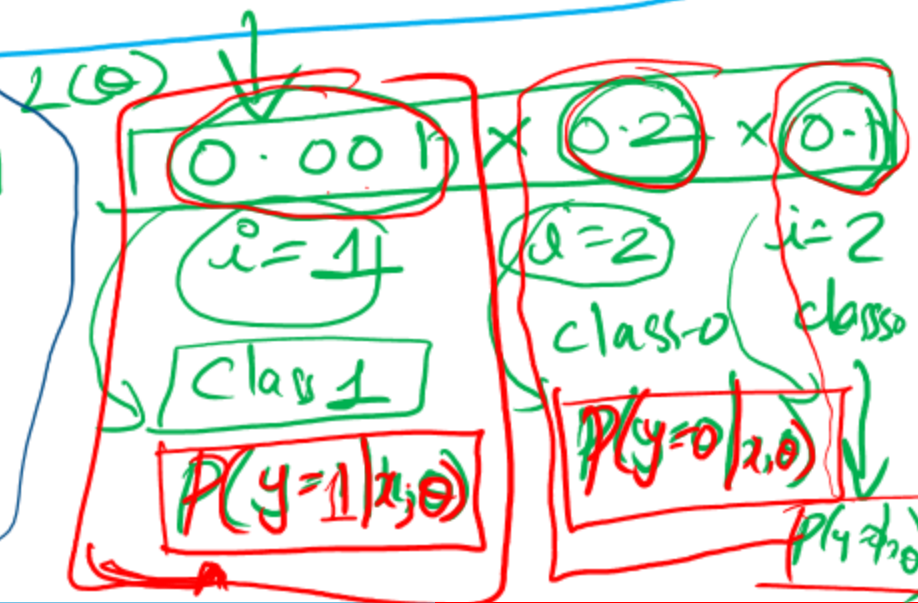


$$L(\theta) = \prod_{i=1}^m h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1 - y^{(i)}}$$

$$h_{\theta}(x^{(1)}) = 0.001$$

$$h_{\theta}(x^{(2)}) = 0.8$$

$$h_{\theta}(x^{(3)}) = 0.9$$



# Maximizing Likelihood

## Maximizing the Likelihood Function in Logistic Regression

Once we have the likelihood function:

$$L(\theta) = \prod_{i=1}^m h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{(1-y^{(i)})}$$

our goal is to find the parameters  $\theta$  that maximize this function. In other words, we want to find the best  $\theta$  such that the predicted probabilities align as closely as possible with the actual labels in the training data.

New  $\theta$   $\rightarrow$

$$h_{\theta}(x^{(1)}) = 0.9$$

$$h_{\theta}(x^{(2)}) = 0.2$$

$$h_{\theta}(x^{(3)}) = 0.1$$

$$L(\theta) = 0.9 \times 0.8 \times 0.9$$

$\downarrow$   $\downarrow$   $\downarrow$   
 $P(y=1|x^{(1)}; \theta)$   $P(y=0|x^{(2)}; \theta)$   $P(y=0|x^{(3)}; \theta)$

## Why Maximize the Likelihood?

- The likelihood function represents the probability of observing the given dataset, assuming the logistic regression model is correct.
- A higher likelihood means that the model's predicted probabilities closely match the actual outcomes.
- By maximizing  $L(\theta)$ , we are choosing parameters  $\theta$  that make the observed data most probable under our model.





$\max(l(\theta))$   
 $0.1 \times 0.2 \times 0.01$   
 $\min(-l(\theta))$



$\log(a \cdot b) = \log(a) + \log(b)$   
 $\log a^y = y \log a$

## Step 1: Taking the Log of the Likelihood

Given the likelihood function:

$$L(\theta) = \prod_{i=1}^m h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{(1-y^{(i)})}$$

we take the log-likelihood:

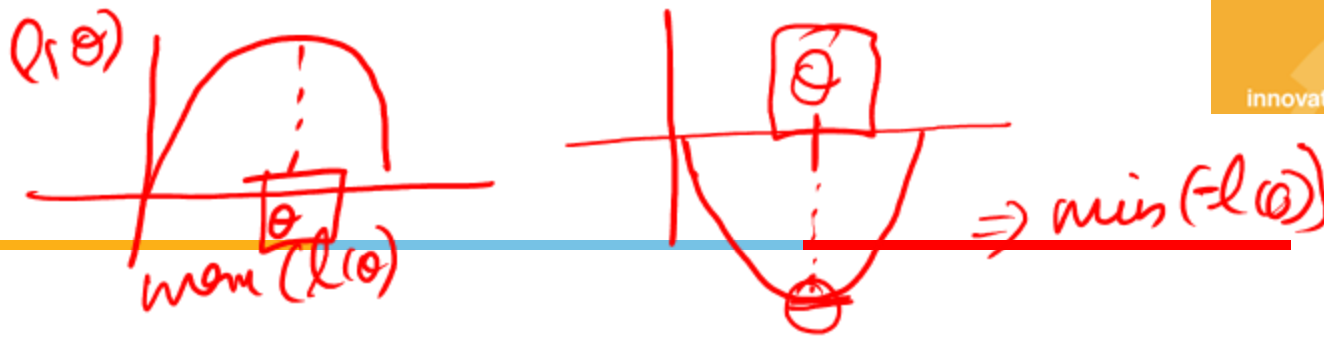
$$\ell(\theta) = \log L(\theta) = \sum_{i=1}^m \left[ y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

To ensure consistency with loss functions commonly used in machine learning, we take the average log-likelihood (by including  $\frac{1}{m}$ ):

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

$\max \ell(\theta)$  is same as  $\min -\ell(\theta)$





### Step 3: Converting to a Minimization Problem

Optimization algorithms typically **minimize** a function rather than maximize it. Instead of maximizing  $\ell(\theta)$ , we minimize its **negative**:

$$J(\theta) = -\ell(\theta)$$

Substituting the log-likelihood:

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

This is the **binary cross-entropy loss** (also called **log-loss**), which is the function we minimize using **gradient descent**.

iterations

$$\theta^{t+1} := \theta^t - \alpha \frac{\partial J(\theta)}{\partial \theta}$$

$$\frac{\partial J(\theta)}{\partial \theta} =$$

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







$$x^{(i)} = \begin{bmatrix} x_1^{(i)} & x_2^{(i)} \end{bmatrix} \quad \theta = \theta_0, \theta_1, \theta_2$$

$$\begin{bmatrix} 1 & x_1^{(i)} & x_2^{(i)} \end{bmatrix}$$

## Gradient Descent for Parameter Estimation

The gradient of the log-loss function with respect to  $\theta_j$  is:

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

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

where:

- $m$  is the number of training examples.
- $h_{\theta}(x^{(i)}) = \frac{1}{1+e^{-\theta^T x^{(i)}}}$  is the predicted probability.
- $y^{(i)}$  is the actual class label (0 or 1).
- $x_j^{(i)}$  is the  $j$ -th feature of the  $i$ -th training example.

$$\frac{\partial J}{\partial \theta_1} = \dots$$

We update  $\theta_j$  using **gradient descent**:

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

$$\frac{\partial J}{\partial \theta_2} = \dots$$

where:

- $\alpha$  is the **learning rate**, controlling how big the update steps are.

# Coding



$$y = m \odot x + b$$

$$\theta^T x$$

$$\begin{bmatrix} \theta_0 & \theta_1 \\ b & m \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix}$$

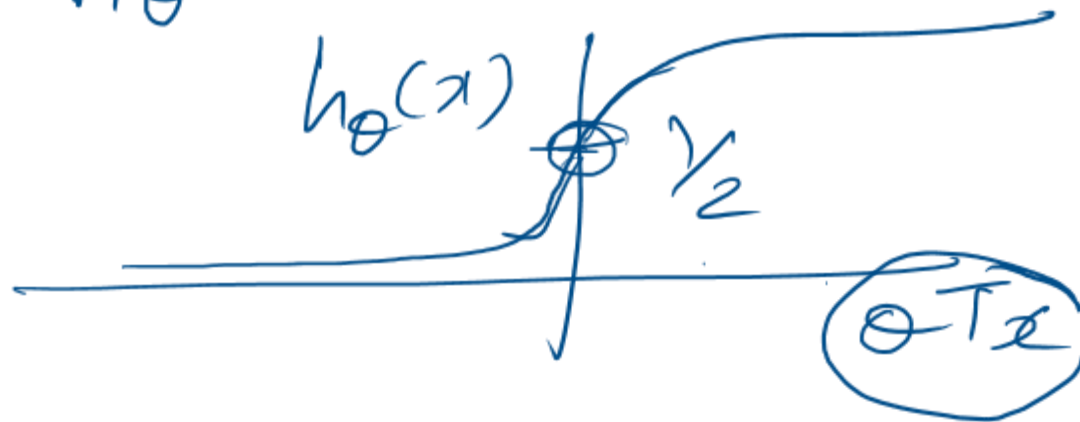
$$m_1 x_1 + m_2 x_2 + \textcircled{b_2}$$

$$\begin{bmatrix} b_2 & m_1 & m_2 \\ \theta^T \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

# Decision Rule

$$h_{\theta}(x_i) \geq 0.5 \text{ class-1}$$

$$h_{\theta}(x) < 0.5 \text{ class-0}$$





$$\theta^T x$$

$$y = mx + b$$



$$mx - y + b$$

$$\theta \begin{pmatrix} b \\ m \\ -1 \end{pmatrix} \begin{pmatrix} 1 \\ x \\ y \end{pmatrix}$$

$$\theta \quad \theta^T x = b + mx - y$$

.

$$p(y=1|x, \theta) = h_{\theta}(x)$$

$$p(y=0|x, \theta) = 1 - h_{\theta}(x)$$

For a test  
data

if  $p(y=1|x; \theta) \geq p(y=0|x; \theta)$

$$p(y=1|x; \theta) \geq 1$$

$$\Rightarrow \left[ \frac{h_{\theta}(x)}{1 - h_{\theta}(x)} \geq 1 \right]$$

$$\frac{\frac{1}{1+e^{-\theta^T x}}}{1 - \frac{1}{1+e^{-\theta^T x}}} \geq 1$$

$$\Rightarrow \frac{1}{1+e^{-\theta^T x}} \frac{1+e^{-\theta^T x}}{e^{-\theta^T x}} \geq 1$$

$$e^{\theta^T x} \geq 1$$

$$\ln e^{\theta^T x} \geq \ln(1)$$

$$\boxed{\theta^T x \geq 0}$$

$$\frac{\frac{1}{1+e^{-\theta^T x}}}{\cancel{1+e^{-\theta^T x}} - 1}$$

$$\frac{1}{1+e^{-\theta^T x}}$$























---

# Perceptron

# Logistic Regression

# SVM

# Perceptron

innovate

achieve

lead

















