

Applied Machine Learning

BITS Pilani
Pilani Campus

Dr. Harikrishnan N B
Computer Science and Information Systems

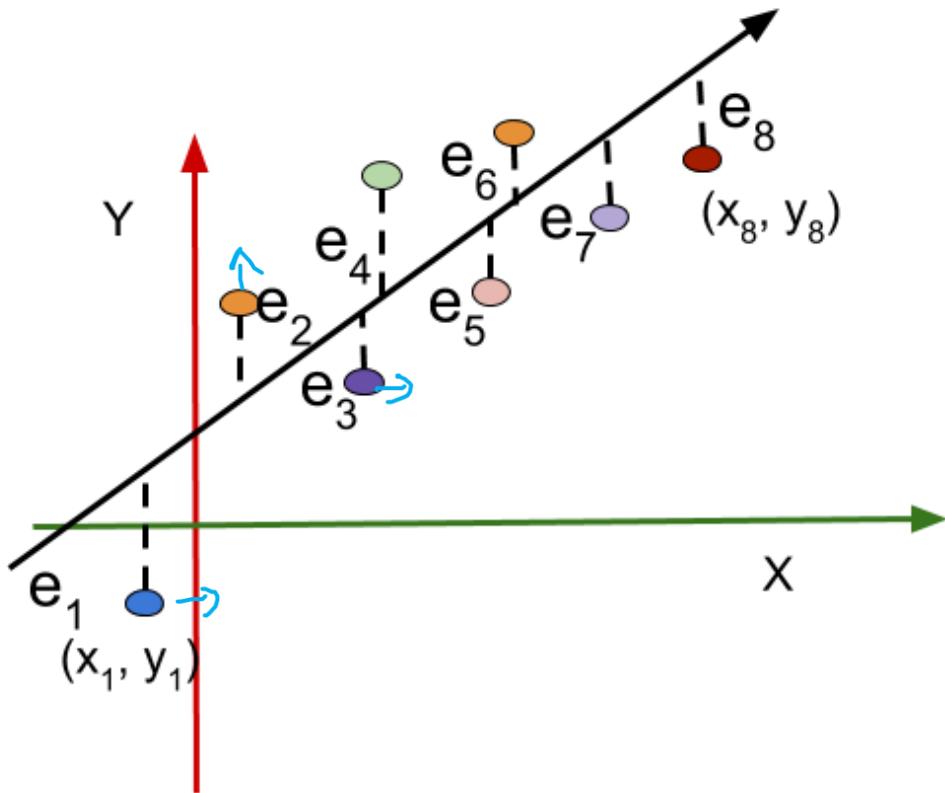


SE ZG568 / SS ZG568, Applied Machine Learning Lecture No. 10

Linear Regression using Optimization

Logistic Regression

Linear Least Square Regression



$$\begin{aligned}y_1 &= mx_1 + c + e_1 \\y_2 &= mx_2 + c + e_2 \\y_3 &= mx_3 + c + e_3 \\y_4 &= mx_4 + c + e_4 \\y_5 &= mx_5 + c + e_5 \\y_6 &= mx_6 + c + e_6 \\y_7 &= mx_7 + c + e_7 \\y_8 &= mx_8 + c + e_8\end{aligned}$$

SG x_1 x_2 no of bed rooms y Price of price

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

$$y = m x + c$$
 ~~$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2$~~
 ~~θ_0~~
 ~~θ_1~~
 ~~θ_2~~

$$y = m x + c$$
 ~~$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2$~~
 ~~θ_0~~
 ~~θ_1~~
 ~~θ_2~~

$$y =$$

$$\boxed{\theta_0 + \theta_1 x}$$

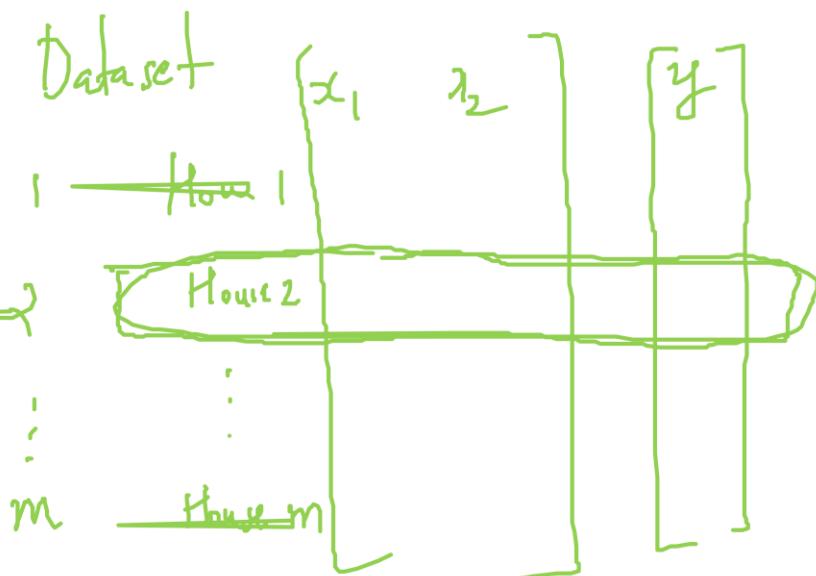
$$h_{\theta}(x) = \theta_0 \cdot 1 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

error =

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

Hour
Second
Z₂
Z₁

Training data = m

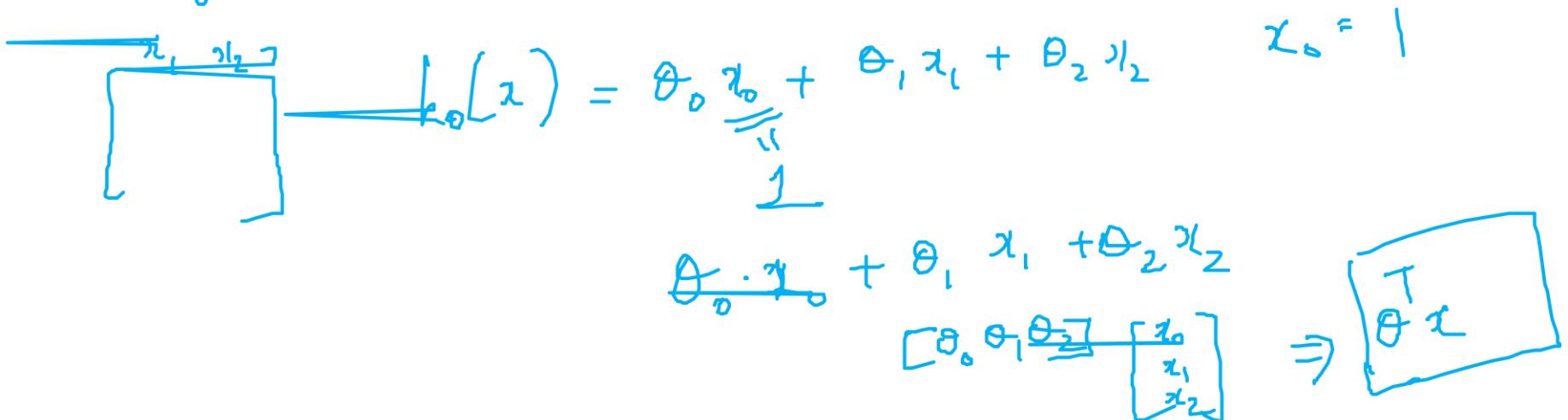


$x^{(i)}$ } ith training data
 $x^{(i)}$ = [$x_1^{(i)}$ $x_2^{(i)}$]
 feature

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

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

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$



$$\text{Error} = \frac{1}{2m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$$

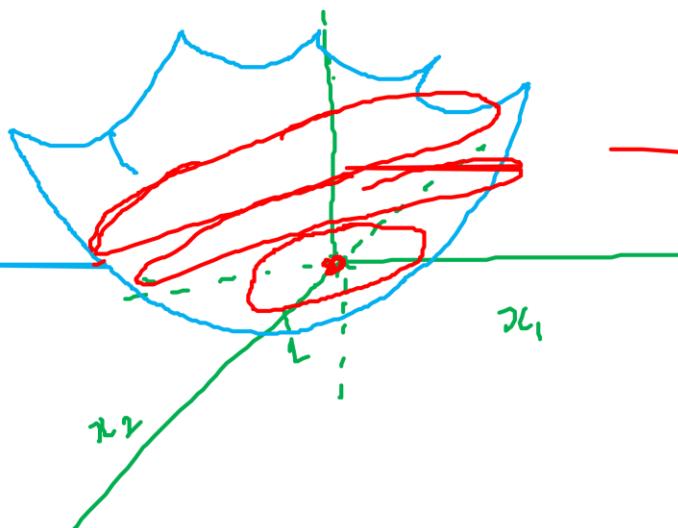
Gradient Descent

A [Puncher]

$$f(x_1, x_2) = x_1^2 + x_2^2$$

$x_1=0, x_2=1$

$f(x_1, x_2) = 1$

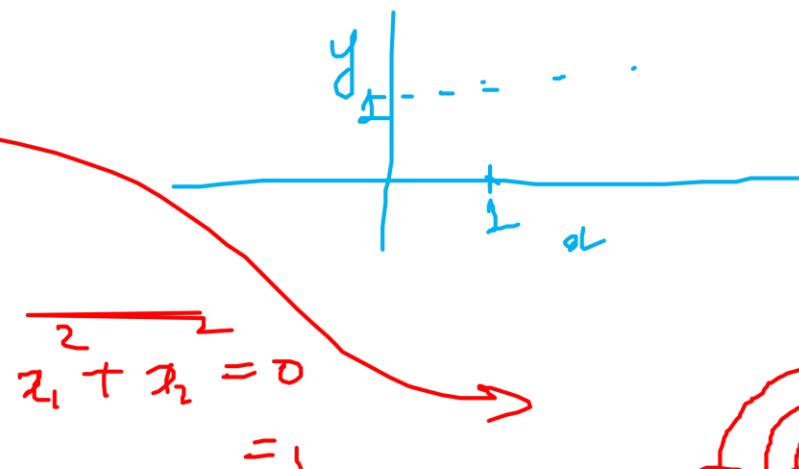


B [Level set]

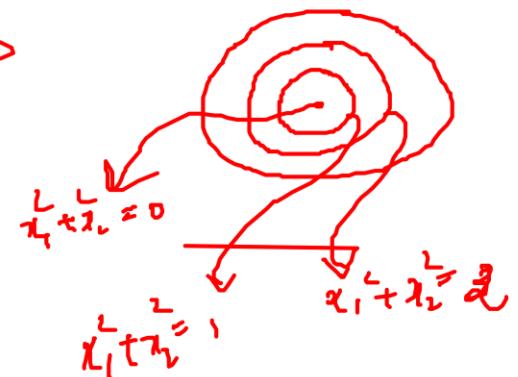
$$x_1^2 + x_2^2 = 5$$

$$f(x) = \max_{y \in \mathbb{R}} L$$

$$m = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$



$$x_1^2 + x_2^2 = 5$$





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

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

where: \underline{x}

- $x = [1, x_1, x_2, \dots, x_n]^T$ (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:

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

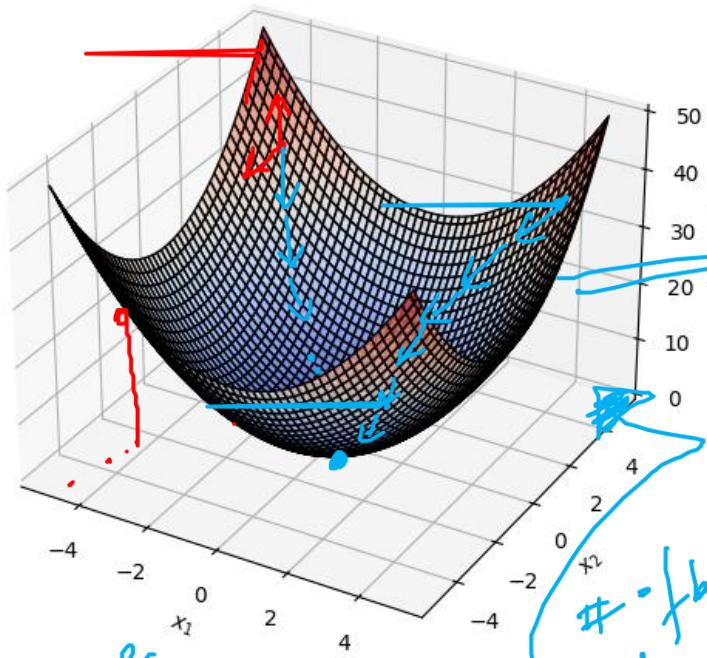


θ_{new}

$$\theta_{\text{new}} = \theta_{\text{old}} - \frac{\partial \text{Error}}{\partial \theta}$$


 $\theta = \text{random initialization}$

$$\theta - \alpha \frac{\partial \text{Error}}{\partial \theta}$$

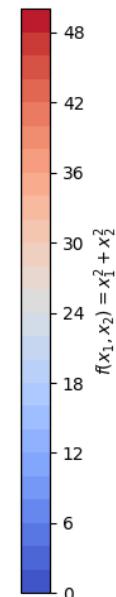
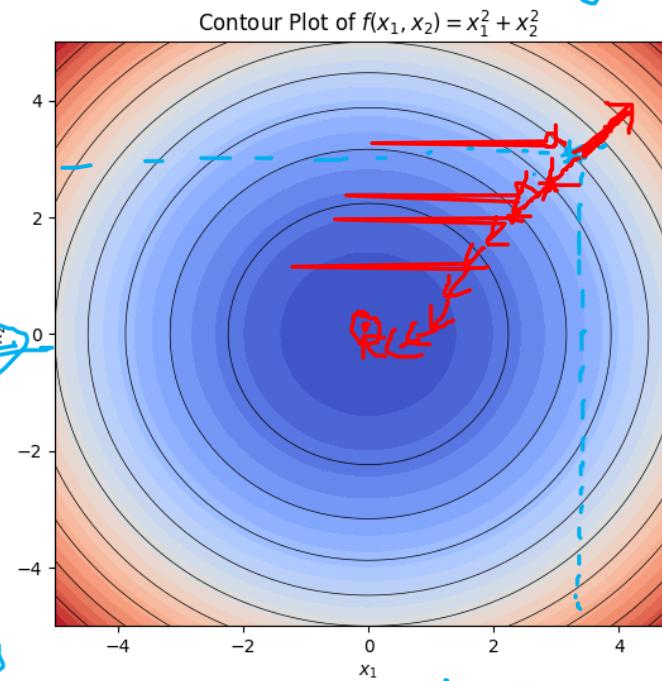
 3D Plot of $f(x_1, x_2) = x_1^2 + x_2^2$


$$\frac{\# \text{ of features}}{2m} \sum_{i=1}^m$$

$$(h_{\theta}(x^{(i)}) - y^{(i)})^2$$

 ~~θ_{new}~~
 $\alpha : \text{learning rate}$
~~Level off~~

$$\theta_{\text{new}} = \theta_{\text{old}} - \alpha \frac{\partial \text{Error}}{\partial \theta_{\text{old}}}$$



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

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

=====

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.

$$\theta_{\text{new}} := \theta_{\text{old}} - \frac{\partial J(\theta)}{\partial \theta}$$

$$\theta_{\text{new}} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}_{n+1} \quad \vec{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\frac{\partial J(\theta)}{\partial \theta_0}, \quad \frac{\partial J(\theta)}{\partial \theta_1}, \quad \frac{\partial J(\theta)}{\partial \theta_2}, \quad \dots, \quad \frac{\partial J(\theta)}{\partial \theta_n}$$

$m=3$, \neq for $n=2$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^3 \left(\theta_0 x_0^{(i)} + \cancel{\theta_1 x_1^{(i)}} + \cancel{\theta_2 x_2^{(i)}} - y^{(i)} \right)^2$$

$$\cancel{\theta_1 x_1^{(i)}} \cancel{+ \theta_2 x_2^{(i)}} \\ f_{\theta(i)} = \frac{1}{2} (y^{(i)} - \hat{y}^{(i)})^2$$

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{2m} \left[\begin{aligned} & \left(\theta_0 x_0^{(1)} + \cancel{\theta_1 x_1^{(1)}} + \theta_2 x_2^{(1)} - y^{(1)} \right) \circledcirc \\ & \left(\theta_0 x_0^{(2)} + \theta_1 \cancel{x_1^{(2)}} + \theta_2 x_2^{(2)} - y^{(2)} \right) \circledcirc \\ & \left(\theta_0 x_0^{(3)} + \cancel{\theta_1 x_1^{(3)}} + \theta_2 x_2^{(3)} - y^{(3)} \right) \circledcirc \end{aligned} \right]$$

$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{2m} \left[\begin{aligned} & 2 \left(\theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} + \cancel{\theta_2 x_2^{(1)}} - y^{(1)} \right) \times x_0^{(1)} + \\ & 2 \left(\theta_0 x_0^{(2)} + \theta_1 x_1^{(2)} + \cancel{\theta_2 x_2^{(2)}} - y^{(2)} \right) \times x_0^{(2)} + \\ & 2 \left(\cancel{\theta_0 x_0^{(3)}} + \theta_1 x_1^{(3)} + \theta_2 x_2^{(3)} - y^{(3)} \right) \times x_0^{(3)} \end{aligned} \right]$$

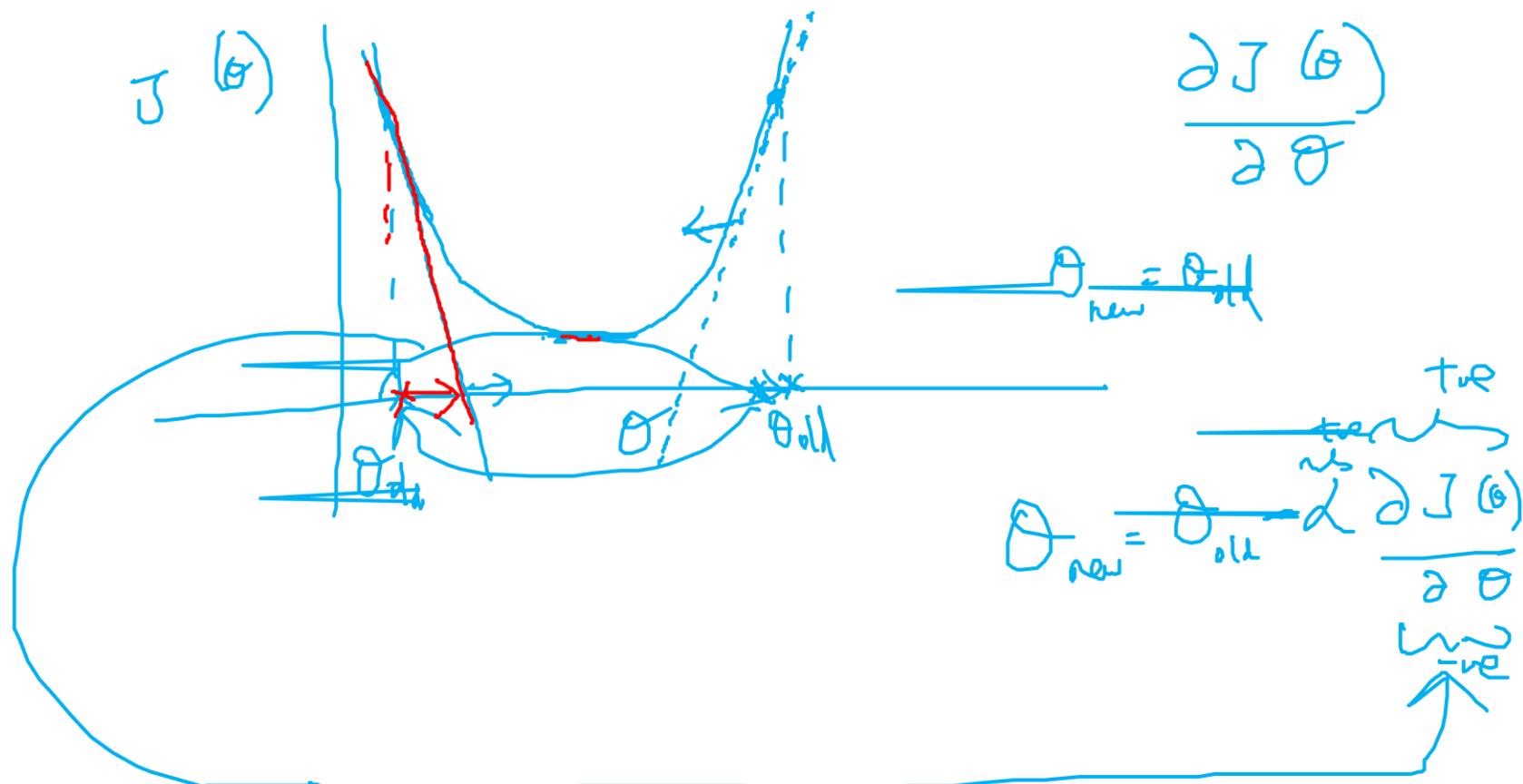
$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m \left(h_\theta(x^{(i)}) - y^{(i)} \right) \times x_0^{(i)}$$

$$\frac{\partial J(\theta)}{\partial \theta} \Rightarrow \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_0} \\ \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{bmatrix} \quad \begin{array}{l} \text{scalar} \\ \text{scalar} \\ \vdots \\ \text{vector} \end{array}$$



$$\theta_{\text{new}} = \theta_{\text{old}} - \frac{\partial J(\theta)}{\partial \theta}$$

-ve



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:

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

Computing the gradient:



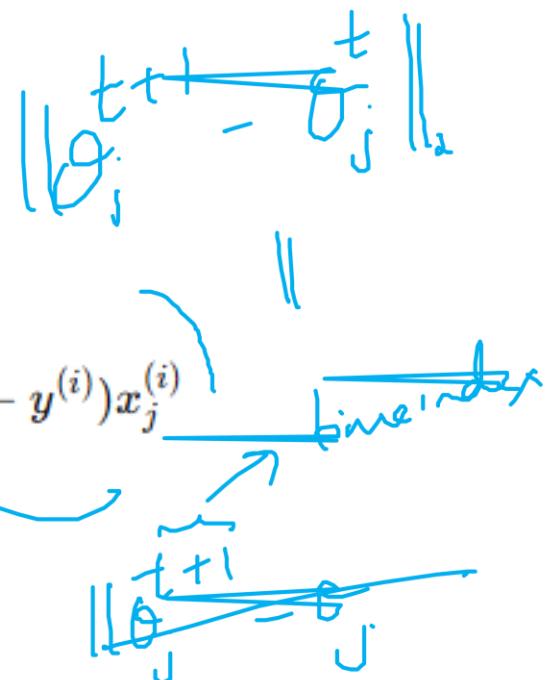
$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Thus, the update step becomes:

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

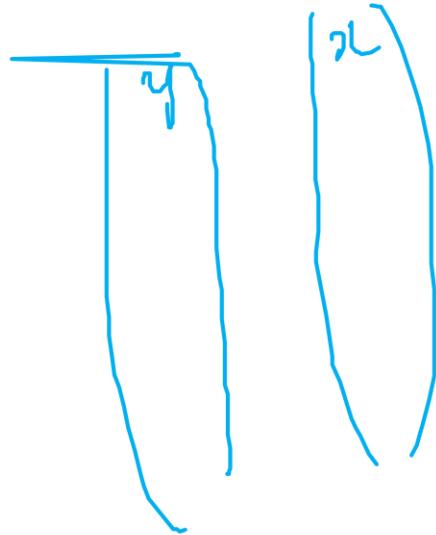
This process repeats until convergence (i.e., when $J(\theta)$ stops changing significantly).

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



Coding from Scratch

$$y = m\theta + \epsilon + \text{noise}$$



$\theta_0 + \theta_1 x_1$

$\theta^T x$ ~~plus~~ ~~noise~~

~~θ~~

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

$$\epsilon = \begin{bmatrix} 1 \\ \epsilon_1 \end{bmatrix}$$

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

Logistic Regression

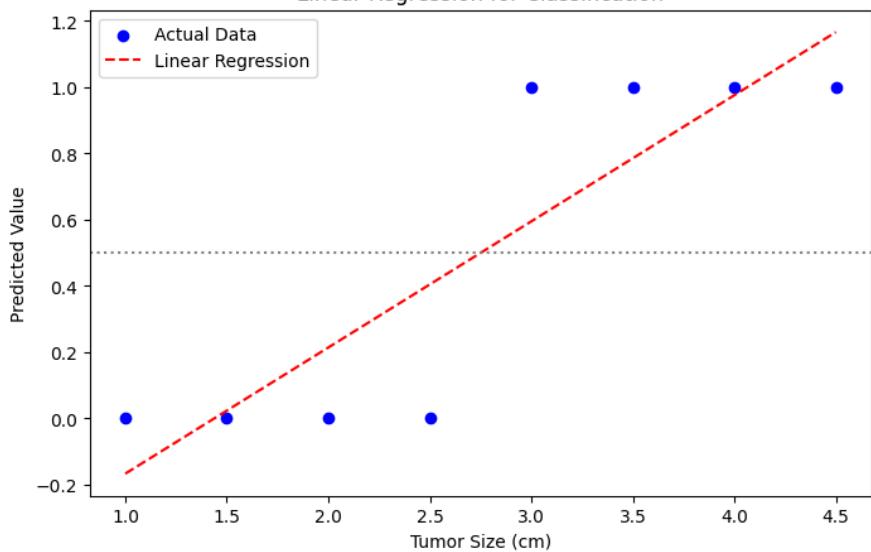
Scenario

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.

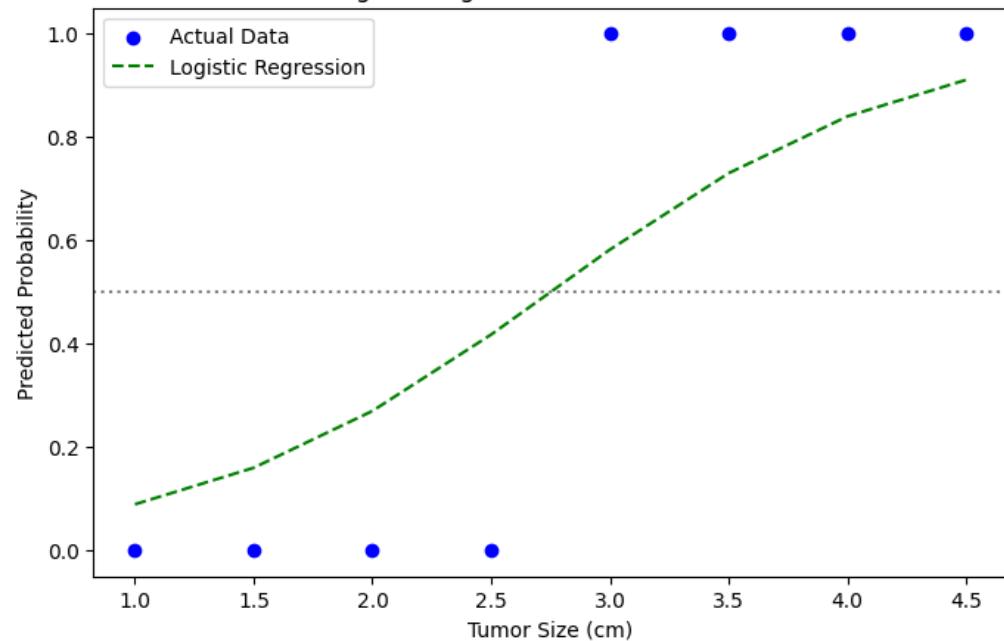
Hypothetical Data

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

Linear Regression for Classification



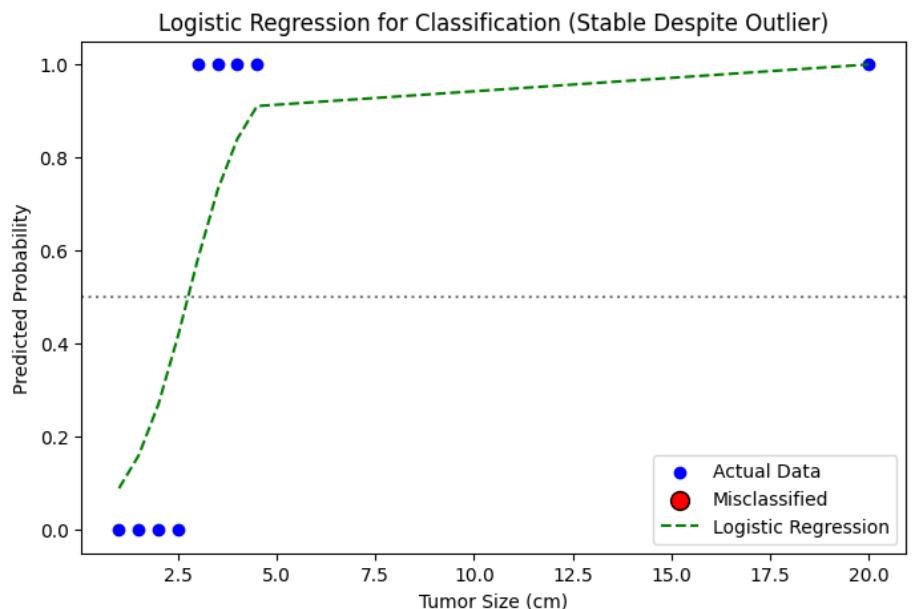
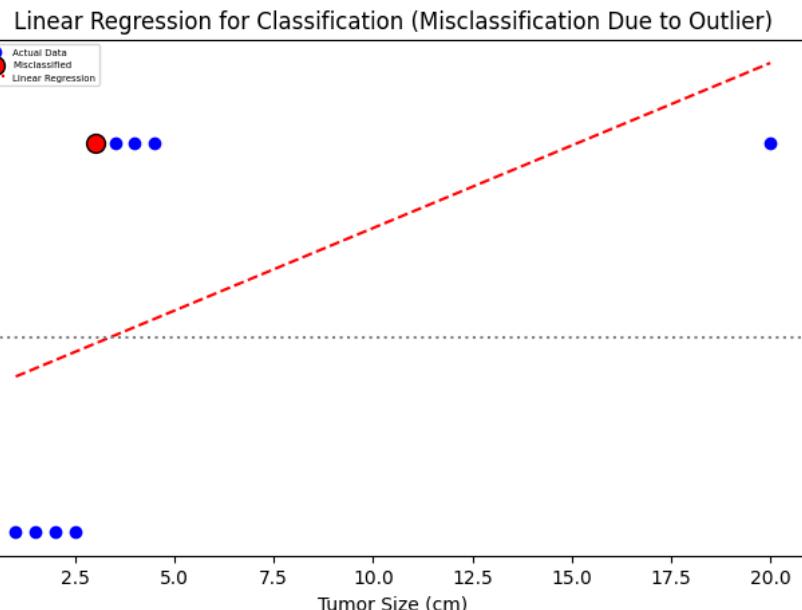
Logistic Regression for Classification



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



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

The hypothesis function in **logistic regression** is given by:

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

where:

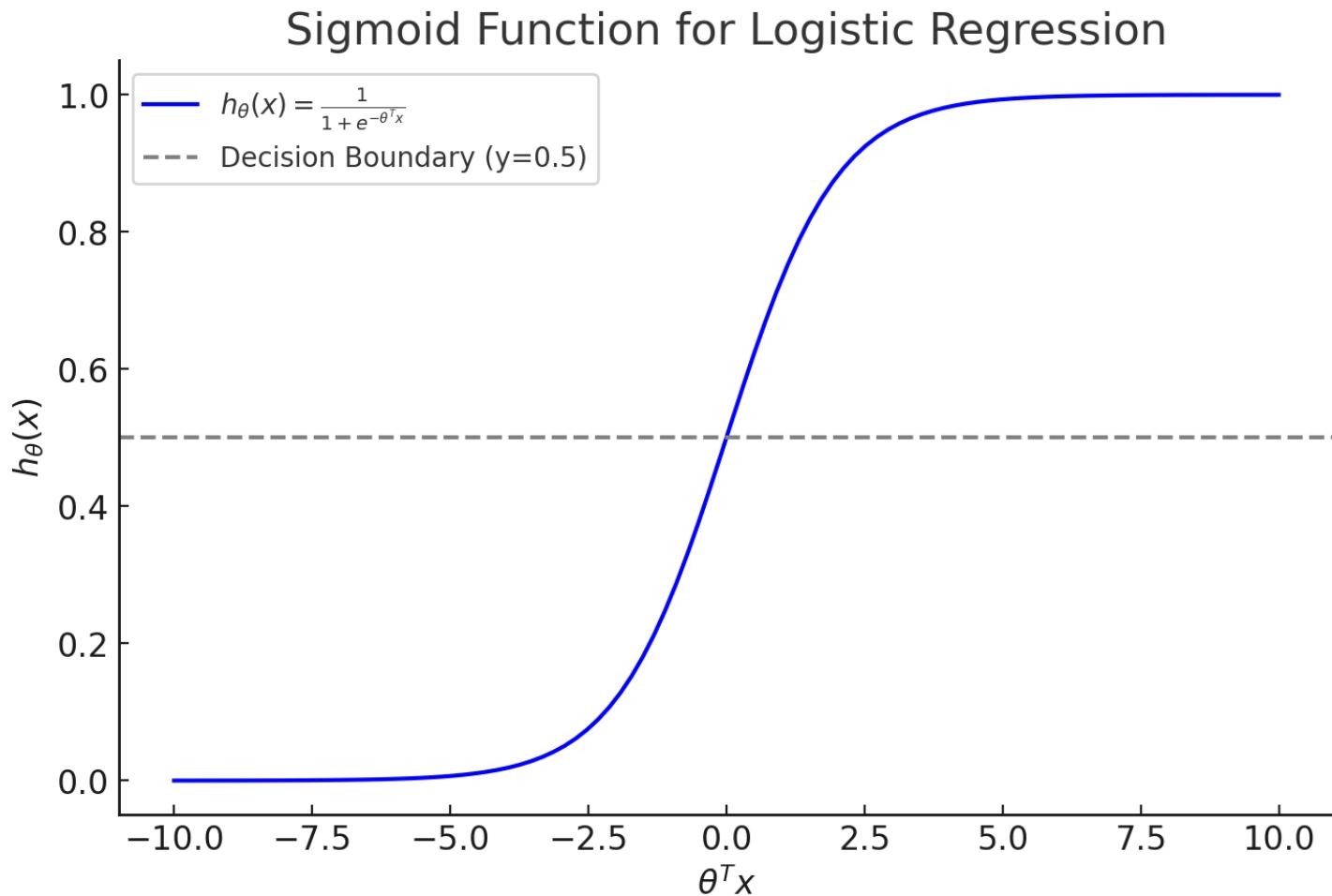
- θ is the parameter vector (weights).
- x is the feature vector (input data).
- $\theta^T x$ represents the **linear combination** of features and parameters.
- $g(z)$ is the **sigmoid function**, defined as:

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

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

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

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



Mathematical Formulation of Logistic Regression

We have already defined the **hypothesis function** as:

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

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

1. Probability Interpretation

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

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

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

Likelihood Function

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:

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

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.

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 θ that **maximize** this function. In other words, we want to find the **best** θ such that the predicted probabilities align as closely as possible with the actual labels in the training data.

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 θ that make the observed data most probable under our model.

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 $\frac{1}{m}$):

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

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**.

Gradient Descent for Parameter Estimation

The gradient of the **log-loss function** with respect to θ_j is:

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

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.

We update θ_j using **gradient descent**:

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

where:

- α is the **learning rate**, controlling how big the update steps are.



Coding



























