

CS60050

MACHINE LEARNING

Logistic Regression

Somak Aditya

Assistant Professor

Sudeshna Sarkar

Department of CSE, IIT Kharagpur

August 11, 2023



Logistic Regression for Classification

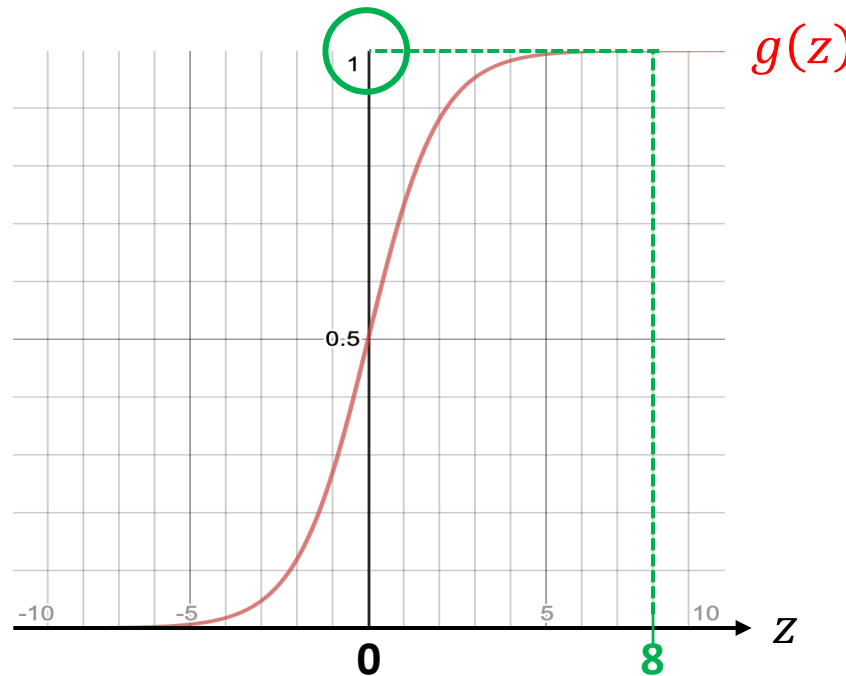
Regression vs. Classification

We want the possible outputs of $f_{\theta}(x) = \theta^T x$ to be discrete-valued

Use an **activation function** (e.g., **sigmoid or logistic function**)

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

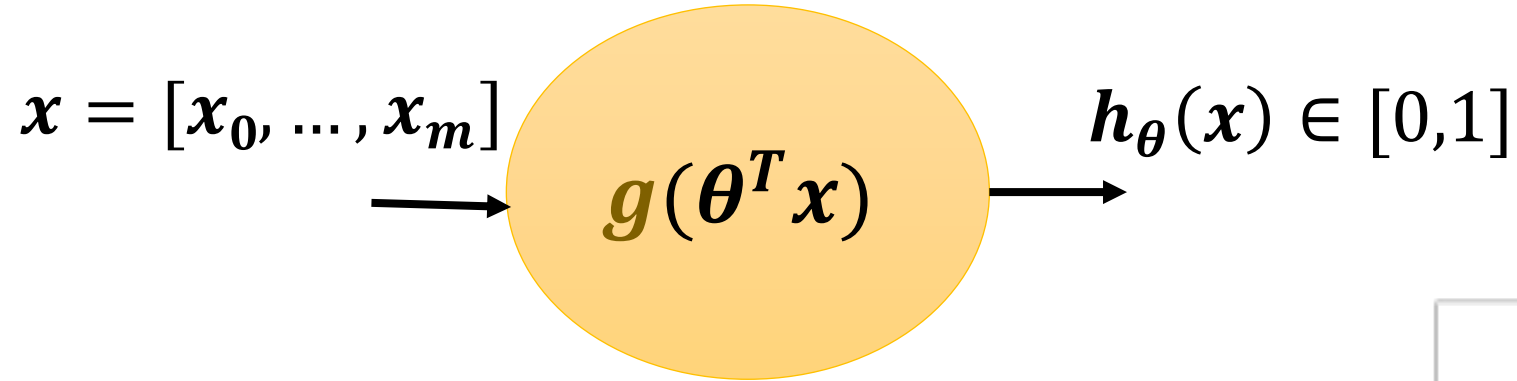
$z \in \mathbb{R}$, but
 $g(z) \in [0,1]$



If $y = 1$, we want $g(z) \approx 1$ (i.e., we want a correct prediction)
For this to happen, $z \gg 0$

If $y = 0$, we want $g(z) \approx 0$ (i.e., we want a correct prediction)
For this to happen, $z \ll 0$

Classification

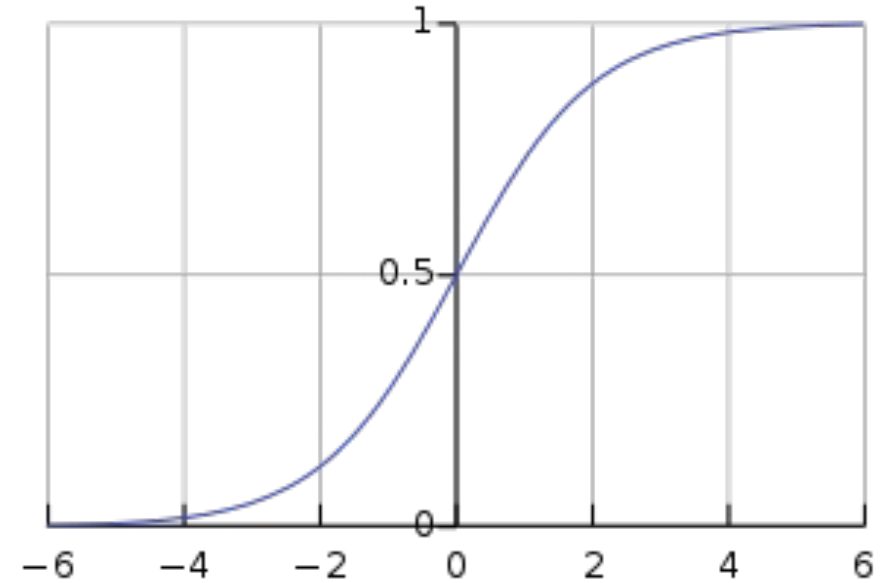


$$h_\theta(\mathbf{x}) = g(\theta^T \mathbf{x})$$
$$g(z) = \frac{1}{1 + e^{-z}}$$

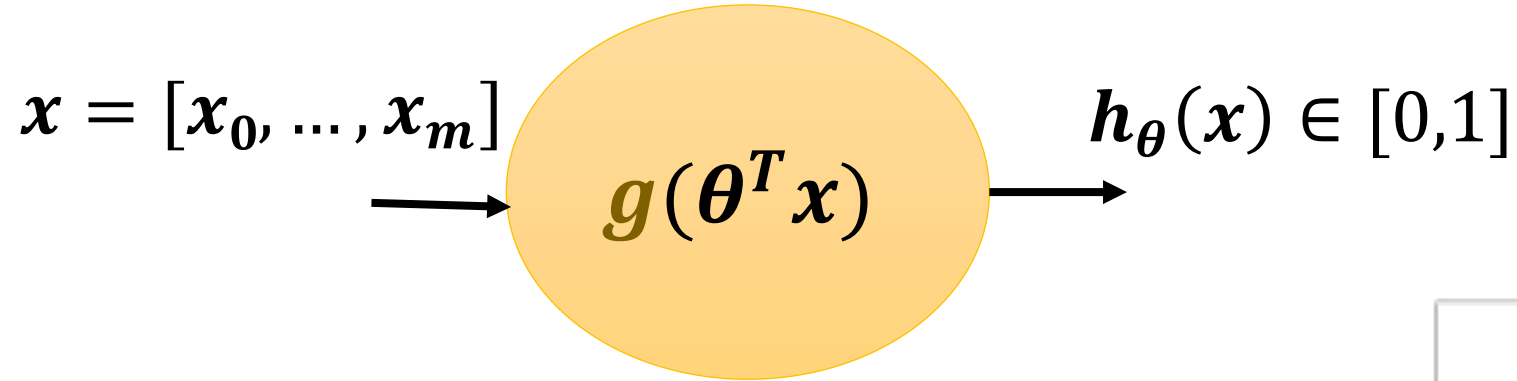
Thresholding:

predict “y = 1” if $h_\theta(\mathbf{x}) \geq 0.5$

predict “y = 0” if $h_\theta(\mathbf{x}) < 0.5$



Classification



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

Thresholding:

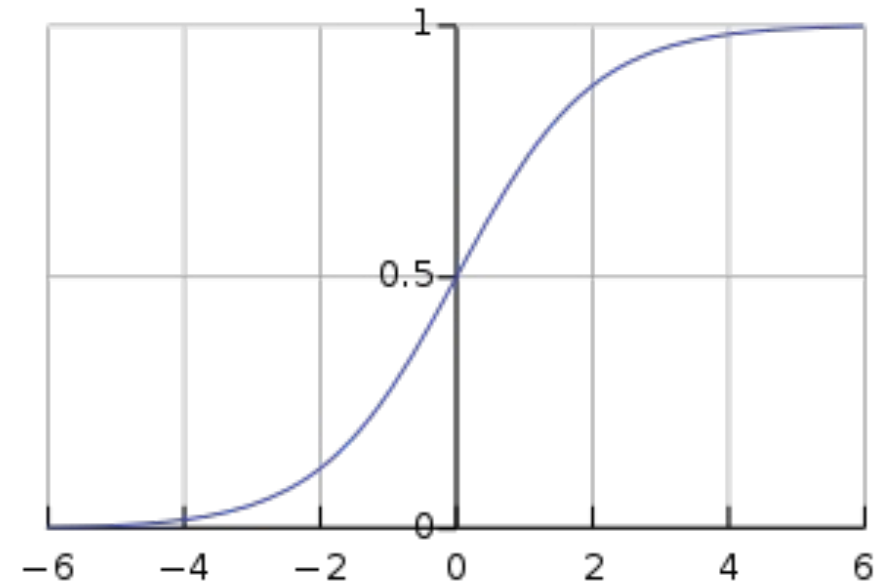
predict “y = 1” if $h_{\boldsymbol{\theta}}(\mathbf{x}) \geq 0.5$

$$\mathbf{z} = \boldsymbol{\theta}^T \mathbf{x} \geq 0$$

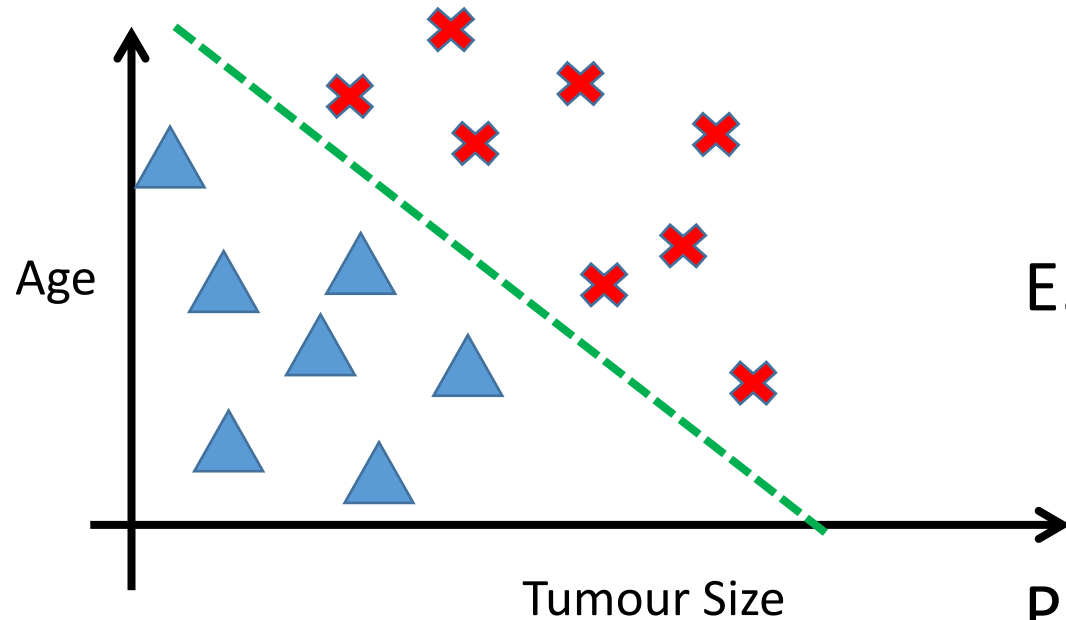
predict “y = 0” if $h_{\boldsymbol{\theta}}(\mathbf{x}) < 0.5$

$$\mathbf{z} = \boldsymbol{\theta}^T \mathbf{x} < 0$$

Alternative Interpretation: $h_{\boldsymbol{\theta}}(\mathbf{x}) =$
estimated probability that $y = 1$ on input \mathbf{x}



Decision boundary



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

E.g., $\theta_0 = -3, \theta_1 = 1, \theta_2 = 1$

Predict “ $y = 1$ ” if $-3 + x_1 + x_2 \geq 0$

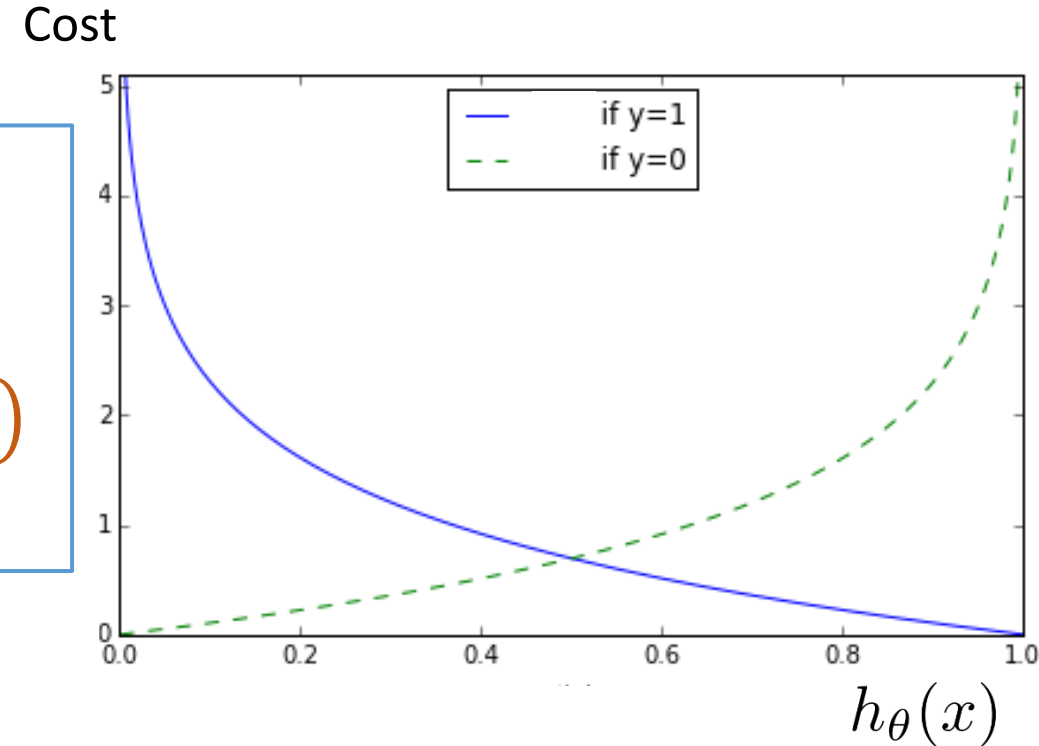
Cost function for Logistic Regression

Logistic Regression

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

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(\mathbf{h}_{\theta}(x^{(i)}), y^{(i)})$$

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



Gradient descent

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

Goal: $\min_{\theta} \text{loss}(\theta)$

Good news: Convex function!

Bad news: No analytical solution

Gradient descent

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

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

Gradient descent

Repeat {

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

}

(Simultaneously update all θ_j)

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

Gradient descent for **Linear Regression**

Repeat {

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

}

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

Gradient descent for **Logistic Regression**

Repeat {

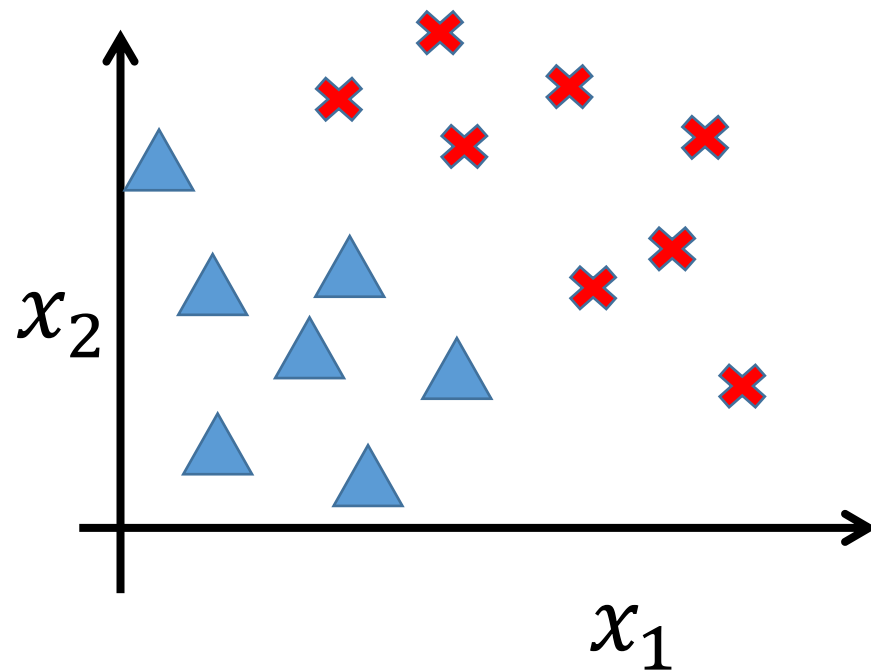
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

}

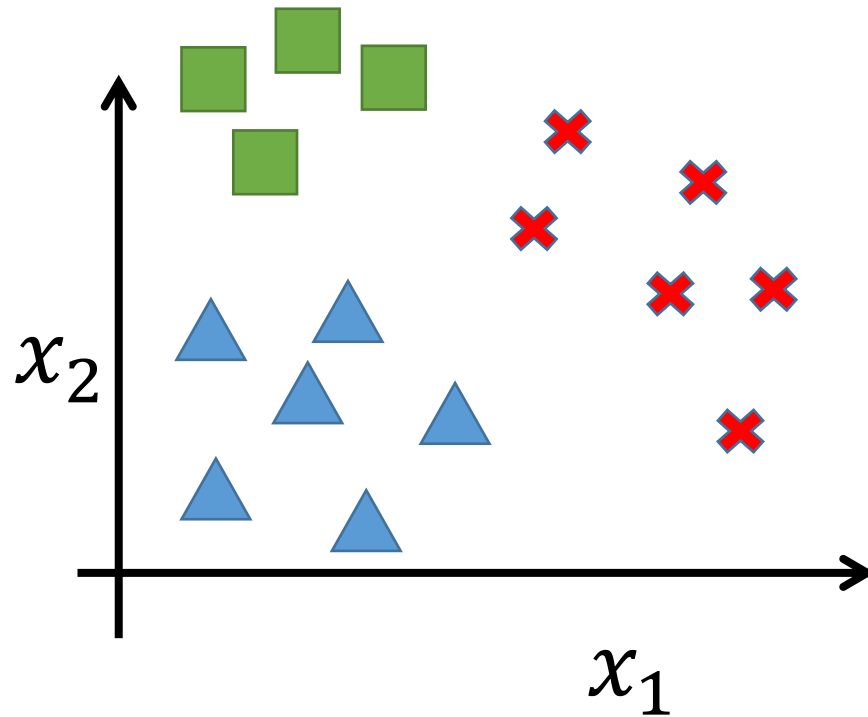
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^{\top} x}}$$

Multiclass classification

Binary classification



Multiclass classification



Multi-class Classification

- Multi-class Classification: y can take on K different values $\{1, 2, \dots, k\}$
- $f_{\theta}(x)$ estimates the probability of belonging to each class

$$P(y = k|x, \theta) \propto \exp(\theta_k^T x)$$

$$\theta = \begin{bmatrix} \vdots & \vdots & \vdots \\ \theta_1 & \theta_2 & \theta_k \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$P(y = k|x, \theta) = \frac{\exp(\theta_k^T x)}{\sum_{j=1}^K \exp(\theta_j^T x)}$$

$$J(\theta) = - \left[\sum_{i=1}^m \sum_{j=1}^K 1\{y^{(i)} = k\} \log \frac{\exp(\theta_k^T x^{(i)})}{\sum_{j=1}^K \exp(\theta_j^T x^{(i)})} \right]$$