

Logistic Regression

⇒ It is an algorithm for performing binary Classification.

Classification

Input

Output

⇒ In binary Classification, Output can only take two values.
 $\{0, 1\}$

⇒ Set of Labeled Input or Input Output pair

⇒ Objective: Learn the mapping between input & Output.

(i.e. Given Input we want to predict the Output)

⇒ Let Input be represented by a set of m input features. (x_1, x_2, \dots, x_n) where $x_i \in \mathbb{R} \quad i = \{1, \dots, n\}$

⇒ Let Output be represented by y where $y \in \{0, 1\}$

⇒ Let us consider a training set of m training examples:

$$\{(x^{(i)}, y^{(i)}) \mid x^{(i)} \in \mathbb{R}^n, y^{(i)} \in \{0, 1\}, i \in \{1, \dots, m\}\}$$

⇒ In Logistic Regression, we are looking for a linear decision boundary that separates the two labels.

Let $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n = 0$ the linear decision boundary parametrized by $(\theta_0, \theta_1, \dots, \theta_n)$.

2
⇒ For simplifying the notation let $x_0 = 1$ be a dummy feature. So decision boundary becomes

$$\Rightarrow \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n = 0$$

$$\Rightarrow \boxed{\theta^T x = 0}, \text{ where } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \text{ and } x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

⇒ x for which $\theta^T x \geq 0$ we will classify it as $y = 1$.

⇒ and for which $\theta^T x < 0$ we will classify it as $y = 0$.

⇒ Given a decision boundary parametrized by θ , we ~~are looking~~ want to assign probability to all x .

(i.e. If the input is x , what is the probability that corresponding $y = 1$)

⇒ So we are looking for a function $h_\theta(x)$ such that,

- $0 \leq h_\theta(x) \leq 1 \quad \forall x \in \mathbb{R}^{n+1}$

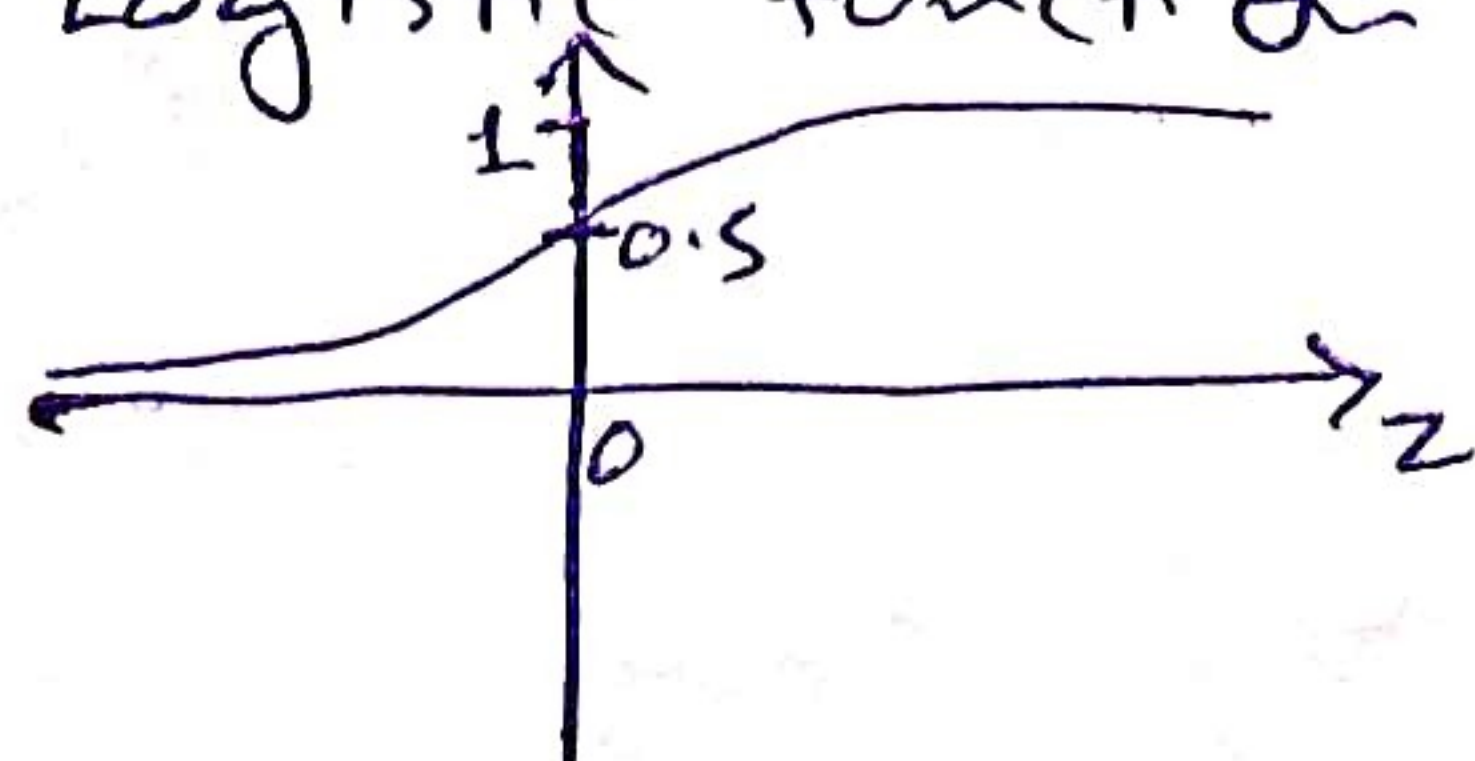
- $h_\theta(x) = 0.5 \quad \forall \theta^T x = 0$

- $h_\theta(x) > 0.5 \quad \forall \theta^T x > 0$ ($h_\theta(x) \uparrow$ as $\theta^T x \uparrow$)

- $h_\theta(x) < 0.5 \quad \forall \theta^T x < 0$ ($h_\theta(x) \downarrow$ as $\theta^T x \downarrow$)

⇒ One such function is Logistic function:

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



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

⇒ For estimating θ given the training set, we can maximize the likelihood of θ .

$$L(\theta) = P(y^{(1)}, \dots, y^{(n)} | x^{(1)}, \dots, x^{(n)}; \theta)$$

$$L(\theta) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta)$$

Assuming all the data
in the training set are
independent

~~log~~

⇒ Let $\ell(\theta) = \log L(\theta)$ be the log likelihood. Maximizing $L(\theta)$ is equivalent to maximizing $\ell(\theta)$ as $\log(x)$ is a monotonic increasing function of x .

$$\ell(\theta) = \sum_{i=1}^n \log P(y^{(i)} | x^{(i)}; \theta)$$

$$= \sum_{i=1}^n \log \left(h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1 - y^{(i)}} \right)$$

$$\ell(\theta) = \sum_{i=1}^n y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)}))$$

⇒ To maximize $\ell(\theta)$ we can use gradient ascent.

$$\theta := \theta + \alpha \nabla_{\theta} \ell(\theta)$$

Learning rate

* Calculating the gradient ($\nabla_{\theta} l(\theta)$)

$$\nabla_{\theta} l(\theta) = \sum_{i=1}^m (y^{(i)} \nabla_{\theta} \log h_{\theta}(x^{(i)}))$$

$$(1 - y^{(i)}) \nabla_{\theta} \log (1 - h_{\theta}(x^{(i)}))$$

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

$$\frac{1}{h_{\theta}(x^{(i)})} * \nabla_{\theta} h_{\theta}(x^{(i)})$$

$$(1 + e^{-\theta^T x^{(i)}}) * \frac{\theta e^{-\theta^T x^{(i)}}}{(1 + e^{-\theta^T x^{(i)}})^2}$$

$$\frac{\theta e^{-\theta^T x^{(i)}}}{(1 + e^{-\theta^T x^{(i)}})^2}$$

$$\frac{1}{1 - h_{\theta}(x^{(i)})} * -\nabla_{\theta} h_{\theta}(x^{(i)})$$

$$\frac{1}{\frac{e^{-\theta^T x^{(i)}}}{1 + e^{-\theta^T x^{(i)}}}} * \frac{-\theta e^{-\theta^T x^{(i)}}}{(1 + e^{-\theta^T x^{(i)}})^2}$$

$$\frac{1 + e^{-\theta^T x^{(i)}}}{\cancel{e^{-\theta^T x^{(i)}}}} * \frac{-\theta \cancel{e^{-\theta^T x^{(i)}}}}{(1 + e^{-\theta^T x^{(i)}})^2}$$

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

$$\nabla_{\theta} l(\theta) = \sum_{i=1}^m \left\{ \frac{\theta}{1+e^{-\theta^T x^{(i)}}} \left(y^{(i)} e^{-\theta^T x^{(i)}} - 1 + y^{(i)} \right) \right\}$$

$$\downarrow$$

$$\frac{\theta}{1+e^{-\theta^T x^{(i)}}} \left(y^{(i)} (1+e^{-\theta^T x^{(i)}}) - 1 \right)$$

$$\nabla_{\theta} l(\theta) = \sum_{i=1}^m \left(y^{(i)} - \frac{1}{1+e^{-\theta^T x^{(i)}}} \right) \theta$$

$$\boxed{\nabla_{\theta} l(\theta) = \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) \theta}$$

\Rightarrow So update rule for gradient ascent is:

$$\boxed{\theta := \theta + \alpha \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) \theta}$$

~~————— X ————— X —————~~