



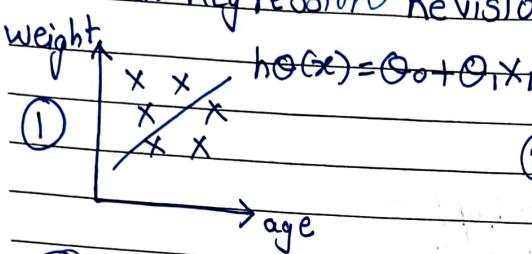
Semester : VI

Subject : CSC601 Data Analytics and Visualization

Academic Year: 2023- 2024

Logistic Regression

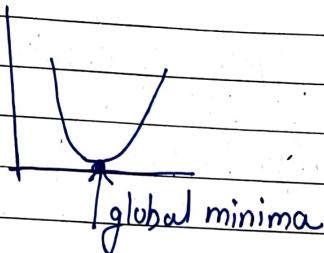
Linear Regression Revision



Cost function

$$\textcircled{2} \quad J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x_i) - y_i)^2$$

③ Gradient decent \rightarrow minimise cost funⁿ to reach to global minima



Linear Regression is used to solve regression problem.

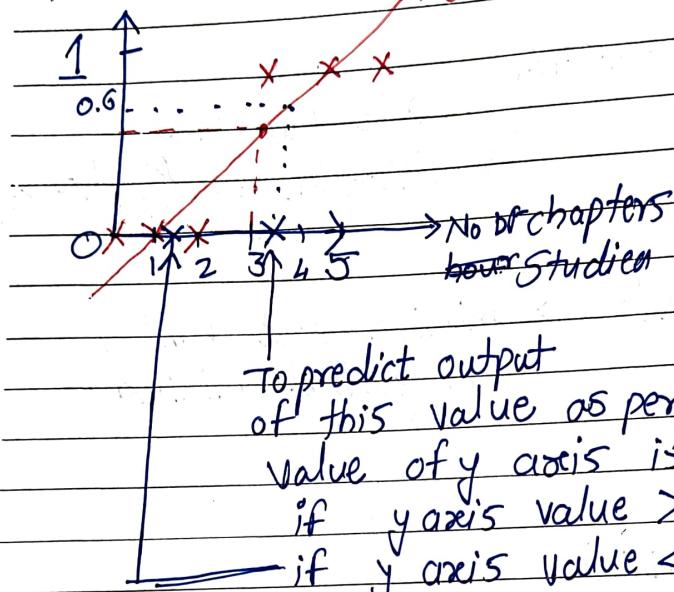
Logistic Regression is used to solve classification problem / Binary classification Problem

No of chapters	Pass/Fail	No of chapters Studied	Pass/Fail
1	XX	1	Pass
2	XX	2	Fail
3	X	0	Fail
4	(X)	1	Fail
5	Pass/Fail	2.5	Pass



Can we use linear regression of solving this problem?

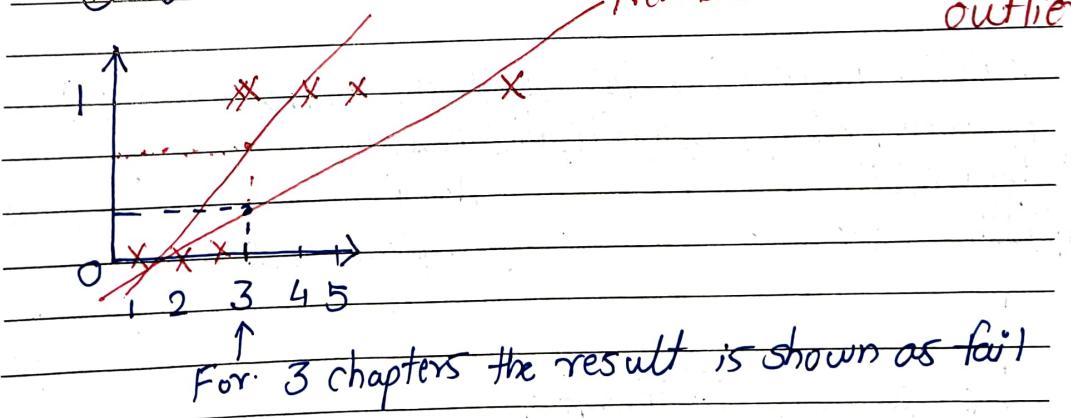
$$h(x) = \theta_0 + \theta_1 x$$



The problem we have here are

① outliers

New Best Fit line due to outliers





(2)

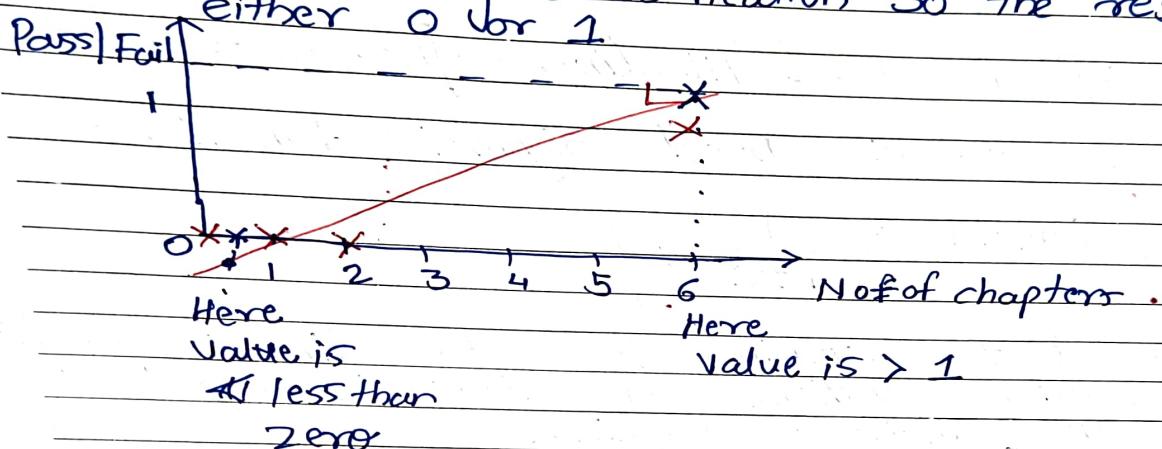
Semester : VI

Subject : CSC601 Data Analytics and Visualization

Academic Year: 2023- 2024

Classification

② This binary classification so the result will be either 0 or 1



② The problem > 1 or < 0 (called squashing)

To address these two problems linear regression can not be used.

① outliers

② output values > 1 or < 0

The classification problems can not be solved using linear regression.



For logistic regression, to address squashing and outliers we need to

The eqn for best fit line - $h\theta(x) = \theta_0 + \theta_1 x$

Also we will need sigmoid function

$$\text{Sigmoid} = \frac{1}{1+e^{-x}}$$

We need to combine these two lines equations to generate equations of logistic regression that can handle squashing.

$$h\theta(x) = g(\theta_0 + \theta_1 x)$$

Sigmoid activation function

$$\text{Sigmoid fun} \rightarrow g = \frac{1}{1+e^{-z}} \quad \text{here } z = (\theta_0 + \theta_1 x)$$

$$h\theta(x) = \frac{1}{1+e^{-(\theta_0 + \theta_1 x)}}$$

Hypothesis of logistic regression



(3)

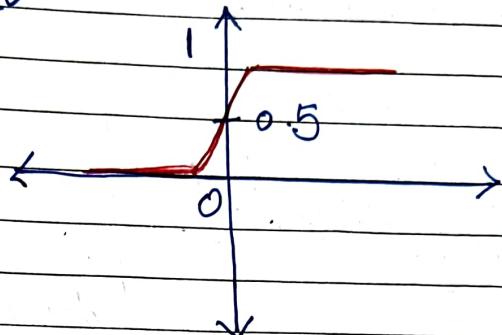
Semester : VI

Subject : CSC601 Data Analytics and Visualization

Academic Year: 2023- 2024

Sigmoid function

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



If $z \leq 0$

$$g(z) \leq 0.5$$

If $z \geq 0$

$$g(z) \geq 0.5$$

The Problem Statement of logistic regression

Training Dataset $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$

$$x \in \{x_0, x_1, \dots, x_m\} \quad y \in \{y_0, y_1, \dots, y_m\}$$

$$y \in \{0, 1\}$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}$$

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

$$\begin{aligned} z &= \theta_0 + \theta_1 x \\ &= \theta^T x \end{aligned}$$



Academic Year: 2023- 2024

Semester : VI Subject : CSC601 Data Analytics and Visualization

Cost Function

Linear Regression

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Here,
 $h_\theta(x) = \theta_0 + \theta_1 x$

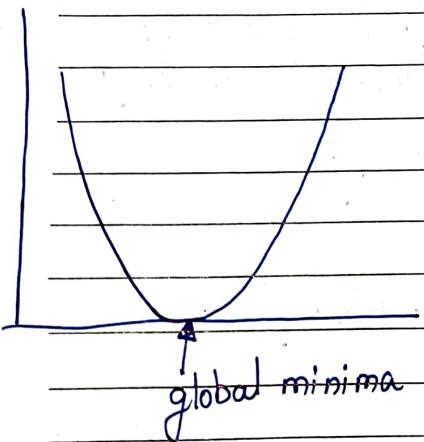
Logistic Regression

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

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

where $\theta^T x = \theta_0 + \theta_1 x$

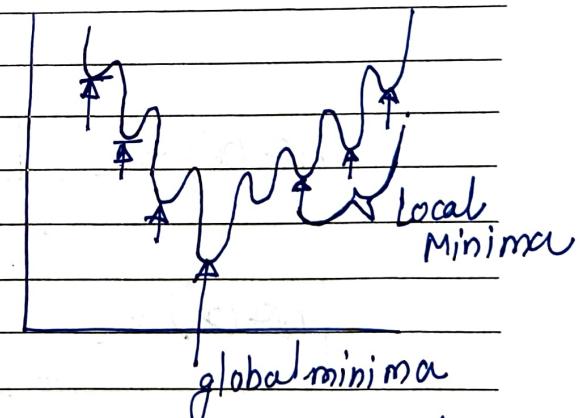
The $h_\theta(x)$ is a convex function



No local minima

The $h_\theta(x)$ is a non convex function

Due to sigmoid function we are generating non convex function



If it point reaches to local minimum it is stucked there only it cannot

reach to global minima after that



Semester : VI

Subject : CSC601 Data Analytics and Visualization

Academic Year: 2023- 2024

So, the cost function of logistic regression

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^i) - y^i)^2$$

$$\text{cost}(h_\theta(x^i), y^i)$$

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

$$z = \theta_0 + \theta_1 x$$

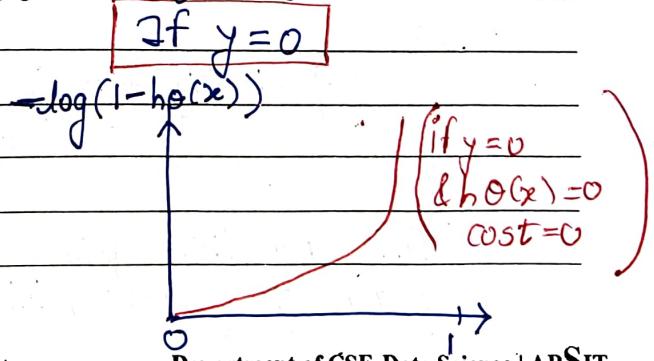
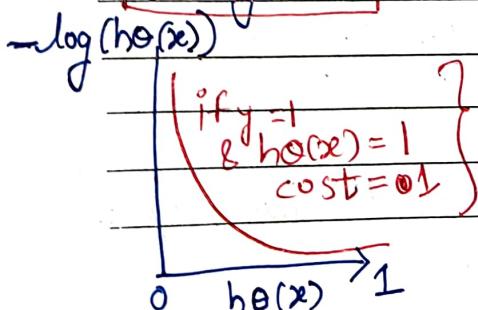
This cost function generates non-convex function so to generate the convex function we need to rewrite the cost function as

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

This function generates a curve like gradient decent.

If $y = 1$

If $y = 0$





Semester : VI Subject : CSC601 Data Analytics and Visualization

$$\text{cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

$$\begin{aligned} \text{if } y &= 1 \\ &= -\log(h_\theta(x)) \end{aligned}$$

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

The cost function of ~~Re~~ logistic regression is

$$J(\theta_0, \theta_1) = \frac{-1}{2m} \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))$$

And convergence repeat function will be

Convergence Repeat

$$\theta_j \approx \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

This function will continue till we reach to global minima.



Semester : VI

Subject : CSC601 Data Analytics and Visualization

Academic Year: 2023- 2024

Example 1 (Logistic Regression)

The dataset of pass or fail in an exam of 5 students is given in the table.

Hours Study	Pass(1)/Fail(0)
29	0
15	0
33	1
28	1
39	1

Use logistic regression as classifier to answer the following questions.

1. Calculate the probability of pass for the student who studied 33 hours

2. At least how many hours student should study that makes he will pass the course with the probability of more than 95%

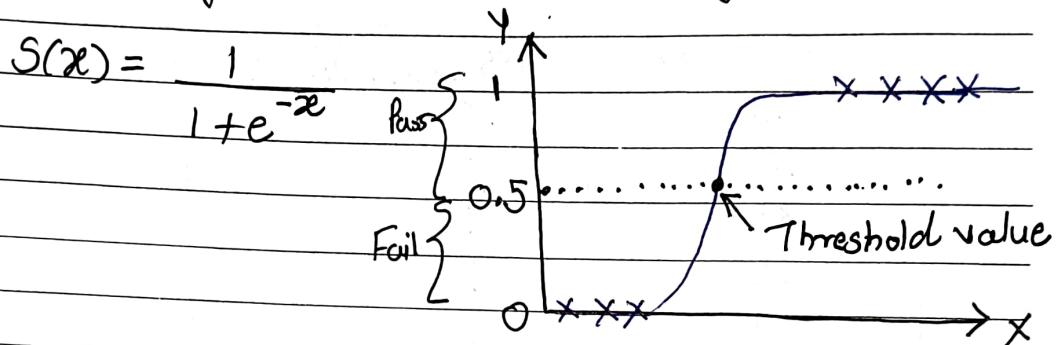
Assume the model suggested by the optimizer for odds of passing the course is,

$$\log(\text{odds}) = -64 + 2 * \text{hours}$$



Solution

We use sigmoid function in logistic regression



i) Calculate the probability of pass for the student who studied 33 hours.

$$\text{h}(x) h(x) = \frac{1}{1 + e^{-z}}$$

$$\text{Given } \log(\text{odds}) = z = -64 + 2 * \text{hours}$$

Let's find value of z ,

$$z = -64 + 2 * (33)$$

$$z = -64 + 66$$

$$\boxed{z = 2}$$

(2)



Semester : VI

Subject : CSC601 Data Analytics and Visualization

Academic Year: 2023- 2024

Now, let's put value of z in equation given below,

$$h_0(x) = p = \frac{1}{1 + e^{-z}}$$

$$p = 0.88$$

Thus, if student studies 33 hours, then there is 88% chance that the student will pass the exam.

- ii) At least how many hours student should study that makes he will pass the course with the probability of more than 95%.

$$h_0(x) = p = \frac{1}{1 + e^{-z}} = 0.95$$

$$0.95(1 + e^{-z}) = 1$$

$$0.95 \cdot e^{-z} = 1 - 0.95$$

$$e^{-z} = \frac{0.05}{0.95}$$

$$e^{-z} = 0.0526$$



Semester : VI Subject : CSC601 Data Analytics and Visualization

Academic Year: 2023-24

We need to know the value of z

$$\text{we have } e^{-2} = 0.0526$$

Take natural log both sides,

$$\ln(e^{-2}) = \ln(0.0526) \quad \textcircled{1}$$

As per natural logarithmic,

$$\ln(e^x) = xe$$

$$\text{So } \ln(e^{-z}) = -z$$

Replace $\ln(e^{-z})$ with $-z$ eq in eqn $\textcircled{1}$

$$-z = \ln(0.0526)$$

$$-z = -2.94$$

$$z = 2.94$$

The given equation of log(Odds) or z

$$z = -64 + 2 * \text{hours}$$

$$2.94 = -64 + 2 * \text{hours}$$

$$2 * \text{hours} = 2.94 + 64$$

$$\text{hours} = \frac{2.94 + 64}{2}$$



PARSHWANATH CHARITABLE TRUST'S
A.P. SHAH INSTITUTE OF TECHNOLOGY
Department of Computer Science and Engineering
Data Science



(3)

semester : VI

Subject : CSC601 Data Analytics and Visualization

Academic Year: 2023- 2024

$$\text{hours} = \frac{66.94}{2}$$

$$\boxed{\text{hours} = 33.47}$$

The student should study at least 33.47 hours, so that he will pass the exam with more 95% probability!