

Ma

Jiff

Regres

Fu

Fun

traini
set

trai

)

①

→

→

→

→

①

Types of learning:

Supervised and unsupervised

② Types of problems

Regression Classification and
clustering

③ Cost function, optimization -

Gradient Descent algorithm.

④ Sampling, decision boundary,
under-fitting and overfitting of models

⑤ Bias-Variance tradeoff

(Supervised).

Regression: Linear, Regularization
problem - Ridge and Lasso.

Classification - Logistic Regression.

(Unsupervised).

Intro - to - Artificial Neural Network

. Model validation and selection

↳ Accuracy

Confidence interval

Confusion Matrix

Precision, Recall

Machine Learning (M.L.)

Types of learning

Supervised Unsupervised

Regression Classification

Machine Learning : Learning from Data

it is about predicting the future

based on the past. M.L. algorithm's

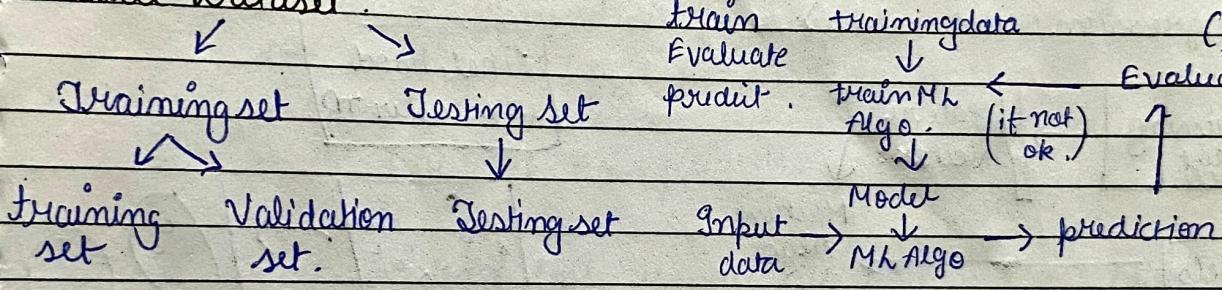
attempt to optimize along a certain

matrix dimension; i.e. minimize error or

interval. Maximize the likelihood of their prediction

→ optimizing an error/cost function.

Full Dataset



Training is a process of making the system able to learn.

→ Training and testing set come from same distribution.

① Supervised and Unsupervised:

(Unlabelled data)

Supervised .

- Training Data
- Both Input and outputs.
- are known
- Classification
- Naive Bayes Algo.

supervised

learning e.g;

election exit poll,
climate .

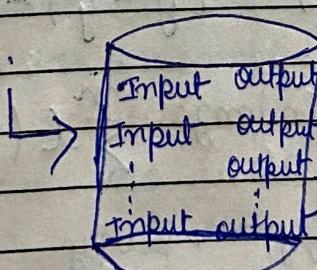
Unsupervised .

e.g; some categories .

→ only Inputs

→ forms clusters

→ Clustering, K-mean



New I/P

↓

Learning → Model

algo

↓

O/P

Training

data .

(Labeled data)

supervised

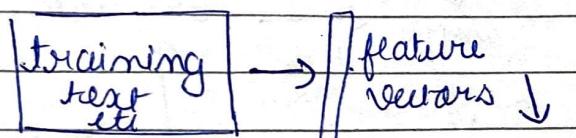
Classification	Regression
Logistic Regression	Linear
Naïve Bayes	Ridge
Linear Discriminant (LDA)	Lasso

- KNN, Decision trees, Random forest
- - Prediction
 - Classification (discrete labels), Regression (real values)

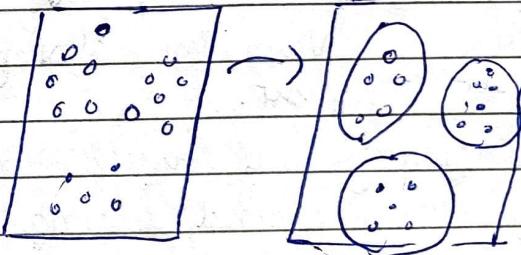
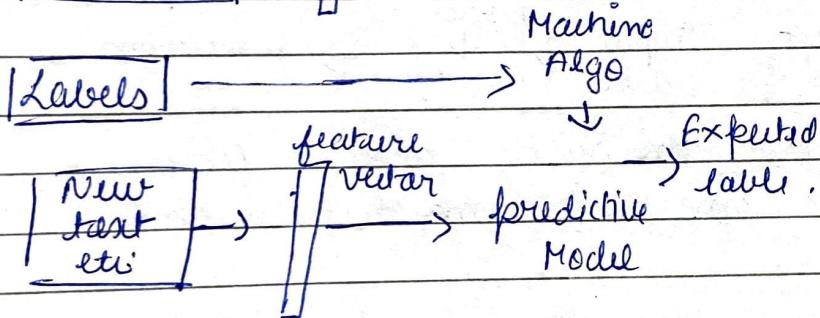
unsupervised

PCA
K-mean clustering
Hierarchical clustering
DB Scan.

- - Clustering
 - $P(E)$ distribution estimation
 - Finding association (features)
 - Dimension reduction,



Likelihood
or cluster ID.



② Key Concepts (cost function, optimization) etc :

Loss function: $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, measures the difference between a prediction and an actual output.

$$\text{Eg: } \{(\hat{y}, y) = (\hat{y} - y)^2\}$$

ML optimization problem: minimize $\sum_0^m l(h_0(x, y))$

every ML algo has this form, just specify

→ what is the hypothesis function (h_0)

→ what is loss function (l)

→ How do we solve the optimization problem.

Loss function v.s Cost function

- The loss function computes the error for a single training eg , while the cost function is the average of the loss functions of the entire training set.
- if we have m training data like this $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$
- $y_1 = \text{output of the model for training eg } x_1$
- $y_2 = \text{expected output/true value for training eg } x_1$
- Loss function $L(y_1, y_2)$ to defines the error/difference between y_1 and y_2 for the single training eg x_1 .
- Loss function refers to an error in model output for an individual sample

* If we want to find loss over all the training ex. present in a training-set , we refer to it as the (cost function). i.e total or average loss over all training examples.

Key differences.

feature	Loss	Cost
Definition	Measures error for a single data point	Measures error (avg) over the entire dataset.
Scope	Individual sample	Whole dataset
Usage	Evaluates single instance performance	Used for optimization and parameter updates.

Example

Loss / Error / Cost - objective function.

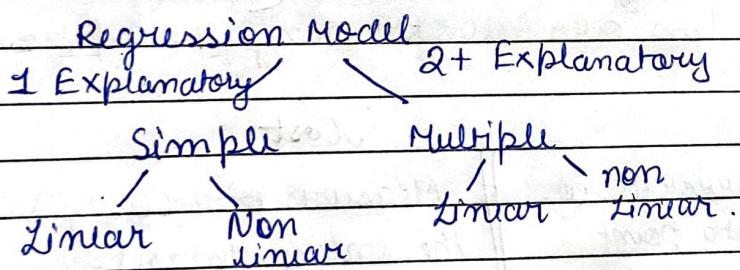
The more general scenario is to define an objective function first, which we want to optimize. This objective function could be to :

1. Min a (MSF), cost or loss function
2. Max the posterior probabilities (eg: Naive Bayes)
3. Max a fitness function (genetic programming)
4. Max the total reward / value function (Reinforcement Learning)
5. Max information gain / Min child node impurities.
6. Max -log -likelihood

Supervised learning (Regression) :

Regression: Relationship b/w one dependent variable (y) and Models explanatory variable(s).

- it use the equation to set up a relationship
- one or more numerical or categorical independent (Explanatory) variables.
- mainly for the prediction and estimation of (continuous numeric variables).



Linear Equation:

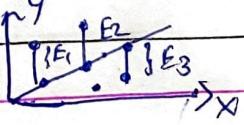
A linear relationship between the predictor / independent variable (x) and the response / dependent variable (y) is assumed.

$$y = \beta_0 + \beta_1 * x \text{ or } y = mx + c$$

Least square: 'Best fit' means difference between the actual y value and predicted \hat{y} values are a min.

→ positive differences off-set negative. (So square errors!)

$$\sum (y_i - \hat{y}_i)^2 = \sum (E_i)^2$$



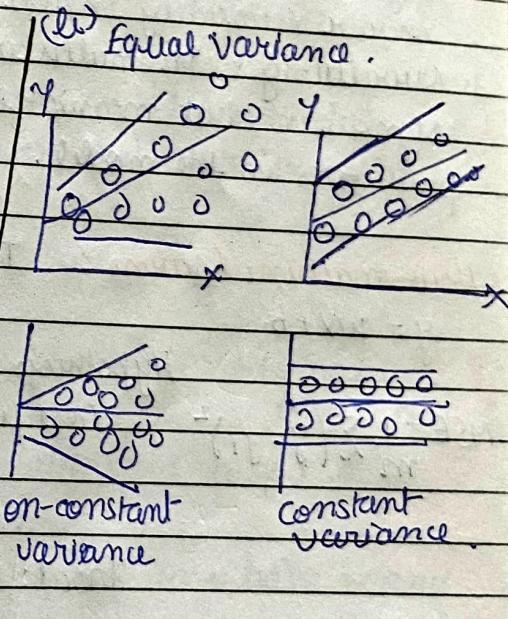
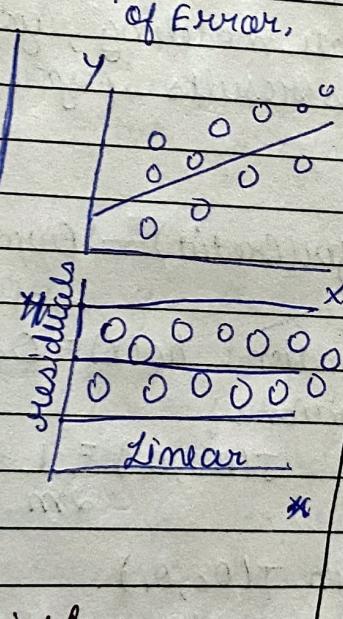
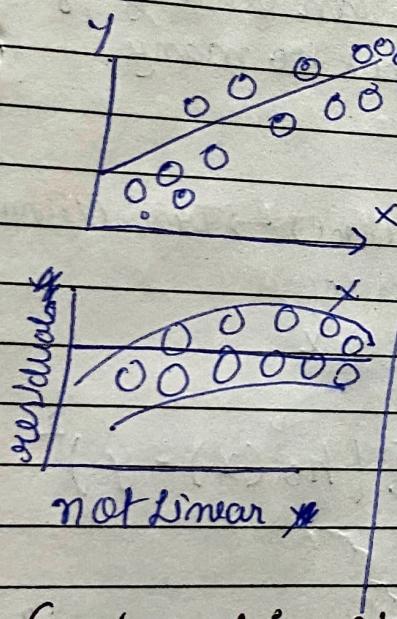
- Assumptions for linear regression:
- Linearity: x and y must be linear.
 - Independence of error
 - Normality of error
 - Equal variance.

Residual Analysis. $c = y_i - \hat{y}$.
 $\text{Residual} = \text{observed} - \text{predicted}$

$y_i = 48$
 $\hat{y} = 34$
 $y_i - \hat{y} = 14$

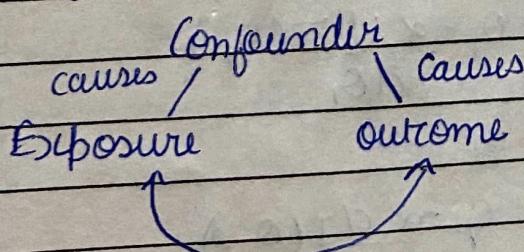
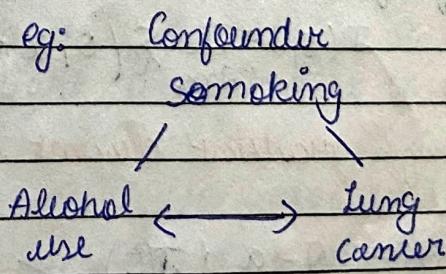
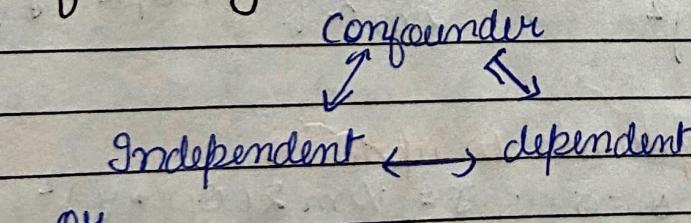
$q = 34$

(a) Linearity and Independence of Error.



(b) Equal Variance.

Confounding Variable:



Distorted association
when failing to
control for confounder.

Multivariate Regression Pitfalls

1. Multicollinearity: two variables that measure the same thing or similar things (eg weight and BMI) are both included in a multiple regression model, they will, in effect cancel each other out and generally destroy your model.
2. Residual confounding: we cannot completely wipe out confounding simply by adjusting for variables in multiple regression unless variables are measured with zero error (usually impossible)
3. Overfitting: In multivariate modeling, you can get highly significant but meaningless results if you put too many predictors in the model.

How machine learns: Infer(Predict) \rightarrow Error(cost) \rightarrow Train(learn)
 $y = wx + b$

Hypothesis function: $h_\theta(x) = \underline{\theta_0 + \theta_1 x}$

$$MSE = \frac{1}{m} \sum (\hat{y} - g_i)^2 \quad \text{cost function: } J(\theta_0, \theta_1) = \frac{1}{2m} \sum (h_\theta(x) - y)^2$$

Goal: Min $J(\theta_0, \theta_1)$

If $\theta_0 = 0$

$$J(\theta_1) = \frac{1}{2m} \sum (\theta_1 x - y)^2$$

$$\text{else: } = \frac{1}{2m} \sum (\theta_0 + \theta_1 x - y)^2$$

Using gradient descent algorithm to update θ :

$$\text{if } \theta_0 = 0 : J = \frac{1}{2m} \sum (\theta_1 x - y)^2$$

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

$$\theta_1 = \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

If $\theta_0 \neq 0$

$$\theta_0 = \theta_0 - \alpha \frac{d}{d\theta_0} J(\theta_0, \theta_1) \quad \theta_1 = \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_0, \theta_1)$$

$$= \theta_0 - \alpha \frac{1}{m} \sum (h_\theta(x) - y) \quad \theta_1 = \theta_1 - \alpha \frac{1}{m} \sum (h_\theta(x) - y) \cdot x$$

Performance Metrics:

1. Mean Squared Error (MSE):

$MSE = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$ Measures the avg. squared difference between actual and predicted value.

2. Root Mean Squared Error (RMSE):

$RMSE = \sqrt{MSE}$ gives an interpretable error in the original unit.

3. Mean Absolute Error (MAE):

$MAE = \frac{1}{n} \sum |y_i - \hat{y}_i|$ the avg. absolute difference between actual and predicted values.

4. R-Squared (r^2)

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} \quad SS_{\text{res}} = \sum (y_i - \hat{y}_i)^2 \quad (\text{residual sum of } \hat{\sigma}^2)$$

$$SS_{\text{tot}} = \sum (y_i - \bar{y})^2 \quad (\text{total sum of } \hat{\sigma}^2).$$

More about Gradient Descent :

Linear relationship: $y = \alpha_0 + \alpha_1 x$

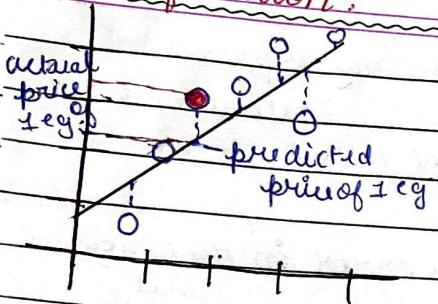
$\begin{cases} y = \text{target} \\ x = \text{feature} \end{cases} ; \begin{cases} \alpha_0 \text{ and } \alpha_1 \text{ are model parameters} \\ \end{cases}$

As per our linear regression model, we need to fit a straight line to it $y = \alpha_0 + \alpha_1 x$, depending on values of α_0 and α_1 , we can have many possibilities, which look promising.

We need to settle the case on a value of parameters α_0 and α_1 corresponding to which straight line fits best to data.

→ For this we need to agree on a metric to judge best fit and then we can choose that straight line which performs best on that metric → { cost function } . To resolve this

Cost function:



Error Term (e): $y_{\text{actual}} - y_{\text{predicted}}$.

$$e_i = y_{\text{actual}} - y_{\text{predicted}}$$

for i^{th} term ;

$$e_i = y_{\text{predicted}} - y_{\text{actual}}$$

Now e_i can be +ve or -ve depending on the value of y_{actual} and $y_{\text{predicted}}$ so square e_i .
So now;

$$\text{Cost function } J = \frac{1}{2m} (e_1^2 + e_2^2 + e_3^2 + \dots + e_m^2)$$

$$= \frac{1}{2m} \sum (y_i - \hat{y})^2 \quad \text{where } m = \text{no of example}$$

$$= \frac{1}{2m} \sum (a_0 + a_1 x_i - \hat{y}_{\text{act}})^2 \quad \left\{ \begin{array}{l} y_i = a_0 + a_1 x_i \\ \hat{y} = \text{predicted} \end{array} \right\},$$

Clearly, cost function (J) is a function of parameter space
 $a = (a_0, a_1)$.

→ The best fitting model would be the one which minimizes any error metric which is cost function.

also { Lesser the distance better is our model }.
between straight line and data points.

→ How to minimize cost function?

→ As we have seen cost function (J) is a function of parameter space $a = (a_0, a_1)$;

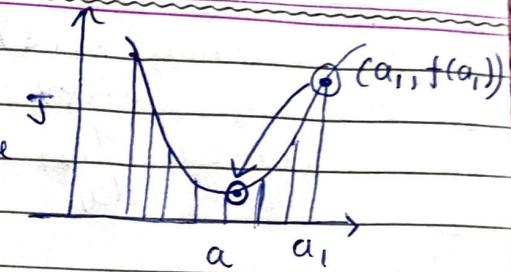
we will keep changing a_0 and a_1 , till we find a combination where cost function is minimized.

→ for this Gradient Descent Algorithm is used.

Gradient Descent Algorithm:

Assume a function $J = f(a)$

This curve represents J
The values of J will assume
for different values of ' a '.

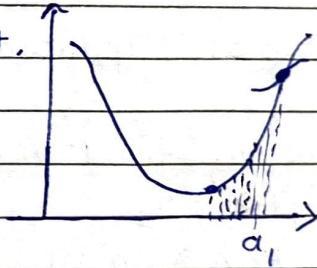


We are at $(a_1, f(a_1))$ as of now, and we want to reach at the minimum of this function.

→ How can we reach the minimum starting from a_1 ?

For this we will find slope at this point.

$$\rightarrow \frac{dJ}{da} \text{ at } a_1 \text{ or } f'(a_1).$$



Let's move a step ' α ' in this direction
to reach $a_1 - \alpha * f'(a_1)$

→ ' α ' is a small fixed quantity in the range of 0.01. Therefore the size of our steps is dependent on the slope $f'(a_1)$.

→ Higher the value of slope (which occurs away from minima)
larger the steps and vice-versa.

→ As we move closer to minimum, the slope decreases and hence our steps towards minimum keeps on getting smaller.
at minimum, slope becomes zero.

→ We perform multiple iteration of this Gradient Descent Algorithm.
The number of iterations performed depends upon

→ nature of function

→ ' α ' (learning rate)

→ and where we start from (a_1) in this case.

Step 1: Calculate $\frac{d(J)}{da_0}$ at current value of parameter a_0

calculate $\frac{d(J)}{da_1}$ at current value of parameter a_1

Step 2:

$$(\text{new}) a_0 = a_0 - \alpha \left(\frac{dJ}{da_0} \right)$$

$$(\text{new}) a_1 = a_1 - \alpha \left(\frac{dJ}{da_1} \right)$$

Step 3: update cost function J with new (a_0 and a_1)

~~Repeat step 1~~

One important thing about gradient descent is learning rate (α):

Keep α small, its value is kept around 0.01 so that it neither makes algorithm too slow nor does it fail to converge.

→ Learning Rate (α) is a hyperparameter for this model even though it doesn't directly affect model like parameters a_0 and a_1 do but it can impact the performance of our model and we need to be cautious about choosing model hyperparameters.

Logistic Regression: (Overview)

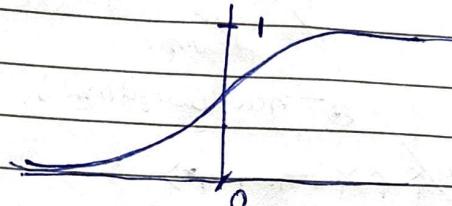
We often use linear regression to generate predictions for (numerical continuous) data eg: height, weight etc.

But what if data can-not be expressed in number eg: color, gender or job role etc.

→ Categorical data consists of discrete values that describe an object. Logistic Regression is a great candidate for ML problems that involve binary classification, which determines the P(E) that each data entry can be classified into one of two different categories.

Logistic regression uses the sigmoid function, also known as the logistic function, to map a linear combination of the input features into a range of 0 to 1.

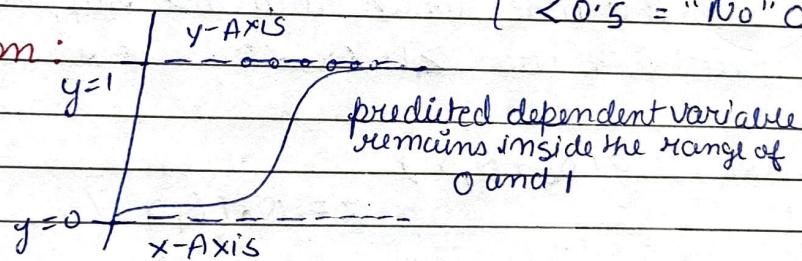
$$f(x) = \frac{1}{1 + e^{-x}}$$



→ We use (Decision Boundary) to determine at which point we will assign each data entry to each category. if (DB) is 0.5

$$\begin{cases} > 0.5 = 1 \\ < 0.5 = \text{"No"} 0. \end{cases}$$

Logistic Regression:



- Technique for classification.
- probabilistic view of classification.
- dependent variable is binary rather than continuous and it can also be applied to ordered categories (ordinal data).

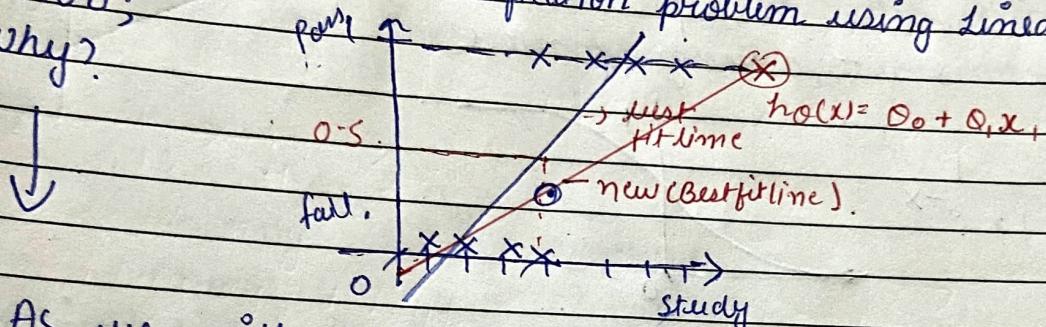
$$(MLE) - \ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 X$$

The term odds. (odds of winning or affecting)
In logistic regression, the dependent variable is a logit, which is the natural log of the odds, that is,

$$\text{odds} = \frac{P}{1-P} \quad \log(\text{odds}) = \text{logit}(P) = \ln\left(\frac{P}{1-P}\right)$$

Logistic Regression Indepth Maths :

→ Can we solve classification problem using Linear Regression
No;
Why?

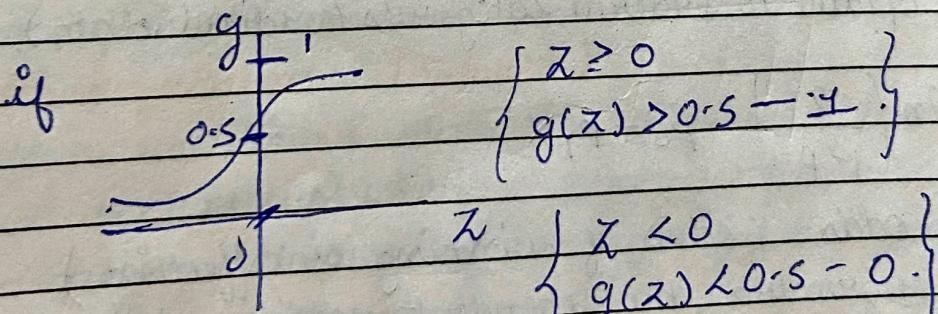


- ① As we will detect any (outlier) in the graph the best fit line will also change and this will lead to mis-value reading, Also in linear regression there will be a straight line which will not support binary output.
- ② output can be > 1 or < 0 also in some case.

$$h_0(x) = g(\theta_0 + \theta_1 x)$$

$$\left\{ g = \frac{1}{1+e^{-z}} ; z = \theta_0 + \theta_1 x \right\}$$

$$\left\{ h_0(x) = \frac{1}{1+e^{-(\theta_0 + \theta_1 x)}} \right\} \quad \begin{array}{l} \text{- hypothesis in} \\ \text{case of logit.} \end{array}$$



(Cost function of Logistic Regression:

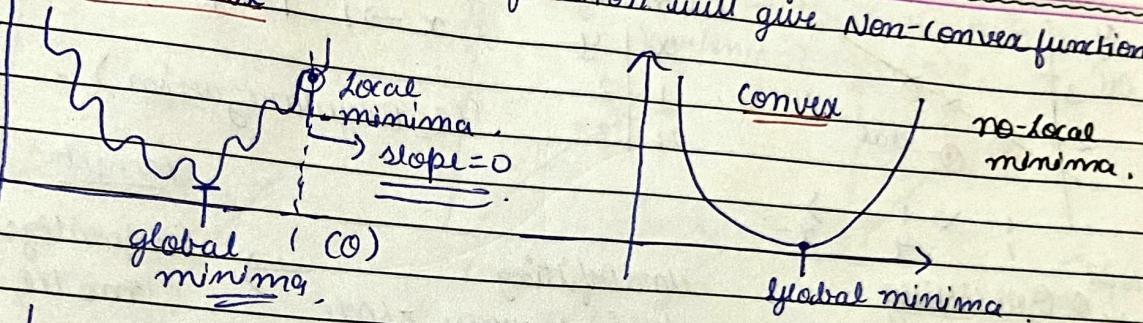
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum (h_0(x)^i - y^i)^2$$

() $\rightarrow h_0(x) = \frac{1}{1+e^{-(\theta_0 + \theta_1 x)}}$

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

In Logistic Regression cost function will give Non-convex function

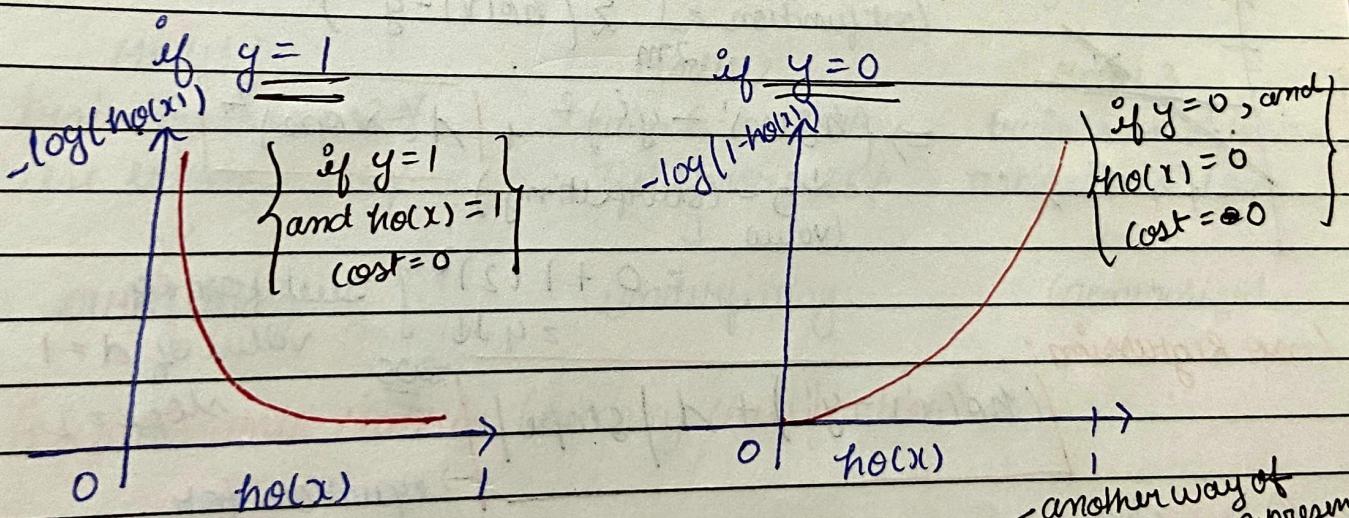
(Q) non-convex



→ to make the function convex

$$\text{cost } J(\theta_0(x^{(i)}), y) = \begin{cases} -\log(\theta_0(x)) & \text{if } y=1 \\ -\log(1-\theta_0(x)) & \text{if } y=0 \end{cases}$$

y ki value will always be 0 or 1.

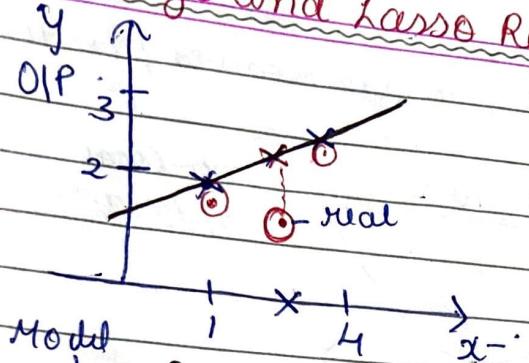


$$\text{cost } J(\theta_0(x), y) = -y \log(\theta_0(x)) - (1-y) \log(1-\theta_0(x))$$

Convergence Repeat:

$$\left\{ \theta_j : \approx \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \right\} \rightarrow \text{tells global minima.}$$

Ridge and Lasso Regression (Overview)



\downarrow Overfitting.

train Accuracy = 90%
test Accuracy = 70%

Low Bias
High Variance

Training set \Rightarrow Linear Regression

x	y
1	2
4	3

(λ_2 Regularization)

Underfitting

train Accuracy = 60%
test Accuracy = 62%

Low Bias and
Low Variance

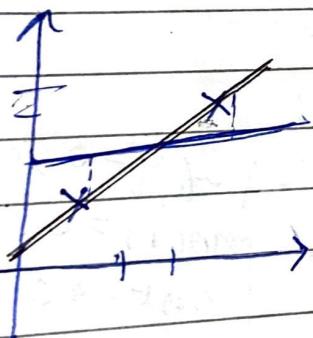
target

Generalized
model

needs

Low Bias and
Low Variance

To prevent overfitting we use Ridge Regression:



$$\text{cost function} = \frac{1}{2m} \sum (h_0(x)^T - y^T)^2$$

$$\Rightarrow (h_0(x)^T - y^T)^2 + \lambda (\text{slope})^2$$

(value)

$$= 0 + 1/2 \cdot 1^2 \quad \begin{cases} \text{Just for eg} \\ \text{value of } \lambda = 1 \end{cases}$$

$$= 4 \cdot 1/2 \cdot 1^2 = 2 \cdot 1 = 2$$

slope = 2.

$$(h_0(x)^T - y^T) + \lambda |\text{slope}|$$

\rightarrow equation for
Lasso.

Lasso Regression:

- \rightarrow ① Overfitting prevented
- \rightarrow ② Helps in feature selection.

(λ_1 Regression Regularization)

Bias and Variance: (Overview)

Regression Problem:

y

x

underfitting
($TE \uparrow$)

High Bias
↓
High variance
↓
Low variance
↓
Both can
be possible.

y

generalized.
 $TE(\downarrow)$

Low Bias
↓
Low variance

y

overfitting
($TE(\downarrow) \approx 0$)

Low Bias
↓
High variance

Bias \rightarrow Training Data
variance \rightarrow Test data.

Classification Problem: = {Confusion Matrix}.

Model 1

training error 2%.

test error = 25%.

↓

overfitting.

Model 2

train E \approx 25%

Test Error \approx 26%

↓

underfitting

Model 3

trainError $< 10\%$

test error $< 10\%$.

↓

generalized.

Performance Matrix (Binary classification).