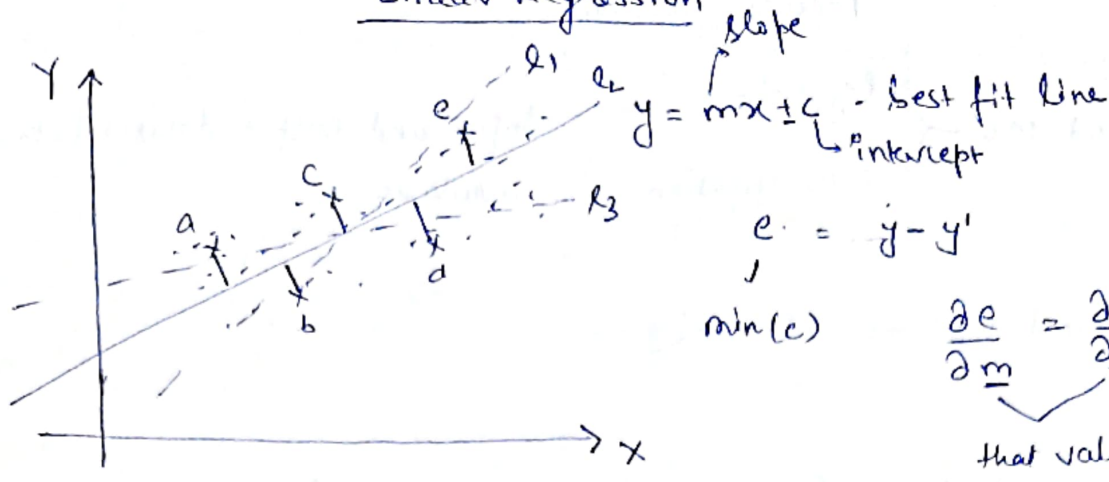


# Linear Regression

2



$$\frac{\partial e}{\partial m} = \frac{\partial e}{\partial c} = 0$$

that value of  $m, c$ .

We will choose a line which is at a min distance from points a, b, c, d, e.  $(d_a + d_b + d_c + d_d + d_e) \rightarrow \text{minimum}$

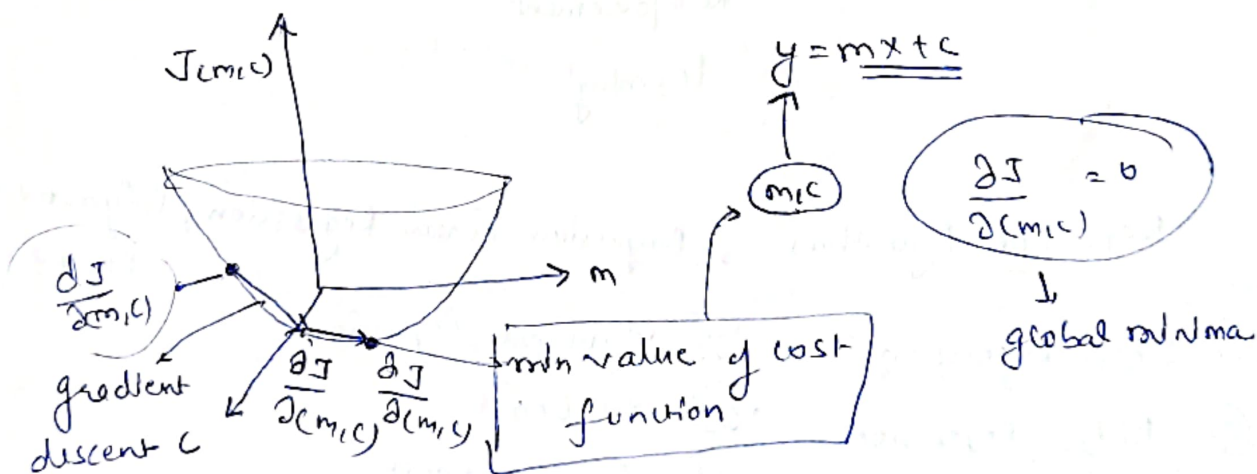
Cost function / error fn  $J(m, c) = \frac{1}{2N} \sum_{i=1}^N (\hat{y} - y)^2$

$\hat{y} = mx + c$

$\hat{y}$  Real Value

minimize this error

Number of all the points



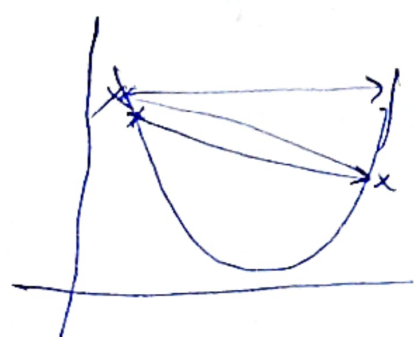
Convergence theorem

$m - \left( \frac{\partial m}{\partial m} \right) \times \alpha$

slope

learning Rate

small

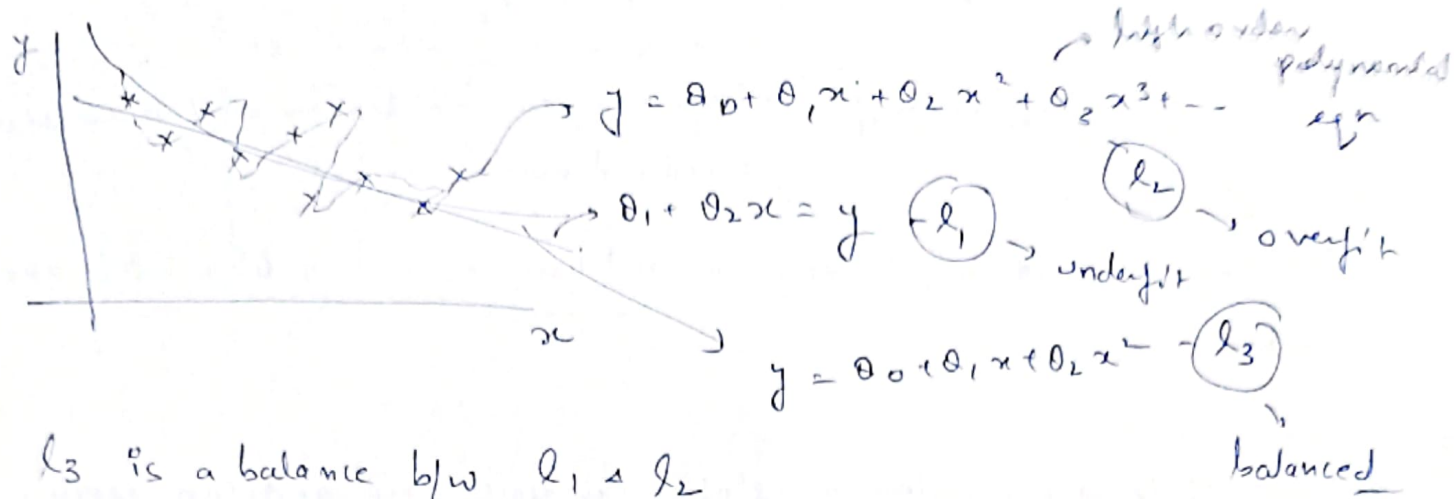


Converging gradient problem

if  $\alpha \uparrow \uparrow$

$\alpha \downarrow \downarrow$   
 slow

# Lasso | Ridge | Elasticnet | L1 L2 Regularization



In order to obtain a balanced fit we need regularisation.

$$\text{mean sq. error} = \frac{1}{n} \sum_{i=1}^n (y_i - y_{\text{pred}})^2$$

$$y_{\text{pred}} = h_{\theta}(x_i)$$

$$\text{mse} = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$

$$+ \lambda \left( \sum_{i=1}^n \theta_i^2 \right)$$

Ridge  
L2 Regularisation  
b/c we are using square

if  $\lambda \uparrow \rightarrow$  horizontal line

$$\lambda (\text{slope})^2$$

This lambda can help define the limits for  $\theta$ .

By adding this value we are trying to ensure that the model gets penalised whenever we have a higher value of  $\theta$  in our predictions.

It tends to reduce the slope steepness. with unit increase in  $x$  there is high increase in  $y$ .

$$\text{mse} = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{i=1}^n |\theta_i|$$

$\rightarrow$  L1 Regularisation  
no square  
Lasso

\* overfitting - high variance, high bias

\* balanced model - low bias and low variance

helps us to do feature selection

(4)

$$y = m_1 x_1 + m_2 x_2 + m_3 x_3 + m_4 x_4 + c$$

$\lambda (m_1 + m_2 + m_3 + m_4)$  → unimportant or less feature will automatically get reduced near to zero.

Let's say fit with  $m_4$  is very less and here  $x_4$  will also be a useless feature.

In case of lasso regression a slope can reach zero and can render a feature completely useless and in case of ridge the slope may reduce but will never get to zero. Elastic net uses both lasso and ridge together to include the impact & benefits of both the varieties.

\* Lasso works best when we have a lot of useless features, as it eliminates all the insignificant ones and keep only the most significant ones.

\* Ridge regression works best when we have a lot of useful variables i.e. when all the variables are significant ones.

\* When we have a model that include millions of parameters, and we don't know whether which variables are useful and which are not.

$$mse = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y})^2 + \lambda_1 \sum | \theta_i | + \lambda_2 \sum (\theta_i^2)$$

elastic net regression

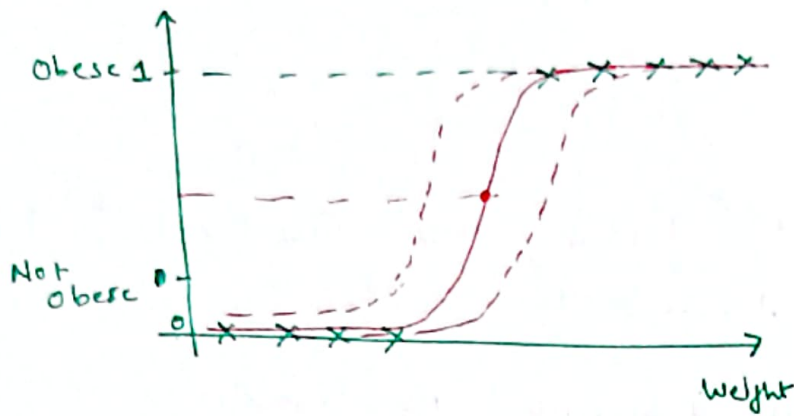
$$\lambda_1 = 0 \quad \lambda_2 \neq 0 \quad \text{Lasso}$$

$$\lambda_1 \neq 0 \quad \lambda_2 = 0 \quad \text{Ridge}$$

$$\lambda_1 = 0 \quad \lambda_2 = 0 \quad \text{mse}$$

↓  
Very good when there exist a lot of correlated parameters.

## \* Logistic Regression



### Binary classification

\* We can also do this classification through linear regression but that will have a high error rate and is prone to outliers. This may not work effectively. Hence sigmoid function is used for classification.

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

\* Logistic regression is applied to a classification problem that is linearly separable, can be divided with the help of straight line.

Cost function =  $\max \sum_{i=1}^n y_i (w_i^T x_i)$    
 Classification by linear regression   
 only parameter to be updated   
 → if I want a linearly separable line   
 +1/-1 distance of point from plane   
 the point   
 If an outlier exists then this will not be possible

Sigmoid fn → logistic (Maximum Likelihood)

$$\text{Cost fn} = \max \sum_{i=1}^n f(y_i w_i^T x_i)$$

$$f(x) = \frac{1}{1 + e^{-x}} \cdot (0-1) \rightarrow \text{removes the effect of outliers in linearly separable algorithms.}$$

- No multicollinearity in data
  - Binary & dichotomous classification
  - Linearly separable data - log odds.
  - Large sample
- logistic regression

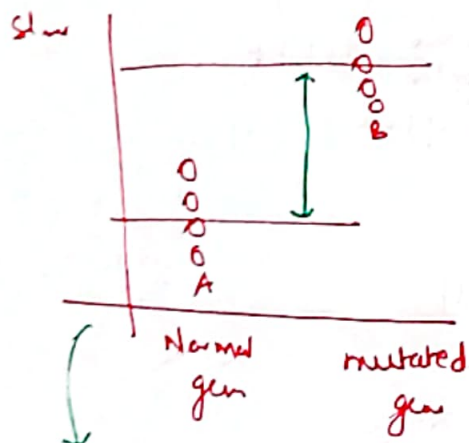
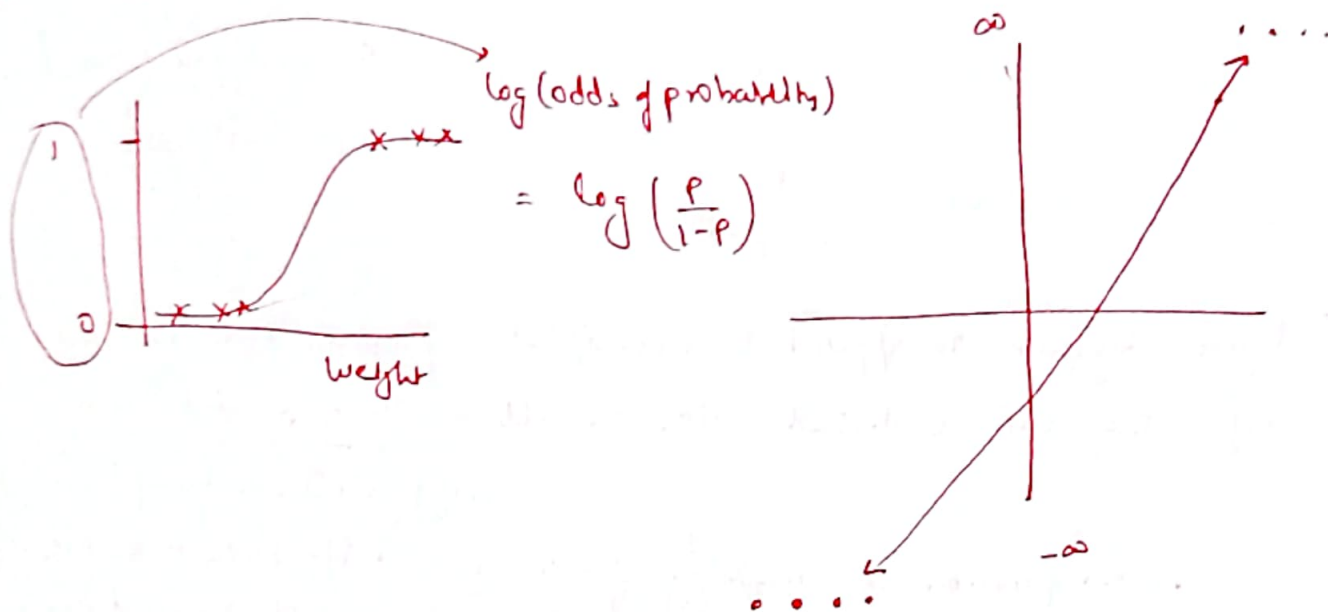


9

~~allison 19-04-15~~

weight is the continuous variable predicting outcome. There can be a categorical variable leading to the prediction as well.

logistic regression is one type of generalised linear model.



$$\text{size} = \text{mean normal}_A \times B_1 + (\text{mean}_B - \text{mean}_A) \times B_2$$

for logistic regression - transform this to a log scale

$$\log(\text{odds gene normal}) = \log\left(\frac{2}{9}\right)$$

$$\log(\text{odds gene mutated}) = \log\left(\frac{7}{3}\right)$$

$$\Rightarrow \text{size} = \log(\text{odds gene normal}) \times B_1$$

$$+ (\log(\text{odds gene mutated}) - \log(\text{odds gene normal})) \times B_2$$

