

4

a) Expressions

① Gradient:

Let our model be

$$P_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}}$$

Loss function for Logistic Regression is

$$L(\beta) = \sum_{i=1}^N \left[y_i \log(P_i) + (1 - y_i) \log(1 - P_i) \right]$$

where:

N = No of samples

P_i = Probability that i^{th} sample having $Y = 1$

\downarrow
Predicted probability

y_i = observed outcome of i^{th} sample.

Now gradient will be $\frac{\partial L(\beta)}{\partial \beta_j}$

$$\Rightarrow \frac{\partial L(\beta)}{\partial \beta_j} = \sum_{i=1}^N \left(\frac{y_i}{P_i} - \frac{y_i - y_i}{(1 - P_i)} \right) \times \frac{\partial P_i}{\partial \beta_j}$$

$$\frac{\partial p_i}{\partial \beta_j} = p_i(1-p_i) x_{ij}$$

$$\Rightarrow \text{Gradient} \quad \frac{\partial L(\beta)}{\partial \beta_j} = \sum_{i=1}^N (y_i(1-p_i) - (1-y_i)p_i) x_{ij}$$

$$\text{Gradient} = \sum_{i=1}^N (y_i - p_i) x_{ij}$$

② Hessian:

$$H_{ij} = - \frac{\partial^2 L(\beta)}{\partial \beta_i \partial \beta_j} = - \frac{\partial}{\partial \beta_i} \left[\frac{\partial L(\beta)}{\partial \beta_j} \right]$$

(a)

$$H_{jk} = - \frac{\partial^2 L(\beta)}{\partial \beta_j \partial \beta_k} = - \frac{\partial}{\partial \beta_j} \left[\frac{\partial L(\beta)}{\partial \beta_k} \right]$$

$$= - \frac{\partial}{\partial \beta_j} \left[\sum_{i=1}^N (y_i - p_i) x_{ik} \right]$$

$$= \sum_{i=1}^N (-1) x_{ik} \frac{\partial p_i}{\partial \beta_j}$$

$$H_{jk} = \sum_{i=1}^N p_i(1-p_i) x_{ij} x_{ik}$$

⇒ Hessian Matrix is (Size $M \times M$) ($M = \text{No of parameters}$)

$$\begin{bmatrix} \sum_{i=1}^N p_i(1-p_i) x_{i1}^2 & \sum_{i=1}^N p_i(1-p_i) x_{i1} x_{i2} & \dots & \sum_{i=1}^N p_i(1-p_i) x_{i1} x_{iM} \\ \sum_{i=1}^N p_i(1-p_i) x_{i2} x_{i1} & \sum_{i=1}^N p_i(1-p_i) x_{i2}^2 & \dots & \sum_{i=1}^N p_i(1-p_i) x_{i2} x_{iM} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^N p_i(1-p_i) x_{iM} x_{i1} & \sum_{i=1}^N p_i(1-p_i) x_{iM} x_{i2} & \dots & \sum_{i=1}^N p_i(1-p_i) x_{iM}^2 \end{bmatrix}$$

③ update rules

$$(B_i)_{n+1} = (B_i)_n - \frac{\partial(L(B))}{\partial B_i}$$

should do this
for all B_i ($i=1$ to m)

Let 'B' be the vector of all co-efficients

For ~~log~~ gradient descent

$B_{n+1} \rightarrow$ updated for n^{th} time

$B_n \rightarrow$ updated for ~~(n-1)~~ $(n-1)^{\text{th}}$ time.

$$B_{n+1} = B_n - [H^{-1} \cdot \nabla L(B_n)]$$

$H^{-1} \rightarrow$ Inverse of Hessian Matrix

$\nabla L(B_n) \rightarrow$ Gradient vector.

All the coefficients
will be updated

Newton-Raphson
method

④ Algorithm :

Let $\epsilon = 10^{-6}$

while ($B_{n+1} - B_n > \epsilon$):

$$B_{n+1} = B_n - [H^{-1} \cdot \nabla L(B_n)]$$

More Mathematical

Let's see a more code related algo

```
import numpy as np
```

```
def sigmoid(z):
```

```
(Returns probability, i.e  $1/(1 + \exp(-z))$ )
```

```
def gradient(X, y, beta):
```

```
(Returns the gradient updated value, i.e  $B_{n+1}$ )
```

```
def hessian(X, beta):
```

```
(Returns the Hessian matrix)
```



```
def newton_raphson(X, y):
```

```
    beta = np.zeros(N)
```

```
    while True:
```

```
        gr = gradient(X, y, beta)
```

```
        he = hessian(X, beta)
```

```
        beta_new = beta - np.linalg.solve(he, gr)
```

```
        if (beta_new - beta <  $10^{-6}$ ):
```

```
            Break.
```

```
        beta_new = beta
```

```
    Return beta_new
```

```
# X is the feature matrix
```

```
# y is the binary outcome vector
```