**Linear regression:-**

$$\hat{y} = h_\beta(x) = \sum_{j=0}^{N} \beta_i x_i$$

→ objective :- Minimize the sum of squared error (SSE)

→ Used for regression problem

**Classification:-**

Linear function + Another function ⟹ classification

$$(i) \; \hat{y} = g(h_\beta(x))$$

where) $g(z) = \dfrac{1}{1+e^{-z}}$

$$\hat{y} = h_\beta(x) \Rightarrow g(\beta^T x)$$

Value between 0 and 1

$$\Rightarrow \dfrac{1}{1+e^{-\beta^T x}}$$ (ii) sigmoidal function

$<0.5 \qquad \geq 0.5$

Note :- $g(z) = 1$ if $z = \infty$

$\qquad g(z) = 0$ if $z = -\infty$

$$g'(z) = \dfrac{1}{1+e^{-z}} \times e^{-z}$$

$$= \dfrac{1}{1+e^{-z}} \times \left(1 - \dfrac{1}{1+e^{-z}}\right)$$

$$= g(z) * (1-g(z)) \Bigg\} \text{——①}$$

(or)

$$h_\beta(x) \cdot (1 - h_\beta(x))$$

* A single trial
* $P(\text{success}) = p$
  $P(\text{failure}) = 1-p$.

Likelihood of a single observation for $p$ given $x$ and $y$:

$$L(p_i / y_i) = p(y_i = y_i) = p_i^{y_i} (1-p_i)^{1-y_i}$$

Given the observations are independent,

$$L(p/y) = \prod_i p(y_i = y_i) = \prod_i p_i^{y_i} (1-p_i)^{(1-y_i)}$$

Log likelihood turns product into sums.

$$\ell(p/y) = \sum_i y_i \log(p_i) + (1-y_i) \log(1-p_i)$$

$$= \sum_i y_i \log(\hat{y}) + (1-y_i) \log(1-\hat{y})$$

Take derivative of $\ell(p/y)$ to maximize log likelihood

$$\text{New } \beta = \text{old or initial } \beta + \alpha \nabla_\beta \ell(p/y)$$

$$\frac{\partial \ell(p/y)}{\partial \beta} = \left[ y(-\frac{1}{\hat{y}}) + (1-y) x \frac{-1}{1-\hat{y}} \right] * \frac{\partial}{\partial \beta}(\hat{y})$$

$$= \left[ \frac{\dot{y}}{g(\beta^T x)} + (1-y) \times \frac{-1}{1-g(\beta^T x)} \right] \frac{-\partial}{\partial \beta} g(\beta^T x)$$

$$= \left[ \frac{y}{g(B^T x)} - (1-y) \times \frac{1}{1-g(B^T x)} \right] \times g(B^T x) \, (1-g(B^T x)) \, \frac{d}{dB} (B^T x)$$

$$\text{②}$$

$$= \frac{\left[ y \, (1-g(B^T x)) \right] - \left[ (1-y) \times g(B^T x) \right]}{g(B^T x) \, (1-g(B^T x))} \times g(B^T x) \, (1-g(B^T x)) * x$$

$$= \left[ y - yg(B^T x) - g(B^T x) + y \, g(B^T x) \right] * x$$

$$= \left[ y - g(B^T x) \right] * x$$

$$= \left( y - h_\beta(x) \right) * x \longrightarrow \text{②}$$

$$\boxed{\therefore \quad \beta_j = \beta_j + \alpha * y - h_\beta(x) * x_j} \longrightarrow \text{③}$$

# Stochastic Gradient Descent :-

1. Initialize $\beta$

2. For randomly selected sample

   $\Rightarrow h_\beta(x) \Rightarrow$ prediction for $x_i$ using the current i/p

   $\Rightarrow$ for each non zero feature of $x_{ij}$

   $$\beta_j = \beta_j * \alpha(y^{(i)} - h_\beta(x^{(i)})) * x_j^{(i)}$$

Ex:- $\beta = (0, 0, 0)$    $\alpha = 0.1$    $x_1 = ((62, 58), 1)$

$\qquad\qquad\qquad\qquad\qquad\qquad x_2 = ((52, 41), 0)$

$h_\beta(x_1) = \dfrac{1}{1 + e^{-0}} = 0.5$

Use the score to fine tune the parameter

$\beta_0 = 0 + 0.1[(1 - 0.5) * 1] = 0 + 0.1 * 0.5 = 0.05$

$\beta_1 = 0 + 0.1(1 - 0.5) * 62 = 0 + 0.05 \times 62 = 3.1$

$\beta_2 = 0 + 0.1(1 - 0.5) * 58 = 0 + 0.05 \times 58 = 2.9$

$$h_\beta(x_2) = \dfrac{1}{1 + e^{-(0.05 * 3.1 * 52 + 2.9 * 41)}}$$

$$= \dfrac{1}{1 + e^{-280.15}} = 1$$

Expected value is zero, fine tune the parameters.

$\beta_0 = 0.05 + 0.1[0 - 1] * 1 = -0.05$

$\beta_1 = 3.1 + 0.1[0 - 1] * 52 = -2.1$

$\beta_2 = 2.9 + 0.1[0 - 1] * 41 = -1.2$

$x_1 \Rightarrow h_\beta(x) = \dfrac{1}{1+e^{-(-0.05-2.1*62-1.2*58)}} = 1.6 \times 10^{-87}$ ③

$\Rightarrow$ class $0$

$x_2 = h_\beta(x) = \dfrac{1}{1+e^{-(0.05-2.1*52-1.2*41)}}$

$= 1.53 \times 10^{-69}$

$= $ class $0$

Accuracy $= \begin{array}{c} 0.5 \Rightarrow 50\% \\ \end{array}$

(e)

$\frac{1}{2}$.

Gradient Descent (GD) :-
————X————

$$\beta_j = \beta_j^{\circ} - \alpha \dfrac{d}{d\beta_j} \, \ell(P/y)$$

1. Start with random

2. Loop until convergence

    2.1 compute gradient

    2.2 Update.

3. Return.

Stochastic GD :-
————X————

1. Start with random

2. Loop until convergence

    2.1. pick the single data point i

    2.2 compute gradient over that single point

    2.3 Update

3. Return.

| gradient descent | stochastic GD |
|---|---|
| – small or medium dataset | – Large dataset |
| – slow in computation. | – Fast in computation |

Newton Raphson:—

$$\beta_1 = \beta_1 - \frac{\partial^2}{\partial \beta^2} l(P/y) \Rightarrow \beta_1 - \frac{\partial}{\partial \beta} l(P/y) *$$

$\longrightarrow$ Reduces number of iterations.

$$\left(\frac{\partial^2}{\partial \beta^2} l(P/y)\right)^{-1}$$

$$\frac{\partial}{\partial \beta} l(P/y) = (y - h_\beta(x)) * x = yx - xh_\beta(x) = yx - x * \frac{1}{1+e^{-\beta x}}$$

$$\frac{\partial^2}{\partial \beta^2} l(P/y) = 0 - x(h_\beta(x)(1-h_\beta(x)) * x)$$

Hessian $\Leftarrow H$ $= -h_\beta(x)(1-h_\beta(x)) * x^2$

$$\beta_1 = \beta_1 - \left(h_\beta(x^{(i)})(1-h_\beta(x^{(i)}) x^{(i)}(x^{(i)})^T\right)^{-1} \frac{\partial}{\partial \beta}(l(P/y)$$

(c)

$$\beta_1 = \beta_1 - H^{-1} \cdot \frac{\partial}{\partial \beta} l(P/y)$$