

05-02-2022

Saturday, February 5, 2022 2:04 PM

Inverse logit func. or Sigmoid()

$$\sigma^z = \frac{e^{A_0 + B^T z_i}}{1 + e^{A_0 + B^T z_i}} \quad \text{or} \quad \frac{1}{1 + e^{-(A_0 + B^T z_i)}}$$

$$\sigma^z = \frac{1}{1 + e^{-B^T x}} \quad //$$

$$\sigma(B^T z) = \frac{1}{1 + e^{-B^T z}}$$

Mom's notation

$$\begin{aligned} \sigma(z) &= \frac{1}{1 + e^{-z}} \\ \frac{\partial (\sigma(z))}{\partial z} \frac{\partial (\sigma(z))}{\partial z} &= - \frac{1}{(1 + e^{-z})^2} \cdot e^{-z} \cdot -1 \\ &= \frac{e^{-z}}{(1 + e^{-z})^2} \\ &= \frac{e^{-z}}{1 + e^{-z}} \cdot \frac{1}{1 + e^{-z}} \\ &\approx \frac{1 + e^{-z} - 1}{1 + e^{-z}} \cdot \frac{1}{1 + e^{-z}} \\ &= \left[ 1 - \frac{1}{1 + e^{-z}} \right] \left( \frac{1}{1 + e^{-z}} \right) \end{aligned}$$

$$\frac{\partial(\sigma(z))}{\partial z} = [1 - \sigma(z)] \sigma'(z) \quad \text{--- A}$$

Maximum Likelihood Estimate of  $\beta$  /  $\hat{\beta}_{MLE}$   
 - Best  $P$  value

Minimizes Prob. of seeing training data

$\hat{\beta}_{MLE} \Rightarrow$  likelihood of seeing some training data.

$$y: \begin{matrix} y \\ 0 \\ 0 \end{matrix} \quad \begin{matrix} \hat{y} \\ 0 \\ 0 \end{matrix} \quad p_1$$

$$| \quad | \quad p_2$$

| |  $p_3 \rightarrow$  Prob.  $> 0.5$  (as high as possible)

| |  $p_4 \rightarrow$  Prob.  $< 0.5$  (as low as possible)

$$| \quad | \quad p_5$$

(Joint prob.)  $= (1-p_1) p_2 \cdot p_3 (1-p_4) p_5$

$$\hat{\beta}_{MLE} = \underset{\beta}{\operatorname{argmax}} \quad P(y|n; \beta)$$

$$= \underset{\beta}{\operatorname{argmax}} \prod_{i=1}^n P(y_i|n_i; \beta)$$

$$L(\beta) = \underset{\beta}{\operatorname{max}} \prod_{i=1}^n P(y_i|n_i; \beta)$$

$$= \underset{\beta}{\operatorname{max}} \prod_{i=1}^n P(y_i=1) \cdot \prod_{i=1}^n P(y_i=0)$$

$$= \underset{\beta}{\operatorname{max}} \prod_{i=1}^n P(y_i=1) \cdot \prod_{i=1}^n (1 - P(y_i=1))$$

$$1 + \sum_{i=1}^n P(y_i=1) - \frac{1}{n}$$

$$\text{Let } \underline{s_i} P(y_i=1) = \frac{1}{1+e^{-\beta x_i}}$$

$$= \max_{\beta} \prod_{i=1}^n s_i \cdot \prod_{i=1}^n (1-s_i)$$

Bernoulli  $\because f(n) = p^n (1-p)^{1-n}$  ] explanation

$$L(\beta) = \max_{\beta} \prod_{i=1}^n s^{y_i} (1-s)^{1-y_i}$$

Log Likelihood ( $L$ )

$$\begin{aligned} \log(L(\beta)) &= LL(\beta) = \log \left( \max_{\beta} \prod_{i=1}^n s^{y_i} (1-s)^{1-y_i} \right) \\ &= \max_{\beta} \sum_{i=1}^n (\log(s)^{y_i} + \log(1-s)^{1-y_i}) \end{aligned}$$

$$LL(\beta) = \max_{\beta} \sum_{i=1}^n (y_i \log(s) + (1-y_i) \log(1-s))$$

Cost function (or) Loss function : Both have slight diff.

Minimise Cost

Minimise Likelihood

$$NLL(\beta) = -\min_{\beta} \sum_{i=1}^n (y_i \log s + (1-y_i) \log(1-s))$$

$$\hat{\beta}_{MLC} = \text{argmax}(LL(\beta)) = \text{argmin}(NLL(\beta))$$

To calc.  $\beta$  val.  $\rightarrow$  we use Optimizers.

Jacobian Nchm  
less accuracy

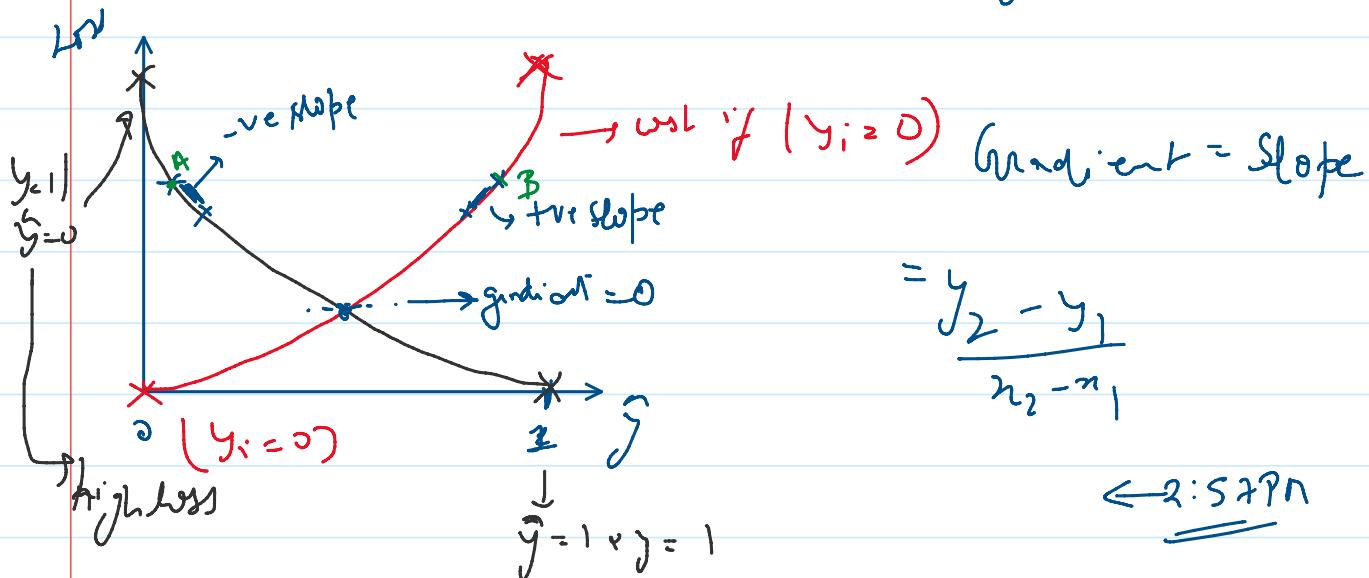
To calc  $\beta$  val  $\rightarrow$  we use Optimizers.

- $\Rightarrow$  1<sup>st</sup> Order - eg: Gradient Descent (Less costly)
- $\Rightarrow$  2<sup>nd</sup> Order - eg: Newton-Raphson (High Cost)  $\downarrow$

Jacobian Matrix  
Hessian Matrix  
Brent Accuracy

Gradient descend - minimize Likelihood func.

1) Descent - Minimise cost / loss func.



{ GD  $\Rightarrow$  in each iteration we take steps proportional to the -ve of the gradient at that pt. }

$$\left\{ \beta = \beta - \alpha \frac{\partial \text{NLL}(\beta)}{\partial \beta} \right\}$$

$\hookrightarrow$  Learning Rate: how fast should model learn

From above graph

$$\left\{ \alpha = 0.1 \text{ to } 0.3 \right\}$$

i) Let  $\beta$  be diff 'A' Gradient = -ve

$$\text{Loss} \propto \frac{1}{\beta} \quad \therefore \beta = \beta - \alpha (-\text{ve grad.}) \\ = \beta + \alpha \text{ grad.} \\ \therefore \uparrow \beta$$

$$\beta \quad = \beta + \alpha \text{ grad.}$$

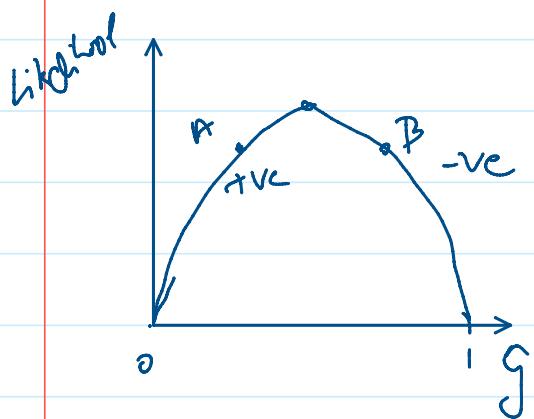
$\therefore \uparrow \beta$

ii) Let  $\beta$  be at 'B' Gradient = +ve

$$\begin{aligned}\beta &= \beta - \alpha (+\text{ve grad.}) \\ &= \beta - \alpha \text{ grad} \\ &= \beta \downarrow\end{aligned}$$


---

### Gradient Ascent



$$\beta = \beta + \alpha \frac{\partial}{\partial \beta^T} LL(\beta)$$

$$\beta = \beta + \alpha (\text{Stp}) \quad \textcircled{A}$$

$$\begin{aligned}\beta &= \beta + \alpha (-\text{ve Stp}) \\ &= \beta - \alpha (\text{Stp})\end{aligned} \quad \textcircled{B}$$

$$LL(\beta) = \sum_{i=1}^n \left[ y_i \underbrace{\log(\sigma(\beta^T n_i))}_{B} + (1-y_i) \underbrace{\log(1-\sigma(\beta^T n_i))}_{C} \right]$$

$$\frac{\partial LL(\beta)}{\partial \beta^T} = \frac{\partial}{\partial \beta^T} (B) + \frac{\partial}{\partial \beta^T} (C)$$

For B :  $\frac{\partial}{\partial \beta^T} \left( \sum_{i=1}^n y_i \log(\sigma(\beta^T n_i)) \right)$  Chain Rule

$$\Rightarrow \sum y_i \frac{\partial \log(\sigma(\beta^T n_i))}{\partial (\sigma(\beta^T n_i))} \cdot \frac{\partial \sigma(\beta^T n_i)}{\partial \beta^T n_i} \cdot \frac{\partial \beta^T n_i}{\partial \beta^T}$$

From A

$$\frac{\partial L(\beta)}{\partial \beta^T} = \sum_i y_i \cdot \frac{1}{\sigma(\beta^T n_i)} \cdot \cancel{\sigma(\beta^T n_i)} \cdot (1 - \sigma(\beta^T n_i)) \cdot n_i$$

$$\Rightarrow \sum_i y_i (1 - \sigma(\beta^T n_i)) n_i$$

For C:

$$\frac{\partial}{\partial \beta^T} \left( \sum_i (1-y_i) \log (1 - \sigma(\beta^T n_i)) \right)$$

$$= \sum_i (1-y_i) \frac{\partial \log (1 - \sigma(\beta^T n_i))}{\partial (1 - \sigma(\beta^T n_i))} \cdot \frac{\partial (1 - \sigma(\beta^T n_i))}{\partial (\sigma(\beta^T n_i))} \cdot \frac{\partial \sigma(\beta^T n_i)}{\partial \beta^T}$$

$$= \sum_i (1-y_i) \frac{1}{1 - \sigma(\beta^T n_i)} (-\sigma(\beta^T n_i)(1 - \sigma(\beta^T n_i))) \cdot n_i$$

$$= - \sum_i (1-y_i) \sigma(\beta^T n_i) n_i$$

Combine B and C

$$\Rightarrow \sum_i [y_i (1 - \sigma(\beta^T n_i)) n_i - (1-y_i) \sigma(\beta^T n_i) n_i]$$

$$\Rightarrow \sum_i (y_i - y_i \sigma(\beta^T n_i) - \sigma(\beta^T n_i) + y_i \cancel{\sigma(\beta^T n_i)}) n_i$$

$$\frac{\partial L(\beta)}{\partial \beta^T} = \sum_{i=1}^n [(y_i - \sigma(\beta^T n_i)) n_i]$$

Gradient Ascent

$$\left\{ \beta = \beta + \alpha \left( \sum_{i=1}^n (y_i - \sigma(\beta^T n_i)) n_i \right) \right\} *$$

⋮ - 1 -

✓

- ↗

$$\hat{y} = \frac{1}{1 + e^{-px}}$$