

Question-4

Poisson regression

Q1) : ML estimation:

$$L = \prod_{i=1}^n P(y_i / \theta_i) \rightarrow \text{Likelihood.}$$

$\theta_i \rightarrow$ parameters for i^{th} data point

Here, we try to fit data to a poisson distribution, with different parameters.

$$P(y_i / \theta_i) = P(y_i / \lambda_i) = \frac{e^{-\lambda_i} (\lambda_i)^{y_i}}{(y_i)!}$$

Here ' λ ' \rightarrow parameter for poisson

But $\lambda_i \rightarrow$ scalar for i^{th} data point.

$$\begin{array}{ccc} & \beta \rightarrow \text{LT} & \\ X_i & \xrightarrow{\quad} & \lambda_i \\ d \times 1 & & 1 \times 1 \end{array}$$

} $d \rightarrow$ dimensions / features of a data point.

\Rightarrow We perform a linear transformation,
& estimate λ_i for a given X_i

$$\lambda_i = \exp(X_i \beta)$$

$\beta \rightarrow d \times 1$ column vector
 $X_i \rightarrow 1 \times d$ row vector

In order to fit data, we need to find parameter λ_i of distribution

⇒ MLE:-

↳ we use MLE to find λ_i (indirectly we need to find β)

$\beta \rightarrow ?$

$$L = \prod_{i=1}^N \frac{e^{-\lambda_i} (\lambda_i)^{y_i}}{(y_i)!}, \text{ where } \lambda_i = \exp(X_i \beta)$$

→ Goal is to maximise 'L'

⇒ log-likelihood (it does not change the loci)

$$\log_e(L) = \sum_{i=1}^N -\lambda_i + y_i \log_e(\lambda_i) - \log(y_i)!$$

Need to maximise $\log_e(L(\lambda_i))$

$$\Rightarrow \log(L(\beta)) = \sum_{i=1}^N -\lambda_i + y_i \log(e^{X_i \beta}) - \log(y_i)!$$

$$= \sum_{i=1}^N -e^{(\bar{X}_i \beta)} + y_i \bar{X}_i \beta - \log(y_i)!$$

* Let $J(\beta)$ be loss function,
Aim is to minimise J , to find ' β '

$$\boxed{J(\beta) = -\log(L(\beta))}$$

log loss function

$$\Rightarrow J(\beta) = \sum_{i=1}^N e^{(x_i^T \beta)} - y_i (x_i^T \beta) - \log(y_i)$$

→ Find β st $J(\beta)$ is min.

$$\Rightarrow \frac{\partial J(\beta)}{\partial \beta_k} = \sum_{i=1}^N e^{(x_i^T \beta)} \cdot x_i^k - y_i (x_i^k) \quad : k \rightarrow k\text{th element in } \beta$$

$$= \sum_{i=1}^N x_i^k (e^{(x_i^T \beta)} - y_i)$$

$$\frac{\partial J(\beta)}{\partial \beta} = X^T (e^{(X\beta)} - y)$$

$$X = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_N \end{bmatrix}_{N \times D} \Rightarrow \bar{x}_i = [x_i^1 \dots x_i^D]_{1 \times D}$$

$$\beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_D \end{bmatrix}_{D \times 1}, \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}$$

For minimising $J(\beta)$.

$$\Rightarrow \frac{\partial J(\beta)}{\partial \beta} = 0 \Rightarrow X^T (e^{(X\beta)} - y) = 0$$

\Downarrow

Need to Solve a non-linear eqn

→ Hence, we use gradient descent to solve it. Since no close form expression exists.

⇒ Gradient Descent

$$\beta_{n+1} = \beta_n - \eta \left(\frac{\partial J}{\partial \beta} \right)$$

$$\beta_{n+1} = \beta_n - \eta \left(X^T (e^{X\beta_n} - y) \right)$$

$\eta \rightarrow$ learning rate

→ When it converges

$$\lambda_i^o = \exp(X_i^o \beta)$$

$$\hat{y}_i^o = \hat{\lambda}_i^o$$

→ Prediction

$$\boxed{\hat{y}_i^o = \exp(X_i^o \beta)}$$

4.4 Regularisation:

⇒ As explained in 4.1:

cost func / loss fn (No regularisation).

$$J(\beta) = \sum_{i=1}^N e^{x_i^T \beta} - y_i (\bar{x}_i^T \bar{\beta}) - \log(y_i!)$$

$$\Rightarrow \beta_{n+1} = \beta_n - \eta \left(\frac{\partial J}{\partial \beta} \right)$$

For L_1 regularisation

* hyperparameter: λ
 $M \rightarrow$ no. of features

$$J(\beta) = \sum_{i=1}^N e^{x_i^T \beta} - y_i (\bar{x}_i^T \bar{\beta}) - \log(y_i!) + \lambda \sum_{j=1}^M |\beta_j|$$

$$\therefore \left. \frac{\partial J}{\partial \beta} \right|_{L_1 \text{ regularised}} = \left. \frac{\partial J}{\partial \beta} \right|_{\text{non regularised}} + \lambda \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{M \times 1}$$

⇒ Hence: for L_1 -regularisation

$$\beta_{n+1} = \beta_n - \eta \left(\left. \frac{\partial J}{\partial \beta} \right|_{\text{non-reg}} - \lambda \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}_{M \times 1} \right)$$

Use of L1-regularisation

- ↳ This tries to shrink the ^{weights} features of unimportant features.
- ⇒ Helping us to know "useful" features.

For L2 regularisation

hyperparameter: λ

$$\Rightarrow J(\beta) = \sum_{i=1}^N e^{x_i^T \beta} - y_i^T (\bar{x}_i^T \beta) - \log(y_i!) + \frac{\lambda}{2} \sum_{j=1}^M |\beta_j|^2$$

$$\Rightarrow \left. \frac{\partial J}{\partial \beta} \right|_{L_2\text{-reg.}} = \left. \frac{\partial J}{\partial \beta} \right|_{\text{non-reg.}} + \lambda(\beta)$$

$$\Rightarrow \beta_{n+1} = \beta_n - \eta \left(\left. \frac{\partial J}{\partial \beta} \right|_{\text{non-reg.}} \right) - \lambda(\beta_n)$$

$$\boxed{\beta_{n+1} = \beta_n(1-\lambda) - \eta \left(\left. \frac{\partial J}{\partial \beta} \right|_{\text{non-reg.}} \right)}$$

⇒ Used for avoiding over-fitting.