Simple Linear Regression

- simple approach to predict quantitative value (response) Y on the basis of a predictor variable (X)

$$\hat{Y} \approx \hat{\beta}_0 + \hat{\beta}_1 X$$
 (preclication)

- We intend to find / estimate the coeefficient, 13°, and 13°, so that we can represent the patter emerged by the data, as close as possible.

We use "MSE" to measure "closeness" / accuracy of the model.

Let $\hat{y}_i = \hat{\beta}_{i+1} + \hat{\beta}_{i} \times i$, represent the prediction for Y based on the ith value of X.

We define

e;=y;-ÿ; to be the ith nesidual

We define,

RSS = Ci^+ Cz^+ Cz^+ + Cz^+ + en

Residualsum of squines

perivation
$$L = \frac{1}{2} \frac{\Sigma}{(\gamma_{i} - \dot{\gamma}_{i})^{2}} \qquad \frac{5L}{J\beta_{0}} = \frac{1}{3} \frac{\Sigma}{J\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i})^{2} = 0$$

$$= \frac{1}{n} \frac{\Sigma}{J\beta_{0}} - \frac{1}{3} \frac{\Sigma}{J\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i})^{2} = 0$$

$$= \frac{1}{n} \frac{\Sigma}{I\beta_{0}} - 2(\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \frac{\Sigma}{I\beta_{0}} (\gamma_{i} - \dot{\beta}_{0} - \dot{\beta}_{1} \dot{x}_{i}) = 0$$

$$= \left[\overline{Y} - \beta_1 \overline{X} = \underline{M}_{30} = \beta_0 \right]$$

$$\frac{\int L}{\int B_i} = \frac{d}{\int B_i} \left(\gamma_i - \beta_0 - \beta_1 \gamma_i \right)^2 = 0$$

$$\rightarrow Y_1 - \overline{Y} + \beta_1 \overline{X} - \beta_1 X_1$$

$$\frac{\partial \beta_{1}}{\partial \beta_{1}} \qquad \frac{\partial \beta_{1}}{\partial \beta_{1}}$$

$$= \underbrace{\frac{1}{n}}_{i=1} \underbrace{\frac{2}{x}}_{i} 2(y_{i} - (\overline{y} - y_{i}\overline{x}) - \hat{y}_{i}\overline{x}) - \hat{y}_{i}x_{i}). (-x_{i} + \overline{x}) = 0$$

$$= \sum_{i=1}^{3} -2 (y_i + \beta_i \bar{x} - \bar{y} - \beta_i x_i) \cdot (\bar{x} - x_i) = 0$$

$$= \sum \left[\left(\begin{array}{c} Y_{1} - \overline{Y} \end{array} \right) - \int_{\overline{X}_{1}}^{x} \left(X_{1} - \overline{X} \right) \right] \left(X_{1} - \overline{X} \right) = 0$$

$$= \sum (Y_i - \overline{Y})(X_i - \overline{X}) = \beta_i^{\Lambda} \sum (X_i - \overline{X})^2$$

$$|\hat{x}| = \sum_{j=1}^{n} \frac{(y_j - \overline{y})(x_j - \overline{x})}{\sum_{j=1}^{n} (x_j - \overline{x})^2}$$

- -> 30 is the value of Y when X=0
- -> |31 represents the increase in y associated un/ one unit increase in X.

The concept of bias (Sample variance)

→ An unbiased estimator is one that, on average, equal to the true population estimate.

For eg - an unbiased mean for the sample would:

- · gn one sample overestimated.
- · In one sample understimated.

Basically, the sample estimate might not be exact in terms of the population, but it will not systematically overestimate or

Standard Error (Var (jí))

-> We know that averages of fi over different samples will be very close to f. However, we don't know how World wille be a single estimat y. To do this, we use standard Errar

$$Var(\hat{H}) = SE(\hat{Y})^2 = \frac{\sigma^2}{n}$$

$$\Rightarrow SE(\hat{H}) = \frac{\sigma}{\sqrt{n}} \Rightarrow standard deviation$$

Similarly, for regression coefficients

SE(
$$\beta$$
₀)² = $\nabla^2 \left[\frac{1}{n} + \frac{\pi^2}{\pi^2} \right]$

Notation of energy variance of energy.

SE(β ₁)² = ∇^2

$$\frac{\pi^2}{\Sigma} (\chi_1 - \overline{\chi})^2$$

$$t^2 = \sqrt{(RSS)/(n-2)}$$
Residual Standard

Standard Error

RECAP
$$\begin{bmatrix}
\hat{\beta}_0 = \overline{y} - \hat{\beta}_1 \overline{x} \\
\\
\hat{\beta}_1 = \Sigma (\underline{y}_1 - \overline{y}) (\underline{x}_1 - \overline{x}) \\
\underline{\Sigma} (\underline{x}_1 - \overline{x})^2$$

$$SE(\hat{\beta})^{2} = \sigma^{2} \rightarrow Von(\mathcal{E}) = \sum_{i=1}^{2} \underbrace{\mathcal{E}_{i}}_{n-2}$$

$$SE(\hat{\beta})^{2} = \sigma^{2} \left[\frac{1}{n} + \overline{x}^{2} \right]$$

$$SE(\hat{\beta})^{2} = \sigma^{2} \left[\frac{1}{n} + \overline{x}^{2} \right]$$

$$SE(\hat{\beta})^{2} = \sigma^{2} \left[\frac{1}{n} + \overline{x}^{2} \right]$$

$$SE(\hat{\beta})^{2} = \sigma^{2}$$

$$\Sigma(\hat{\lambda}_{i} - \overline{x})^{2}$$

Assumption:
$$\varepsilon \sim N(0, \sigma^2)$$

This implies that for a given value of $x = x_i$, the value of y follows $N \sim (\beta_0 + \beta_1 x_i, \sigma^2)$

We use standard error to estimate the confidence intermeds. For LR:

 $\beta_1 \pm 2.SE(\beta_1) \Rightarrow$ there is 95% chance that the interval $[\beta_1 - 2.SE(\beta_1)]$ will contain the true value of the parameter.

Hypothesis Testing



Null Hypothesis: Typically assumes no relation blue any of the features, For LR (single beature):

Ho: There's na relationship b/u X and Y i.e \\ \beta_i = 0

Alternate Hypothesis: Assumes there's some relationship b/we X and Y. Ha: 3, #0

To test the hypothesis we calculate t-statistic
$$1 + \frac{1}{2} = \frac{1}{2} = \frac{1}{2} = 0$$
 represents the difference blue the estimated slope and the hypothesized value of $\frac{1}{2}$.

the larger the abs. value of t, the more evidence we have against Ho. If $SE(\beta_i)\uparrow\Rightarrow t\downarrow$, giving everidence for Ho.

Accuracy of the model

RSE (Residual Standard Error)

$$RSE = \sqrt{\frac{1 \cdot RSS}{n-2}}$$

$$RSE = \sqrt{\frac{1}{n-2} \cdot \frac{2}{(\gamma_i - \hat{\gamma}_i)^2}}$$

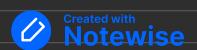
$$\sqrt{\frac{1}{n-2} \cdot \frac{2}{i=1}}$$

RZ This is blue 0 and 1, and independent of the scale of Y.

$$R^{2} = 1 - RSS$$

$$TSS$$

$$TSS = \sum_{i=1}^{n} (y_{i} - \overline{y})^{2} \quad (variance)$$



3.7 Multiple Linear Regression

$$Y = |_{30} + |_{31}X_{1} + |_{32}X_{2} + \cdots + |_{3p}X_{p} + \varepsilon$$

$$RSS = \sum_{i=1}^{n} (Y_{i} - (\beta_{0} + \beta_{1}X_{1} + \cdots + \beta_{p}X_{p} + \varepsilon))$$

1. Is there a relationship b/w the Response and Predictors? NWU Hypothesis (H_0) : $\hat{\beta}_1 = \hat{\beta}_2 = \cdots = \beta_P = 0$ (There's no relationship b/w predictors and response.

Alternate (Ma): atleast one Bi is non-zero; 1 « i « P

Bias-Variance Tradeoff

While its important to minimise training MSE, it's equally important to get a minimized MSE on test points.

$$E(y_0 - f(x_0))^2 = Vol(f(x_0)) + [Bias(f(x_0))] + Vol(E)$$
(Exp. test MSE)

Variance refers to the amount by which is would change if we estimated it using different sets. Essentially, if a method has higher variance, it would mean that a small change in the data can get large changes in F.

Bias refers to the inability of the model to accurately represent the patterns in the dataset.

- -> easy to obtaing lave bias and higher variance or high variance and lave bias
- → The challerge is to find an ideal midway blue the true, which should ensure lave bias and veriance. O Notewise