

MULTIPLE LINEAR REGRESSION (MLR)

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik} + \epsilon_i$$

Response surface: hyperplane in $k+1$ dimensions

$$\begin{bmatrix} Y \\ Y_1 \\ Y_2 \\ \vdots \\ Y_N \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1k} \\ 1 & X_{21} & X_{22} & \dots & X_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{N1} & X_{N2} & \dots & X_{Nk} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{bmatrix}$$

$$Y = X\beta$$

$$X^T Y = X^T X \beta$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\text{Moore-Penrose pseudoinverse}$$

$$\hat{Y} = X\hat{\beta} = X(X^T X)^{-1} X^T Y$$

$$\hat{Y} = HY$$

Partial correlation

Effect of X_1 on Y without X_2

$$s_{Y_{X_1, X_2}} = \frac{s_{Y_{X_1}} - s_{Y_{X_2}} s_{Y_{X_2, X_1}}}{\sqrt{(1-s_{Y_{X_1}}^2)(1-s_{Y_{X_2}}^2)}}$$

Qualitative/Categorical Variables

One-hot encode, introduce new dummy variables for each.

For N categories, $n-1$ dummy vars.
Encode 0 as something
OTHERWISE: perfect multicollinearity

REGRESSION

SLR:

$$Y = \beta_0 + \beta_1 X + \epsilon$$

MLR:

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_n X_n + \epsilon$$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

Non-linear:

$$Y = \beta_0 + \beta_1 X_1 + \frac{1}{\beta_2 + \beta_3 X_2} + \epsilon$$

OLS ASSUMPTIONS

- Regression model is linear w.r.t parameters
- Independent var X is deterministic
- Residuals ϵ are normally distributed
- Conditional expected value of residual is zero: $E[\epsilon|X] = 0$
- Variance of residuals is constant: homoscedasticity
- Residuals are uncorrelated

SOLUTION: OLS

Provides Best Linear Unbiased Estimate (BLUE): $E[\hat{\beta} - \beta] = 0$

OLS: minimise SSE

$$SSE = \sum_{i=1}^n \epsilon_i^2$$

Min: $Y = \beta_0 + \beta_1 X_1 + \epsilon$

$$\Rightarrow SSE = \sum_{i=1}^n (Y - \beta_0 - \beta_1 X_1)^2$$

SLR PROCESSES

- Collect data
- Define functional form of regression
- Estimate model parameters
- Train-test split
- Descriptive analysis
- Test for validation
- Display

OLS PROOF

Linear Algebra

$Ax = b$

$$\begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

b does not lie on $C(\alpha)$ so we find best approx by project

$\hat{\beta}_0 = \bar{Y} - \bar{\beta}_1 \bar{X}$

$$\hat{\beta}_1 = \frac{\bar{Y}Y - \bar{X}\bar{Y}}{\bar{X}^2 - (\bar{X})^2}$$

$$\hat{\beta}_1 = \frac{\text{cov}(XY)}{\text{var}(X)} = \frac{n \bar{X} \bar{Y}}{\text{var}(X)}$$

MLE

Assume errors $\sim N(0, \sigma^2)$

$$\Rightarrow Y \sim N(\beta_0 + \beta_1 X, \sigma^2)$$

$$f(y|x; \beta_0, \beta_1, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y-\beta_0-\beta_1 x)^2}{2\sigma^2}}$$

Likelihood and maximise

$$\text{Log likelihood} \rightarrow \text{minimise. Sub: } \beta_0 + \beta_1 x_1 \text{ as g. you will get SSE.}$$

Calculus

Set partial derivatives of $SSE = 0$ and solve

$$SSE = \sum_{i=1}^n (Y - \beta_0 - \beta_1 X_i)^2$$

Gradient descent

$$SSE = \frac{1}{2} \sum_{i=1}^n (Y - \beta_0 - \beta_1 X_i)^2 = \frac{1}{2} (Y - \hat{Y})^2$$

$$\frac{\partial}{\partial \beta_0} (SSE) = -2 \sum_{i=1}^n (Y - \beta_0 - \beta_1 X_i)$$

$$\frac{\partial}{\partial \beta_1} (SSE) = -2 \sum_{i=1}^n X_i (Y - \beta_0 - \beta_1 X_i)$$

Now update:

$$\beta_0 = \beta_0 - \eta \frac{\partial (SSE)}{\partial \beta_0}$$

$$\beta_1 = \beta_1 - \eta \frac{\partial (SSE)}{\partial \beta_1}$$

VALIDATION OF SLR

- R^2 : coeff. of determination measures proportion of variance explained by model
- $SST = \sum_{i=1}^n (y_i - \bar{Y})^2$
- $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{Y})^2$
- $SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$
- $R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{Y})^2}$

Parameter Interpretation

$$Y = \beta_0 + \beta_1 X \equiv \beta_0 + \beta_1 X_1$$

β_0 : expected value of Y when X is zero (it could be anything)

β_1 : rate of change of Y w.r.t X (change in Y for unit X here (it could be anything, since for β_1)

Residual analysis

- Normally distributed
- P-P plot b/w observed cumulative probability vs expected cumulative probability
- Same: $e_1 < e_2 < \dots < e_n < e_0$
- Obs prob = $\frac{k}{n} \text{ for } e_k$
- Expected prob = $\Phi(\frac{e_k - \bar{e}}{s_e})$

Homoscedasticity

variance should be independent of X

$s.e.$ vs X

Functional form of regression

any pattern in p vs b need to be model w.r.t then need to be transformed

transform variables (log, square, etc.)

add new vars (non-linear model)

Durbin-Watson

Z -score

$$Z = \frac{\bar{Y} - \bar{Y}}{s_e} > 3 \rightarrow \text{outlier}$$

Mahalanobis Distance

1. transform col into uncorrelated vars

2. Scale col to make variance 1

3. Finally calc Euclidean distance

$MD = \sqrt{(X - \bar{X})^T C^{-1} (X - \bar{X})}$ w.r.t some distribution

$MD > X_t^*$ critical w.r.t dist = no of X :

$MD > X_t^* \Rightarrow \text{outlier}$

SLR VALIDATION

- R^2 score been there done that
- Hypothesis testing with β_1 , $\beta_1 = 0 \rightarrow$ no relationship, $H_0: \beta_1 = 0$ $H_A: \beta_1 \neq 0$ two-tailed t-test
- $H_0: \beta_1 = 0 \rightarrow$ no relationship, $H_0: \beta_1 \neq 0$ two-tailed t-test
- $\hat{\beta}_1 \rightarrow$ sampled from a normal distn of all $\hat{\beta}_1$
- $s.e.(\hat{\beta}_1) \rightarrow$ SD of distn
- $t\text{-statistic} = \frac{\hat{\beta}_1 - \beta_1}{s.e.(\hat{\beta}_1)} = \frac{\hat{\beta}_1}{s.e.(\hat{\beta}_1)}$

Partial corr

corr b/w Y and X_i after removing effect of X_j on both X_i and Y

$$s_{Y_{X_i, X_j}} = \frac{s_{Y_{X_i}} - (s_{Y_{X_j}})(s_{X_i, X_j})}{\sqrt{(1-s_{Y_{X_i}}^2)(1-s_{X_i, X_j}^2)}}$$

(3) Hypothesis testing for whole model

(4) Residual analysis

→ Scales of diff vars very → P-P plot

→ homoscedasticity ↓

plot standardised residuals against standardized \hat{Y}

any pattern of homoscedastic variance depends on indep. var

ensure functional form

→ residual plot pattern = improve by transform, adding new, on non-linear

(5) Outlier analysis

MLR VALIDATION

- Adjusted R^2 (w.r.t multiple determination)**
- $R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$
- $\text{Adjusted } R^2 = 1 - \frac{SSE / (n-k-1)}{SST / (n-1)} = 1 - \frac{SSE / (n-1)}{SST / (n-1)}$
- $\text{Adjusted } R^2 = 1 - (1 - R^2) \frac{(n-1)}{(n-k-1)} \rightarrow \text{adj } R^2 = R^2$

(2) Hypothesis test for neg. coefficient

Same idea as SLR.

$$S_e = \sqrt{\frac{\sum_{i=1}^{n-1} (y_i - \hat{y}_i)^2}{n-1}}$$

$$S_e(\hat{\beta}_j) = \sqrt{\frac{\sum_{i=1}^{n-1} (\hat{y}_i - \hat{y})^2}{n-1}}$$

$$t\text{-statistic} = \frac{\hat{\beta}_j - \beta_j}{S_e(\hat{\beta}_j)}$$

(3) Significance of entire model: ANOVA

$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$

$H_A: \text{otherwise}$

$F\text{-statistic} = \frac{SSW/k}{SST/(n-1)}$

(4) Significance of part k th model

Full: $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$

Reduced: $H_0: \beta_1 = \beta_2 = \dots = \beta_{k-1} = 0$

$H_0: \beta_k = \beta_{k+1} = \dots = \beta_n = 0$

$H_1: \text{otherwise}$

$F\text{-statistic} = \frac{(SSR - SSE)/(k-1)}{SSE/(n-k-1)} \rightarrow \text{adj } R^2 = \frac{SSE/(n-k-1)}{SST/(n-1)}$

(5) Residual analysis

Same deal as SLR.

(2) Hypothesis test for individual vars - t-test

$\hat{\beta} = (X^T X)^{-1} X^T Y$

linear func follows Normal distn

$\hat{\beta}_j \rightarrow$ derived from sample

So for some X_j :

$H_0: \beta_j = 0$ (no relationship)

$H_A: \beta_j \neq 0$

two-tailed t-test.

(3) Verifying entire model using F-test

F-test statistic = $\frac{MSR/k}{MSE/(n-k-1)} = \frac{SSR/k}{SSE/(n-k-1)}$

$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$

$H_A: \text{not all } \beta_j = 0$

NOTE: this does not tell us which β are non-zero

(4) Verifying parts of the model with F-test

Full model

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

Reduced model

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_{k-1} X_{k-1}$$

Objective of partial F-test: check additional var to x_j for significance

$H_0: \beta_{k+1} = \beta_{k+2} = \dots = \beta_n = 0$

$H_A: \text{some non-zero}$

F-test statistic = $\frac{(SSE - SSE_f)/(k-n)}{SSE_f/(n-k-1)}$

(5) Residual analysis

→ normal distn residuals \rightarrow P-P plot

→ homoscedasticity ↓

plot std. residuals against std. y

pattern = heteroscedasticity

→ ensure functional form (linear instead of log-linear)

(6) Multicollinearity & Variance Inflation Factor

indep vars corr b/w each other

Say $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$.

$x_1 = x_0 + \beta_1 x_1$

Let R^2 be R^2 value for second model

$VIF = \frac{1}{1 - R^2}$

TIP: and by which β std. error inflates \rightarrow std. error inflated

$t_{\text{actual}} = \frac{\hat{\beta}}{s.e.(\hat{\beta})} \cdot VIF$

threshold (VIF) = 4

MULTICOLLINEARITY: VARIANCE INFLATION FACTOR

$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$

To check for corr b/w x_1, x_2 :

$X = [x_0 \ x_1 \ x_2] \rightarrow R^2$ here: R^2

$VIF = \frac{1}{1 - R^2}$

Auto-correlation: Durbin-Watson Test

$t_{\text{actual}} = \frac{\hat{\beta} - \beta^0}{s.e.(\hat{\beta})} \times VIF$

$t_{\text{actual}} = \frac{\hat{\beta} - \beta^0}{s.e.(\hat{\beta})} \times \frac{1}{1 - R^2}$

$\rightarrow \epsilon_t \text{ and } \epsilon_{t-1} \text{ are related} \rightarrow \text{std. error inflated} \rightarrow t\text{-statistic inflated}$

Say $\rho = \text{corr b/w } \epsilon_t, \epsilon_{t-1}$

$H_0: \rho = 0$

$H_A: \rho \neq 0$

$D = \frac{\sum (e_t - e_{t-1})^2}{\sum e_t^2} = 2 \left(1 - \frac{\sum e_t e_{t-1}}{\sum e_t^2} \right)$

$D \in [0, 4]$

Remove var

PCA/SVD Lasso/Ridge

$D < D_L: \text{positive corr}$

$D > D_U: \text{no corr}$

$D_L < D < D_U: \text{inconclusive}$

$(A-D) < D_U: \text{negative corr}$

$(A-D) > D_U: \text{no negative corr}$

$D_L < (A-D) < D_U: \text{inconclusive}$

$D = 2 \Rightarrow \text{no autocorr}$