

MLR.

Theory.

$$Y_i = \beta_0 + \beta_1 X + \beta_2 X + \dots + \beta_p X^p + \epsilon_i$$

$$E\{Y\} = \beta_0 + \beta_1 X + \beta_2 X + \dots + \beta_p X^p$$

obs: linear model means that variables have an additive effect.

General Linear Regression Model can be

1) \Rightarrow X are continuous variables

2) \Rightarrow X are qualitative variables (gender, color)

3) \Rightarrow polynomial regression

\rightarrow some $X_i = X_{i-j}^R$ where $j \geq 1$

4) \Rightarrow transformed variables

$\rightarrow X'_i = f(X_i)$ and

$$Y_i = \beta_0 + \beta_1 X_1 + \dots + \beta_k f(X_k) + \dots + \beta_p X_p + \epsilon_i$$

5) \Rightarrow interaction effects

$\rightarrow X'_i = X_i \cdot X_j$ with $i \neq j$

$$Y_i = \beta_0 + \beta_1 X_1 + \dots + \beta_i X_i + \beta_j X_j + \beta_{j+1} X_i X_j + \epsilon_i$$

Matrix formulation:

$$Y = X\beta + \epsilon$$

$$\sigma^2(\epsilon) = \sigma^2 \cdot I$$

$$\left\{ \begin{array}{l} E\{Y\} = X\beta \\ \sigma^2\{Y\} = \sigma^2 I \end{array} \right.$$

$$Y = Xb \Rightarrow X'Y = X'Xb$$

$$\Rightarrow (X'X)^{-1}(X'Y) = b$$

$$\hat{y} = Xb$$

$$e = Y - \hat{y} = Y - Xb = HY$$

$$H = X(X'X)^{-1}X'$$

$$e = (I - H)Y$$

$$\sigma^2[e] = \sigma^2(I - H)$$

$$s^2[e] = MSE(I - H)$$

ANOVA for MLR

$$SST = \sum (y_i - \bar{y})^2 = Y'Y - \left(\frac{1}{n}\right)Y'JY$$

$$SSE = \sum (y_i - \hat{y}_i)^2 = Y'Y - b'X'Y$$

$$SSR = \sum (\hat{y}_i - \bar{y})^2 = b'X'Y - \left(\frac{1}{n}\right)Y'JY$$

$$E\left[\frac{SSE}{n-p}\right] = E[MSE] = \sigma^2$$

$$E\left[\frac{SSR}{p-1}\right] = E[MSE] = \sigma^2 + \frac{1}{2} \left[\sum_{j=1}^p \beta_j (x_{ij} - \bar{x}_j)^2 + \sum_{j,k} \beta_j \beta_k \sum (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k) \right]$$

$$\left. \begin{matrix} MSE \sim \chi^2(n-p) \\ MSR \sim \chi^2(p-1) \end{matrix} \right\} \Rightarrow F = \frac{MSR}{MSE} \sim F(p-1, n-p)$$

Coefficient of Multiple Determination

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

$$0 \leq R^2 \leq 1$$

$$\text{Adjusted } R^2 = 1 - \frac{MSE}{MST} = 1 - \frac{n-1}{n-p} \frac{SSE}{SST}$$

Inference about Regression parameters

MSE & ML estimators are unbiased.

$$E\{b\} = \beta$$

$$\sigma^2\{b\} = \begin{bmatrix} \sigma^2\{b_0\} & \sigma\{b_0, b_1\} & \dots & \sigma\{b_0, b_{p-1}\} \\ \sigma\{b_{p-1}, b_0\} & \sigma\{b_{p-1}, b_1\} & \dots & \sigma^2\{b_{p-1}\} \end{bmatrix}$$

$$\sigma^2\{b\} = \sigma^2 (X^T X)^{-1}$$

$$S^2(b) = MSE (X^T X)^{-1}$$

$$\frac{b_k - \beta_k}{S(b_k)} \sim t(n-p)$$

Assumption of linear regression.

- ① The relationship is linear
- ② There is no multicollinearity (predictors are independent algebraically \Rightarrow invertible)
 \rightarrow the inverse of $X^T X$ exists.
- ③ the residuals are normally distributed
- ④ residuals are not correlated (no auto-correlation)
sequential
- ⑤ variance is constant (homoskedasticity).

steps in developing linear regression models.

- ① plot.
- ② fit regression using information from plots.
- ③ check significance
- ④ check computational issues.
 - outliers for LSE
 - difference in magnitude → standardized model.
- ⑤ check assumptions.
- ⑥ fix or mitigate violation of assumptions or computational issues.
- ⑦ redo step 2-6 with a different regression equation.
- ⑧ select the final model from the models that have a solution in step 6.
- ⑨. Validate the model on independent data.

Heter et. al presents a simplified version in Fig 8.1

- ① Data collection and preparation
- ② Reduction of explanatory or predictor variables
- ③ Model refinement and selection
- ④ Model validation.

In this class we focus on model refinement & selection.

Plot is offering you a first insight on the relationship between predicted and predictor, usually a one-to-one relationship (very rare predicted vs. 2 predictor).

Brute force assessment = all possible regressions.

$$Y = Xb$$

$$\Rightarrow$$

$$Y = X_1$$

$$Y = X_2$$

$$Y = X_P$$

$$Y = X_1, X_2$$

$$Y = X_1, X_3$$

$$\vdots$$

$$Y = X_1, \dots, X_P$$

$$Y = X_2, X_3$$

$$\vdots$$

$$Y = X_{P-1}, X_P$$

$$Y = X_1, X_2, \dots, X_P$$

$$\text{Total \# regressions} = 2^P = 0 + C_P^1 + C_P^2 + \dots + C_P^{P-1} + C_P^P = 2^P$$

Selection among candidates using various statistics

① R_P^2 or $SS\bar{E}_P$.

② MSE_P or R_a^2

③ C_P criterion

④ $PRESS_P$ criterion.

$$\textcircled{1} \quad R^2 = 1 - \frac{SSE}{SST}$$

$$R_p^2 = 1 - \frac{SSE_p}{SST} \quad \text{stands for the } p \text{ variable(s).}$$

R_p^2 is not intended to identify which subset is "the best" (maximizes or minimizes a criterion)

R^2 increases with the # variables

R_p^2 is intended for finding the point where adding more explanatory variables is not worthy.

②. MSE_p or R_a criterion.

$$R_a^2 = 1 - \frac{n-1}{n-p} \frac{SSE}{SST}$$

$$MSE_p = \frac{SSE_p}{n-p-1}$$

where p is the # variables

Example.

$$MSE_1 = \frac{SSE(X_1)}{n-2}$$

$$MSE_{1,2} = \frac{SSE(X_1, X_2)}{n-3}$$

Interpretation with $n-3$ screen plot.

③. C_p -criterion.

C_p is concerned with total mean square error of the n fitted values for each subset regression model.

error in fitted value is: $\hat{Y}_i - \mu_i = \hat{Y}_i - E[\hat{Y}_i]$

$$MSE_{\hat{Y}_i - \mu_i} = E[(\hat{Y}_i - \mu_i)^2] = E[(E[\hat{Y}_i] - \mu_i) + (\hat{Y}_i - E[\hat{Y}_i])]^2 =$$

$$= (E[\hat{Y}_i] - \mu_i)^2 + \sigma^2[\hat{Y}_i]$$

$$\sum [(E[\hat{y}_i] - \mu_i)^2 + \sigma^2(\hat{y}_i)] = \sum (E[\hat{y}_i] - \mu_i)^2 + \sum \sigma^2(\hat{y}_i)$$

$$\Gamma_p = \frac{MSE \hat{y}_i - \mu_i}{\sigma^2}$$

an unbiased estimator of Γ_p is C_p

$$C_p = \frac{SS \bar{E}_p}{MSE(X_1, \dots, X_{p-1})} - (n - 2p)$$

If there is no bias then $E\{\hat{y}_i\} = \mu_i$ and

$$E\{C_p\} \approx p$$

NOTE C_p identifies subset of X for predicting Y

② PRESS_p criterion.

PRESS = prediction of sum of squares.

PRESS measures how well the use of fitted values for a subset model can predict the observed y .

- Steps:
- ① In PRESS the i th observation is deleted and the regression is estimated from the $n-1$ observations.
 - ② Predict the i th observation with the fitted regression
 - ③ compute PRESS statistics:

$$PRESS_p = \sum (y_i - \hat{y}_{i(i)})^2$$

where $\hat{y}_{i(i)}$ is the predicted y_i from the model with n th case.

Decision: Small PRESS indicates good model.

AUTOMATIC SEARCH PROCEDURES for variable reduction

Identify the "best" subset of variables from a large # of predictors.

① Forward selection procedure.

is an algorithm that is based on F test.

Step 1. Fit a SLR for each p variable

$$y = b_0 + b_1 x$$

$$F^* = \frac{MSE(X_k)}{MSE(X_k)}$$

where for SLR $MSE(X_k) = SSR(X_k)$

The variable with LARGEST F^* is selected.

Step 2. Fit all regression models with three selected variable in step 1.

let say X_{S1} was selected.

$$y = b_0 + b_1 X_{S1} + X_k$$

$$\text{compute } F^* = \frac{MSE(X_k | X_{S1})}{MSE(X_{S1}, X_k)}$$

where $MSE(X_k | X_{S1})$ is MSE with the variable X_k and X_{S1} (X_k given X_{S1})

1 $MSE(X_{S1}, X_k)$ is MSE for the regression with 2 variables X_{S1} and X_k

$$F^* = \left(\frac{b_k}{S(b_k)} \right)^2$$

Among all models select the one with the largest F^* , let say X_{S2}

step 3 assess if any variable should be dropped from the model.

This is done to eliminate variable whose participation in the model is "covered" by other variables or combination of variables.

If no variable is dropped then both X_{S1} & X_{S2} are kept in the model

Step 4. Add a new variable, and execute the regression $Y = b_0 + b_1 X_{S1} + b_2 X_{S2} + b_3 X_K$
among all $F = \frac{RSS(X_K | X_{S1}, X_{S2})}{MSE(X_{S1}, X_{S2}, X_K)}$

select the largest (if significant).

step 5 → Go to step 3.
→ stop if no new variable is added or dropped

Step 6. → Go to step 4.
→ stop if no new variable is added.

Forward selection procedure.

same as stepwise but only add variable, do not drop them (looks only forward)

Backward selection procedure.

start with all variables in and drop one variable at a time.

there is no addition of variables.