

# Regression Analysis

## (1) Linear Model:

A regression model is linear model:  $Y = X\beta + \varepsilon$ .

Where  $X_{n \times p}$  is usually full rank  $p$ .  $C(X)$  contains categorical variables and at least one continue variables. Thus,  $X^T X$  is nonsingular.

## (1) Simple linear Regression:

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad (1 \leq i \leq n, \quad E(\varepsilon_i) = 0, \quad \text{Var}(\varepsilon_i) = \sigma^2)$$

$$\Rightarrow X = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \Rightarrow \hat{\beta} = \frac{1}{\sum_1^n (X_i - \bar{X})^2} \begin{pmatrix} n\bar{Y} \sum_1^n X_i^2 - n\bar{X} \sum_1^n X_i Y_i \\ -n\bar{X} \bar{Y} + n \sum_1^n X_i Y_i \end{pmatrix}$$

$$\text{More formally, } \hat{\beta}_1 = \frac{\sum_1^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_1^n (X_i - \bar{X})^2} = \frac{\widehat{\text{Cov}}(X, Y)}{\widehat{\text{Var}}(X)}$$

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

which follows from  $M = (m_{ij}), \quad m_{ij} = \frac{1}{n} + \frac{(X_i - \bar{X})(X_j - \bar{X})}{\sum_1^n (X_i - \bar{X})^2}$

## (2) Multiple Linear Regression:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \varepsilon_i, \quad (1 \leq i \leq n, \quad \begin{cases} E(\varepsilon_i) = 0 \\ \text{Var}(\varepsilon_i) = \sigma^2 \end{cases} \text{ i.i.d.})$$

More generally, consider  $Y = X\beta + \varepsilon, \quad X \in \mathbb{R}^{n \times p}$ .

Note that  $C(M) = C(X)$  has rank  $p$ .

Decompose into sum of  $p$  orthogonal subspaces of  $\dim = 1$ .

$$C(M) = \sum_{i=1}^p C(M_i), \quad r(M_i) = 1, \quad M_i M_j = 0, \quad i \neq j$$

Def:  $Z_i = (X_1, X_2, \dots, X_i)$   $1 \leq i \leq p$ .  $Z_0 = 0$ .

Let  $P_{Z_i} = Z_i (Z_i^T Z_i)^{-1} Z_i^T$ . Then  $M_i = P_{Z_i} - P_{Z_{i-1}}$ .  $1 \leq i \leq p$ .

It is easy to check  $\{M_i\}_i^p$  is what we want.

Lemma:  $SSR(X_j | X_1, \dots, X_{j-1}) = Y^T M_j Y$ .

Prmk: It can be interpreted as regression sum of squares due to add  $X_i$  when  $X_1, \dots, X_{i-1}$  are already in model.

$$\Rightarrow SSR(X) = \sum_i^p SSR(X_i | X_1, \dots, X_{i-1})$$

Prmk: The orthogonal Decomposition of  $C(X)$  depends on the order of variables  $X_i$  being fit into the model. So the decomposition may not be unique (c.f. Unbalanced ANCOVA)

## (2) Best Linear Prediction:

One of the main goals and use of regression model is for prediction but not only to make inference on parameters.

That is, predict  $Y$  on basis of information of  $X_1, X_2, \dots, X_p$ . Find  $f(X)$  to minimize:

$$E(Y - f(X))^2 \text{ which need joint dist of } (X, Y).$$

i) Dist of  $(X, Y)$  is known:

Thm. Let  $m(x) = m(Y|X)$ . Then  $m(x) = \operatorname{argmin} E(Y - f(x))^2$

Pf: Check  $E((Y - m(x))(m(x) - f(x))) = 0$

Prop.  $\operatorname{Cov}(Y, f(x)) = \operatorname{Cov}(m(x), f(x))$

Pf: Check  $E((Y - m(x))f(x)) = 0$

Thm. (Criteria of Best Predictor)

$\tilde{\eta}(x)$  is any predictor. Then  $\tilde{\eta}(x) = m(x)$  n.s.

$\Leftrightarrow \operatorname{Cov}(f(x), Y - \tilde{\eta}(x)) = 0$  for any function  $f$ .

and  $E(\tilde{\eta}(x)) = m_Y$ .

Pf:  $(\Leftarrow)$  It's from prop. above

$(\Rightarrow)$  Prove:  $\sigma^2(\tilde{\eta}(x) - m(x)) = 0$ . Since  $E(\tilde{\eta}(x) - m(x)) = 0$ .

$$\text{LHS} = \operatorname{Cov}(\eta - m(x), \tilde{\eta}(x) - m(x)) - \operatorname{Cov}(\eta - \tilde{\eta}(x), \tilde{\eta}(x) - m(x))$$

$$= \textcircled{1} + \textcircled{2} = 0. \quad \textcircled{1} = 0 \text{ by prop. } \textcircled{2} = 0 \text{ by condition}$$

ii) List of  $(X, Y)$  is unknown:

We only can find best linear predictor of  $Y$  if we only know means, vars of  $X, Y$  and  $\operatorname{Cov}(X, Y)$ .

i.e. minimize  $E(Y - f(x))^2$ .  $f(x) = \alpha + X^T \beta$ .

Prbk: If  $(X, Y)$  is multipol. Then the best prediction is linear predictor:

$$E(Y|X) = m_Y + \Sigma_{YX} \Sigma_{XX}^{-1} (X - m_X)$$

Thm.  $\beta_x$  is solution of  $\Sigma_{xx}\beta = \Sigma_{xy}$ . Then:  $\hat{E}(Y|X) =$

$\bar{m}_y + (X - \bar{m}_x)^T \beta_x$  is best linear predictor

Pf: Show:  $E((Y - \hat{E}(Y|X))( \hat{E}(Y|X) - f(x))) = 0$

where  $f(x) = \eta + (X - \bar{m}_x)^T \beta^T$ .

$$\Rightarrow E(Y - \alpha - X^T \beta)^2 = E(Y - \hat{E}(Y|X))^2 + E(\hat{E}(Y|X) - \alpha)^2$$

Remark: BLP is unique, i.e. to minimize:

$E((\hat{E}(Y|X) - \eta - (X - \bar{m}_x)^T \beta)^2)$ .  $\eta = \bar{m}_y$  and  $\beta$  must be solution of  $\Sigma_{xx}\beta = \Sigma_{xy}$ .

Note it's:  $(\bar{m}_y - \eta)^2 + E((X - \bar{m}_x)^T (\beta_x - \beta))^2$

To minimization require:  $\eta = \bar{m}_y$ .  $E(\square)^2 = 0$

Since  $E(\square)^2 = (\beta_x - \beta)^T \Sigma_{xx} (\beta_x - \beta)$ .  $\Leftrightarrow$

$$\Sigma_{xx} (\beta_x - \beta) = 0 \Leftrightarrow \Sigma_{xx} \beta_x = \Sigma_{xx} \beta$$

iii) Apply to multiple linear regression:

Denote:  $S_{xx} = \frac{X^T(I - P_n)X}{n-1}$ ,  $S_{xy} = \frac{X^T(I - P_n)Y}{n-1}$

$$\hat{m}_x = \frac{1}{n} J_n^T X, \quad \hat{m}_y = \frac{1}{n} J_n^T Y.$$

$\Rightarrow$  Best linear prediction of  $Y_i$ :  $\hat{Y}_i = \hat{m}_y + (X_i - \hat{m}_x)^T \hat{\beta}_x$

where  $\hat{\beta}_x$  is solution of  $S_{xx} \hat{\beta}_x = S_{xy}$

Remark:  $\hat{\beta}_x$  is LSE of  $Y_i = \alpha + (X_i - \hat{m}_x)^T \beta + \epsilon_i$ ,  $1 \leq i \leq n$ .

$$E(\epsilon_i) = 0, \text{ Var}(\epsilon_i) = \sigma^2, \text{ i.e. } Y = (J_n(I - P_n)X) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + \epsilon.$$

We can obtain LSE of  $\alpha$ :  $\hat{\alpha} = P_n Y = \bar{Y}$ .

Prop. C (Criteria of BLP)

$\tilde{\eta}(x)$  is any linear predictor. Then  $\tilde{\eta}(x) = \hat{E}(Y|X)$  a.s.

$\Leftrightarrow \text{Cov}(f(x), Y - \tilde{\eta}(x)) = 0$  for any linear function  $f$  and  $E(\tilde{\eta}(x)) = M_Y$ .

Pf: Set  $f(x) = \eta + (X - M_X)^T \beta$ .

$$\Rightarrow \text{Cov}(f(x), Y - \tilde{\eta}(x)) = \beta^T \Sigma_{XY} - \beta^T \Sigma_{XX} \beta_x = 0$$

$\Leftrightarrow$  Write  $\tilde{\eta}(x) = M_Y + (X - M_X)^T \delta$  for some  $\delta$ .

$$\therefore \text{Cov}(f(x), Y - \tilde{\eta}(x)) = \beta^T \Sigma_{XY} - \beta^T \Sigma_{XX} \delta = 0$$

iv) For vector  $Y = (Y_1, \dots, Y_p)$ :

If we have data  $X = (X_1, \dots, X_2)$

Thm.  $f(x) = (E(Y_1|\vec{x}), \dots, E(Y_p|\vec{x}))$  minimizes:

$$E((Y - f(x))^T (Y - f(x)))$$

Thm.  $f_i(x) = M_{Y_i} + (X - M_X)^T \beta_i^*$ .  $\Sigma_{XY_i} \beta_i^* = \Sigma_{XX}$ . Then

$\vec{f}(x)$  is best linear predictor of  $\vec{Y}$ .

(3) Coefficient of Determination:

① Multiple Correlation Coefficient:

Consider:  $Y = J_n \beta_0 + X_{n \times p} \beta + \varepsilon$ .  $r(X) = p$ .

Written in  $Y = J_n \delta_0 + (X - \frac{J_n}{n}) X \beta + \varepsilon$

Where  $\delta_0 = \beta_0 + \frac{\sum_{i=1}^n x_i}{n} \beta$ .

Def.  $R^2 = \frac{SSR}{SST_{TOT} - C}$  is Coefficient of Determination

where  $SST_{TOT} = Y^T Y$ ,  $C = n \bar{Y}^2$ ,  $SSR = Y^T (M^* - \frac{J_n^T}{n}) Y$

$M^*$  is  $P_{C \perp J_n}$  orth-normal proj.

Rmk: i)  $SYY = SST_{TOT} - C = Y^T (I - P_n) Y = \sum_{i=1}^n (Y_i - \bar{Y})^2$

ii)  $SSR = SYY - SSE$ .

iii)  $R^2 = 1 - \frac{SSE}{SYY}$ ,  $0 \leq R^2 \leq 1$ .  $R^2$  is interpreted

as the proportion of total variability in  $Y$  explained by indep variable.  $(X_1, \dots, X_p)$

It can measure the predictive ability of a model.  $R^2 \uparrow$ . The better fit of model.

To motivate use of  $R^2$  to assess fit:

$R^2$  is actually estimate of the square of multiple correlation coefficient.

Def: The multiple correlation coefficient between  $Y_{(1)}$

and  $\vec{X} = (X_1, \dots, X_p)$  is  $\max_{\text{corr}} \{ \text{Corr}^2(Y, \gamma + X^T \beta) \}$ .

Rmk: It's  $\max_{\beta} \frac{( \sum_{yx} \beta )^2}{\beta^T \sum_{xx} \beta \cdot \sigma_{yy}} = \max_{\beta} \frac{( \sum_{yx} \sum_{xx}^{-\frac{1}{2}} \sum_{xx}^{\frac{1}{2}} \beta )^2}{\sigma_{yy} \cdot \beta^T \sum_{xx} \beta}$

$= \frac{\sum_{yx} \sum_{xx}^{-1} \sum_{xy}}{\sigma_{yy}}$ , So  $\text{MCC}(Y, \vec{X}) = \frac{\sum_{yx} \sum_{xx}^{-1} \sum_{xy}}{\sigma_{yy}}$

Its estimation is  $\frac{S_{YX} S_{XX}^{-1} S_{XY}}{\hat{\sigma}_{YY}^2}$ .  $\hat{\sigma}_{YY}^2 = Y^T (I - P_n) Y = SY$ .

$S_{YX}^2 = Y^T (I - P_n) X$ ,  $S_{XX}^2 = X^T (I - P_n) X$ ,  $S_{XY}^2 = S_{YX}^T$ .

So  $R^2 = \hat{MCC}(Y, \vec{X})$ .

Apply in Test:

Note that  $\frac{SSR}{SSE} = \frac{R^2}{1-R^2}$ . Test:  $H_0: \vec{\beta} = 0$

Then  $F = \frac{Y^T (M^* - \frac{1}{n} J_n) Y / p}{Y^T (I - M^*) Y / (n-p-1)} = \frac{n-p-1}{p} \frac{F^2}{1-R^2} \sim F_{(p, n-p-1)}$

② Partial Correlation Coefficients:

i) We're also interested in the correlation between 2 variables condition on a set of variables that are already in model.

Denote:  $\rho_{Y_1, Y_2}$  is partial correlation coefficient of  $Y_1, Y_2$  given  $X_1, \dots, X_{p-1}$ .

Remark:  $\rho_{Y_1, Y_2}$  is measure of linear relation between  $Y_1, Y_2$  after taking the effect of  $\vec{X}$  out of  $Y_1, Y_2$ .

Note that BLP of  $Y = (Y_1, Y_2)^T$  given  $X$  is:

$\hat{E}(Y|X) = M_Y + \beta_X^T (X - M_X)$ ,  $\Sigma_{XX} \beta_X = \Sigma_{XY} \in M^{(p) \times p}$ .

$\Rightarrow$  Take effects of  $X$  out of  $Y$  by looking at:

the prediction error:  $Y - (M_Y + \beta_X^T (X - M_X))$

Def:  $\rho_{YX}$  is correlation of two components of

$$Y - (M_Y + \beta_X^T (X - M_X)).$$

Rmk:  $\text{Cov}(Y - M_Y - \beta_X^T (X - M_X)) =$

$$\text{Cov}(Y - M_Y) + \beta_X^T \text{Cov}(X - M_X) \beta_X - \square - \square$$

$$= \Sigma_{YY} + \beta_X^T \Sigma_{XX} \beta_X - \Sigma_{YX} \beta_X - \beta_X^T \Sigma_{XY}$$

$$= \Sigma_{YY} - \beta_X^T \Sigma_{XX} \beta_X \quad (\text{Note } \Sigma_{XX} = \Sigma_{XX} \Sigma_{XX}^{-1} \Sigma_{XX})$$

$$= \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY} \Rightarrow \text{Calculate } \rho_{YX}.$$

ii) In sample case:

If we have sample:  $Y = (Y_1, Y_2) = \begin{pmatrix} Y_{11} & Y_{12} \\ \vdots & \vdots \\ Y_{n1} & Y_{n2} \end{pmatrix}$

$X = (X_{ij})_{n \times p}$ . The estimate of  $\Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$  is

$$S_{YY} - S_{YX} S_{XX}^{-1} S_{XY} = \frac{1}{n-1} Y^T (I - M^*) Y$$

$$\text{So } \hat{\rho}_{YX} = r_{Y1X} = \frac{Y_1^T (I - M^*) Y_2}{(Y_1^T (I - M^*) Y_1)^{1/2} (Y_2^T (I - M^*) Y_2)^{1/2}}$$

Def:  $r_{Y1X}$  is sample partial correlation coefficient

iii) Hypothesis Test:

Consider fitting:  $Y_1 = \beta_0 Y_0 + X Y_1 + Y_2 Y_2 + \varepsilon$ .

where  $\varepsilon \sim N(0, \sigma^2 I)$

Since  $C \perp (I - M^*) Y_2$  is  $C^\perp (J_n, X) \cap C (J_n, X, Y_1)$

$$\text{Then: } SSR (Y_2 | J_n, X) = Y_1^T P_{C \perp (I - M^*) Y_2} Y_1$$

$$= Y_1^T (I - M^*) Y_2 (Y_2^T (I - M^*) Y_2)^{-1} Y_2^T (I - M^*) Y_1$$

$$\text{With } SSE (J_n, X) = Y_1^T (I - M^*) Y_1$$

$$\Rightarrow r_{y_1|x}^2 = \frac{SSR (Y_2 | J_n, X)}{SSE (J_n, X)}$$

To test:  $H_0: \beta_{Y_2|X} = 0 \Leftrightarrow H_0: Y_2 = 0 \Leftrightarrow \text{COV}(Y_1, Y_2 | X) = 0$

$$F = \frac{Y_1^T (P_{C \perp (J_n, X, Y_2)} - P_{C \perp (J_n, X)}) Y_1 / 1}{Y_1^T (I - P_{C \perp (J_n, X, Y_2)}) Y_1 / (n - p - 2)}$$

$$= \frac{Y_1^T P_{C \perp ((I - P_{C \perp (J_n, X)}) Y_2)} Y_1}{Y_1^T (I - P_{C \perp (J_n, X)} - P_{C \perp ((I - P_{C \perp (J_n, X)}) Y_2)}) Y_1 / (n - p - 2)}$$

$$= \frac{r_{y_1|x}^2}{(1 - r_{y_1|x}^2) / (n - p - 2)} \underbrace{M_0}_{F(1, n - p - 2)}$$

### ② Squared Predictive Correlation:

i) S.P.C of predictor  $\tilde{y}(x)$  is  $\text{Corr}(Y, \tilde{y}(x))$

The higher S.P.C. The more predictive  $\tilde{y}(x)$  is.

Thm.  $\text{Corr}(Y, \tilde{y}(x)) \leq \text{Corr}(Y, m(x)) =: R^2$

Pf:  $\sigma_{Y\tilde{y}}^2 = \sigma_{m\tilde{y}}^2 \leq \sigma_{mm} \sigma_{\tilde{y}\tilde{y}}$ . By Schwarz.

$$\text{And } R^2 = \frac{\sigma_{Ym}^2}{\sigma_{Y1} \sigma_{mm}} = \frac{\sigma_{mm}^2}{\sigma_{Y1} \sigma_{mm}} = \sigma_{mm} / \sigma_{Y1}$$

Rmk: In fact, high SPC can also be attained by bad predictors. Since  $\tilde{y}(x)$  is high correlated with  $Y$  doesn't mean  $\tilde{y}(x)$  is close enough to  $Y$ .

## ii) Linearized Predictor:

If we have a predictor  $\tilde{y}(x)$ . We can construct a linear predictor which is at least as good as

$$\tilde{y}(x) : \hat{y}(x) = \mu_Y + \frac{\sigma_{Y\tilde{y}}}{\sigma_{\tilde{y}\tilde{y}}} (\tilde{y}(x) - \mu_{\tilde{y}}) \quad \text{i.e.}$$

the BLP based on  $\tilde{y}(x)$  rather than  $X$ .

Rmk: Note  $\sigma_{\hat{y}\hat{y}} = \sigma_{\tilde{y}\tilde{y}} / \sigma_{\tilde{y}\tilde{y}} = \sigma_{Y\hat{y}}$ . We obtain:

$$\text{Corr}^2(Y, \hat{y}(x)) = \sigma_{\hat{y}\hat{y}} / \sigma_{YY}$$

Def: We measure the goodness of such pred.  $\hat{y}(x)$

$$\text{by } : E(Y - \hat{y}(x))^2 = \sigma_{YY} - \sigma_{\hat{y}\hat{y}}$$

prop. For two linearized predictor  $\hat{y}_1(x), \hat{y}_2(x)$ .

$\hat{y}_2(x)$  is better  $\Leftrightarrow$  SPC of  $\hat{y}_2(x)$  is higher

Pf:  $\text{Corr}^2(Y, \hat{y}_1(x)) \leq \text{Corr}^2(Y, \hat{y}_2(x)) \Leftrightarrow$

$$\sigma_{\hat{y}_1\hat{y}_1} / \sigma_{YY} \leq \sigma_{\hat{y}_2\hat{y}_2} / \sigma_{YY} \Leftrightarrow$$

$$\sigma_{YY} - \sigma_{\hat{y}_1\hat{y}_1} \leq \sigma_{YY} - \sigma_{\hat{y}_2\hat{y}_2} \quad \text{by Def.}$$