

## Simple Linear Regression

- same**
- The distribution of  $X$  is arbitrary (and perhaps  $X$  is even non-random)
  - $Y = \beta_0 + \beta_1 X_1 + \varepsilon$ . If  $X_1 = x_1$ , then  $Y = \beta_0 + \beta_1 x_1 + \varepsilon$  for some coefficients  $\beta_0, \beta_1$  and random noise  $\varepsilon$ .
  - $\varepsilon \sim N(0, 1)$ ,  $E(\varepsilon) = 0$ ,  $\text{Var}(\varepsilon) = \sigma^2$

④  $\varepsilon$  is independent of  $x_1$  and independent across observations.  $\varepsilon$  is uncorrelated with  $x_1$ . Linear correlation  $\Rightarrow$  those having 0 correlation might dependent non-linearly.

## Optimal Linear Predictor

$$(\hat{\beta}_0^*, \hat{\beta}_1^*) = \underset{(\beta_0, \beta_1)}{\operatorname{arg\min}} E_{x,y} [(Y - (\beta_0 + \beta_1 x))^2]$$

$$\begin{cases} \hat{\beta}_0^* = E(Y) - \hat{\beta}_1^* E(x) \\ \hat{\beta}_1^* = \frac{\text{Cov}(x, Y)}{\text{Var}(x)} \end{cases} \Rightarrow \begin{cases} \hat{\beta}_0 = \bar{y} - \frac{s_{xy}}{s_{xx}} \bar{x} \\ \hat{\beta}_1 = \frac{s_{xy}}{s_{xx}} \end{cases}$$

## Least Square Estimates

$$S(LP) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{ii})^2 = \frac{1}{n} (Y - X\hat{\beta})^T (Y - X\hat{\beta})$$

$$\begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \begin{bmatrix} n & \sum_{i=1}^n x_{ii} \\ \sum_{i=1}^n x_{ii} & \sum_{i=1}^n x_{ii}^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_{ii} y_i \end{bmatrix} = (X^T X)^{-1} X^T Y$$

$$E(\hat{\beta}_1) = \hat{\beta}_1, E(\hat{\beta}_0) = \hat{\beta}_0, \text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{s_{xx}}, \text{Var}(\hat{\beta}_0) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}}{s_{xx}} \right)$$

$$\begin{cases} \hat{\beta}_1 = \frac{\sum_{i=1}^n x_{ii} - \bar{x}}{s_{xx}} y_i - \sum_{i=1}^n c_i = 0 \\ \hat{\beta}_0 = \sum_{i=1}^n \left( \frac{1}{n} - \frac{\bar{x}}{s_{xx}} \right) y_i - \sum_{i=1}^n b_i = 1 \\ \sum_{i=1}^n b_i x_i = 0 \end{cases} \quad \hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\begin{aligned} \text{Var}(\hat{\beta}|X) &= A \text{Var}(Y|X) A^T = A \sigma^2 I A^T = \sigma^2 A A^T = \sigma^2 (X^T X)^{-1} \\ \begin{bmatrix} \text{Var}(\hat{\beta}_0|X) & \text{Cov}(\hat{\beta}_0, \hat{\beta}_1|X) \\ \text{Cov}(\hat{\beta}_1, \hat{\beta}_0|X) & \text{Var}(\hat{\beta}_1|X) \end{bmatrix} &= \frac{\sigma^2}{n s_{xx}} \begin{bmatrix} \sum_{i=1}^n x_{ii}^2 & -\sum_{i=1}^n x_{ii} \\ -\sum_{i=1}^n x_{ii} & n \end{bmatrix} \end{aligned}$$

① The Law of Total Expectation

$$E(\hat{\beta}_1) = E [E_{Y|X} (\hat{\beta}_1 | X_1, \dots, X_M)] = E(\beta_1) = \beta_1$$

② The Law of Total Variance

$$\text{Var}(\hat{\beta}_1) = E [ \text{Var}_{Y|X} (\hat{\beta}_1 | X_1, \dots, X_M) ] + \text{Var} [ E_{Y|X} (\hat{\beta}_1 | \dots) ]$$

## Multiple Linear Regression [Multiple Regression]

① There are  $k$  predictors  $X_1, X_2, \dots, X_k$ . We make no assumption about their distribution. They may or may not be dependent.

② There is a single response  $Y$ .

③ Linear model  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon = X\beta + \varepsilon$  for some constants  $\beta_0, \beta_1, \dots, \beta_k$

④ The noise has  $E(\varepsilon) = E(\varepsilon|X) = 0$  and  $\text{Var}(\varepsilon) = \text{Var}(\varepsilon|X) = \sigma^2$

$\beta_0$ : The expected value of  $Y$  when  $X_1 = \dots = X_k = 0$   
 $\beta_j$ : If we select two sets of the cases from the distribution of the data, where  $x_j$  differs by 1, we expect  $Y$  to differ by  $\beta_j$  on average

$$\hat{\beta} = (X^T X)^{-1} X^T Y \quad \text{① } n \geq p$$

② The predictors must be linearly independent (Check column rank of  $X$ :  $k \leq p$ )

$$E(\hat{\beta}) = \beta, \text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1} \quad \hat{\sigma}^2 = \frac{\|Y - X\hat{\beta}\|^2}{n-p}$$

$$\begin{cases} \varepsilon \sim N(0, \sigma^2 I_n) \\ \varepsilon \text{ is independent of } X \end{cases} \Rightarrow Y|X \sim N(X\beta, \sigma^2 I_n) \quad \hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$$

① Estimate

$$\hat{\beta}_j \pm K \sqrt{\hat{\sigma}^2 \cdot C_{jj}}$$

$$C = (X^T X)^{-1}, K = t_{\alpha/2} \sqrt{n-p}$$

② Mean Response

$$\hat{m}(x_0) \pm K \cdot \sqrt{\hat{\sigma}^2 \cdot x_0 (X^T X)^{-1} x_0^T}$$

③ New Observation

$$\hat{m}(x_0) \pm K \cdot \sqrt{\hat{\sigma}^2 (I + x_0 (X^T X)^{-1} x_0^T)}$$

## Residuals

$$\begin{aligned} e_i &= y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{ii} \\ \sum_{i=1}^n e_i x_{ii} &= 0 \\ \sum_{i=1}^n e_i^2 &= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{ii})^2 = S_{\text{Res}} \\ \sum_{i=1}^n e_i &= 0 \end{aligned}$$

$$\hat{\sigma}^2 = \frac{S_{\text{Res}}}{n-p} = M S_{\text{Res}}$$

residual mean square

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

H is symmetric  
H is idempotent  $H^T H = H$   
 $I_n - H$  is idempotent

$$S_{\text{Res}} = (Y - X\hat{\beta})^T (Y - X\hat{\beta}) = Y^T (I_n - H) Y$$

If  $v$  is a random vector,  $E(v) = \mu$ ,  $\text{Var}(v) = \Sigma$

For a constant matrix  $A$ ,  $E(v^T A v) = \text{trace}(A\Sigma) + \mu^T A \mu$

$$E(S_{\text{Res}}|X) = \sigma^2 (n - \text{trace}(H))$$

$$\text{trace}(H) = \text{trace}(X(X^T X)^{-1} X^T) = \text{trace}(X^T X (X^T X)^{-1}) = p$$

standard error (se) =  $\sqrt{\text{Var}(\hat{\beta}_j)}$

$$\text{estimated standard error (ese)} = \begin{cases} \text{ese}(\hat{\beta}_0) = \sqrt{\frac{\sigma^2}{n} + \frac{\bar{x}^2}{s_{xx}}} \\ \text{ese}(\hat{\beta}_1) = \sqrt{\frac{\sigma^2}{s_{xx}}} \end{cases}$$

## AN-SLR

$$\text{① } \hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{s_{xx}}) \quad \hat{\beta}_0 \sim N(\beta_0, \sigma^2 (\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}})) \quad \frac{(n-2)\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-p}$$

$$\text{② } T_1 = \frac{\hat{\beta}_1 - \beta_1}{\text{ese}(\hat{\beta}_1)} \sim t_{n-p} \quad T_0 = \frac{\hat{\beta}_0 - \beta_0}{\text{ese}(\hat{\beta}_0)} \sim t_{n-p}$$

A  $1-\alpha$  confidence interval  $CI(\beta_1) = [\hat{\beta}_1 \pm k \cdot \text{ese}(\hat{\beta}_1)]$

$$\text{③ } \begin{cases} H_0: \beta_1 = c & \text{We reject } H_0 \text{ if } |T_1| \geq k \equiv t_{\alpha/2, n-p} \\ H_1: \beta_1 \neq c & T_1 = \frac{\hat{\beta}_1 - c}{\text{ese}(\hat{\beta}_1)} \end{cases}$$

When  $H_0$  is attained  $\Rightarrow \beta_1$  is statistically undistinguishable from 0

$$\text{④ } CI(m(x_0)) = [\hat{m}(x_0) \pm k \sqrt{MS_{\text{Res}} (\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{s_{xx}}) }]$$

$$P(Y) = [\hat{m}(x_0) \pm k \sqrt{MS_{\text{Res}} (1 + \frac{(x_0 - \bar{x})^2}{s_{xx}}) }]$$

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SS_{\text{Res}}}{SST} \rightarrow \hat{\beta}_1 s_{xy} = \hat{\beta}_1 s_{xx}$$

$$R^2 = 0 \text{ if } \hat{\beta}_1 = 0, R^2 = 1 \text{ if } SS_{\text{Res}} = 0$$

$$\frac{\hat{\beta}_1^2 s_{xx}}{\hat{\beta}_1^2 s_{xx} + SS_{\text{Res}}} = \frac{\hat{\beta}_1^2 \frac{s_{xx}}{n}}{\hat{\beta}_1^2 \frac{s_{xx}}{n} + \frac{SS_{\text{Res}}}{n}} \rightarrow \frac{\hat{\beta}_1^2 \text{Var}(x)}{\hat{\beta}_1^2 \text{Var}(x) + \sigma^2}$$

① We can manipulate the data that

$$\begin{cases} R^2 = 0 \text{ but the model is correct (Var}(x) \text{ small}) \\ R^2 = 1 \text{ but the model is wrong (Var}(x) \text{ large}) \end{cases}$$

② only can be compared with same dataset with same, untransformed response variable.

$$R^{\text{adj}} = 1 - \frac{SS_{\text{Res}} / (n-p)}{SST / (n-1)} \quad \begin{aligned} R^2 &\rightarrow 1 \text{ when } P \rightarrow n \\ R^{\text{adj}} &= R^2 \text{ when } P=1 \end{aligned}$$

③ Low  $R^2$  and low p-value ( $p\text{-value} < \alpha$ )

The model doesn't explain much of variation of the data but the model is significant

(Including  $x_1$  into the model is necessary but not sufficient)

④ Low  $R^2$  and high p-value

The Model doesn't explain much of the data and  $x_1$  is not useful. (worse)

⑤ High  $R^2$  and low p-value

Model is sufficient,  $x_1$  is useful (best)

**F-test**  $y^T(I_{n-H})y = y^T(I_{n-H})y + y^T(H-H)y \quad (H_H = \frac{1}{n} \mathbf{1} \mathbf{1}^T)$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad \text{"systematic variation"}$$

in  $y$  due to  $x$ -linear aggregate measure of malfitted regression line relationship

 $SST / n^2 \sim \chi^2_{n-1} \quad SS_{\text{Res}} / n^2 \sim \chi^2_{n-p} \quad SSR / n^2 \sim \chi^2_{p-1}$ 

**ANOVA** reject null if  $F_0 > F_{\alpha, p-1, n-p}$  ( $P < \alpha$ )

Source	SS	df	MS	F
Regression	SSR	p-1	MSR = $\frac{SSR}{n-p}$	$\frac{SSR / (p-1)}{SS_{\text{Res}} / (n-p)}$
Residual	SS_{\text{Res}}	n-p	MS_{\text{Res}} = $\frac{SS_{\text{Res}}}{n-p}$	
Total	SST	n-1		

$(T_0)^2 = F_0 \quad [\text{P-value same } P(F > F_0) = P(T > T_0)]$

1. If retain  $H_0: \beta_1 = 0$  (don't find any significant share of variance associated with the regression).

- (a) The intercept-only model is correct
- (b)  $\beta_1 \neq 0$ , but the test don't have enough power to detect it  $\hat{\beta}_1$
- (c) The relationship is non-linear, but the best linear approximation has zero slope.

2. If reject  $H_0$  (doesn't mean the SLR model is right, only that the non-zero slope linear model predict better than the intercept-only models).

## Model Adequacy Checking

$\text{residual } e = y - \hat{y} = (I_{n-H})y$

- ① The distribution of  $e$  should have a center around 0.  $E[e|X] = 0_n$ .
- ② The covariance between  $\hat{y}$  and  $e$  is zero.  $\text{Cov}(\hat{y}, e) = 0_{n \times n}$ .
- ③ The variance of  $e$  is roughly constant.  $\text{Var}(e|X) = \sigma^2 (I_{n-H})$
- ④ MLR:  $e$  is Gaussian with mean 0 and variance  $\sigma^2 I_n$ .  $e_i \sim N(0, \sigma^2)$
- ⑤ quantile and quantile plot:  $F^{-1}(p) = \Phi^{-1}(p) + \mu$ .  $F \sim N(\mu, \sigma^2)$ .  $\Phi \sim N(0, 1)$
- $\Rightarrow \hat{F}^{-1}\left(\frac{i}{n}\right) = e_{(i)} \sim \Phi^{-1}\left(\frac{i}{n}\right)$

## Transformation

- ① Variance-stability transformation  $\Rightarrow \begin{cases} E[y] = \bar{y} & \text{Poisson} \\ E[y^2] = E[y] + E[y^2] - E[y]^2 & Y = \sin(Y) \\ E[y^2] = E[y]^2 & \text{Binomial} \end{cases}$
- ② linearity

**Leverage**  $\hat{y}_{ii} = \sum_{j=1}^n h_{ij} \hat{y}_j$  the amount of leverage exerted by the  $i$ th observation  $y_i$  on the  $i$ th fitted value  $\hat{y}_i$ .

$h_{ii} = x_{ii} (X^T X)^{-1} x_{ii}^T \quad \text{large } h_{ii} \Rightarrow i\text{th observation has high leverage}$

## Influence

**Cook's Distance**  $D_i = \frac{(\hat{y}_{(i)} - \hat{y})^T M (\hat{y}_{(i)} - \hat{y})}{P \cdot MS_{\text{Res}}} \quad \begin{matrix} \xrightarrow{\text{fitted estimate with}} \\ \text{ith observation removed} \end{matrix}$

$= \frac{e_i^2}{P} \frac{h_{ii}}{(1-h_{ii})^2} \frac{1}{P}$

(standard residual)<sup>2</sup> leverage

$\hat{y}_{(i)} - \hat{y} = \hat{y}_{(i)} - \hat{y} \Rightarrow D_i = \frac{(\hat{y}_{(i)} - \hat{y})^T (\hat{y}_{(i)} - \hat{y})}{P \cdot MS_{\text{Res}}}$

Squared distance that the fitted value moves when the  $i$ th observation is removed.

## Model Selection

$\text{generalization error } G = E[(Y - \hat{m}(X))^2]$

$G = E\left[\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{y}_i)^2\right] \quad \text{testing error}$

$= E\left[\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{y}_i)^2\right] + \frac{2}{n} \sigma^2 P \quad \text{training error } \frac{MS_{\text{Res}}}{n}$

- ① Mallow's CP ↓

$C_P = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{y}_i)^2 + \frac{2}{n} \hat{\sigma}^2 P$

- ② adjusted R<sup>2</sup> ↑

$R^2_{\text{adj}} = 1 - \frac{MSE \cdot \frac{n}{n-P}}{SST / n} = MSE + MSE \cdot \frac{P}{n} \xrightarrow{n \gg P} MSE + \sigma^2 \cdot \frac{P}{n} = G$

- ③ AIC & BIC ↓

$AIC = n \cdot \ln(MSE) + 2P = n \cdot \ln(\frac{SS_{\text{Res}}}{n}) + 2P$

$BIC = n \cdot \ln(MSE) + \ln(n) \cdot P$

$\left\{ \begin{array}{l} H_0: \beta_{01} = 0 \text{ or } x_1 \\ H_1: \beta_{01} \neq 0 \text{ or } x_1 \end{array} \right. \Rightarrow \left\{ \begin{array}{l} H_0: Y = x_{01} \beta_{01} + \varepsilon \text{ reduced model} \\ H_1: Y = x_{01} \beta_{01} + x_{11} \beta_{11} + \varepsilon \text{ full model} \end{array} \right.$

$F_0 = \frac{SSR(B_0) / (P-1)}{SSR(B) / (n-P)}$

$SSR(B) - SSR(B_0)$

## Sequential F test

$\left\{ \begin{array}{l} H_0: Y = \beta_0 + \varepsilon \\ H_1: Y = \beta_0 + \beta_1 x_1 + \varepsilon \end{array} \right. \Rightarrow F_1 = \frac{SSR(B_1) / (P-1)}{SSR(B_0) / (n-1)}$

$\left\{ \begin{array}{l} H_0: Y = \beta_0 + \beta_1 x_1 + \varepsilon \\ H_1: Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon \end{array} \right. \Rightarrow F_2 = \frac{SSR(B_2) / (P-2)}{SSR(B_1) / (n-2)}$

$\overline{SSR}(B_3, B_2, B_1, B_0) = \overline{SSR}(B_3, B_2, B_1, B_0) - \overline{SSR}(B_0)$

$= \overline{SSR}(B_1, B_0) + \overline{SSR}(B_2, B_1, B_0) + \overline{SSR}(B_3, B_2, B_1, B_0)$

$\text{global model } \overline{SSR}(B_{\text{full}}) = \overline{SSR}(B)$

$\overline{SST} = \overline{SS_{\text{Res}}} + \overline{SSR}$ 
 $\sum y_i^2 = \sum (y_i - \hat{y}_i)^2 + \sum \hat{y}_i^2 \Rightarrow \overline{SSR} = \frac{1}{n} \sum y_i^2 = \hat{\beta}^T \bar{X}^T Y$

## Penalization

$\min_B \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}))^2 \text{ such that } \|B\|^2 = \beta_1^2 + \beta_2^2 + \dots + \beta_k^2 \leq t \quad \text{ridge penalization}$

$\Rightarrow \min_B \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_k x_{ik})^2 + \lambda \|B\|^2$

$\text{no-intercept model } \hat{B} = (X^T X + \lambda I)^{-1} X^T Y$

① Control model complexity

② Deal with high collinearity between predictors

③ When  $n < p$  (MLR will fail)

## Factor Predictor

$\text{with intercept } X^{(1)} = \begin{cases} 1 & x=l \\ 0 & \text{otherwise} \end{cases} \quad \text{for } l=0 \dots M-1$ 

baseline level:  $l=M$

$E[Y|X_l] = \beta_0 + \sum_{i=1}^M P_{il} X_i^{(1)} \quad \begin{matrix} \nearrow P_l = E[Y|X=1] - E[Y|X=M] \\ \text{effect contrast between levels } 1 \dots M-1 \end{matrix}$ 
 $\beta_0 = E[Y|X=M] \quad \text{baseline level effect}$

$\text{no intercept } X^{(1)} = \begin{cases} 1 & x=l \\ 0 & \text{otherwise} \end{cases} \quad \text{for } l=0 \dots M$

$E[Y|X_l] = \sum_{i=1}^M P_{il} X^{(1)} \quad P_{il} = \begin{cases} \beta_0 & l=M \\ \beta_l + \beta_0 & l=1 \dots M-1 \end{cases} \quad \begin{matrix} \text{effect of level } M-1 \dots 1 \\ m=1 \dots M \end{matrix}$

## Interaction

$\text{continuous } E[Y|X] = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 \cdot X_2 \quad \text{interaction}$

modification of the effect of outcome of  $X_1$  on the presence of  $X_2$  (LVV)

$\text{factor } \text{① } E[Y|X] = \beta_0 + \sum_{m=1}^{M-1} \beta_m X^{(m)} + \sum_{k=1}^{M-1} \beta_{1k} X^{(1)} + \sum_{m=1}^{M-1} \sum_{k=1}^{M-1} \beta_{mk} X^{(m)} X^{(k)}$

$\text{② } E[Y|X] = \beta_0 + \sum_{m=1}^{M-1} \beta_m X^{(m)} + \sum_{k=1}^{M-1} \beta_{1k} X^{(1)} + \sum_{m=1}^{M-1} \sum_{k=1}^{M-1} \beta_{mk} X^{(m)} X^{(k)}$

$\text{Var}(X \pm Y) = \text{Var}(X) + \text{Var}(Y) \pm 2 \text{Cov}(X, Y)$

$E[X^2] = [E(X)]^2 + \text{Var}(X)$

$E(XY) = \text{Cov}(X, Y) + E(X)E(Y)$

$\text{Cov}(X, Y) = E[(X - E(X))(Y - E(Y))]$

$S_{XX} = \sum_{i=1}^n (X_{ii} - \bar{X})^2 = \sum_{i=1}^n (X_{ii} - \bar{X}) X_{ii} = \sum_{i=1}^n X_{ii}^2 - n \bar{X}^2$

$S_{XY} = \sum (X_{ii} - \bar{X})(Y_{ii} - \bar{Y}) = \sum (X_{ii} - \bar{X}) Y_{ii} = \sum X_{ii} Y_{ii} - n \bar{X} \bar{Y}$