

# STAT 331



# Personal Notes

---

Marcus Chan

Taught by Peter Balka  
UW CS '25

---



# Chapter 1: Introduction

## REGRESSION

Q: In regression modelling, we attempt to explain or account for variation in a response variate ( $y$ ) by using a model to describe the relationship between  $y$  and one or more explanatory variates ( $x_1, x_2, \dots$ )

## SUMMARIES OF THE DATA

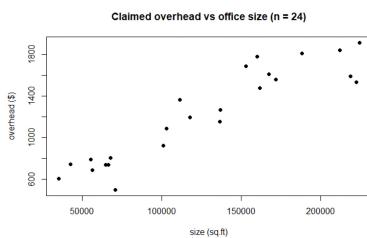
Q: A simple LR model involves:

- ① A single explanatory variate;
- ② A single response variate.

e.g Overhead data example:

response ( $y$ ): claimed overhead (\$)  
explanatory ( $x$ ): office size (sq.ft)

Q: We can summarise the data using a scatter-plot.



Q: To get a numerical summary of the data, we can use the "sample correlation coefficient".

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} = \frac{s_{xy}}{\sqrt{s_{xx}s_{yy}}}$$

Note  $-1 \leq r \leq 1$  and that  $r$  is unitless.

Q:  $r$  tells us the relative strength of the linear relationship.

## THE SIMPLE LR MODEL

Q: we can describe the observed behavior of the response with a model that includes both

- ① a "deterministic component" that describes the variation in  $y$  accounted for by the functional form of the underlying relationship between  $y$  &  $x$ ; &

e.g with the overhead data, the det. comp. is

$$\mu = \beta_0 + \beta_1 x.$$

where  $\mu$  = the mean value of  $y$  for a given value of  $x$ .

- ② an "error term"  $\epsilon$  that describes the random variation in  $y$  not accounted for by the underlying relationship with  $x$ .

Putting these together yields the "simple LR" (or SLR) model:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

where

- ①  $\beta_0$  = the "intercept" parameter

- ②  $\beta_1$  = the "slope" parameter

- ③  $i$  = the index that denotes the observation number.

( $x_1, \dots, x_n$  is explanatory data;  $y_1, \dots, y_n$  is response data).

\* note  $\beta_0 + \beta_1 x_i$  is deterministic &  $\epsilon$  is random.

# THE NORMAL SLR MODEL

We typically assume in SLR that

$$\epsilon_i \sim \text{iid } N(0, \sigma^2), \quad i=1, \dots, n$$

for some variance  $\sigma^2$ .

This yields the normal model

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \sim \text{iid } N(0, \sigma^2).$$

Assumptions needed to use this model:

- ① the functional form (ie linear) of the relationship between  $y$  &  $x$  is correctly specified by the deterministic component of the model;
- ② errors follow a normal distribution;
- ③ errors have a constant variance  $\sigma^2$  (ie "homoskedasticity"); &
- ④ errors are independent.

## LEAST SQUARES ESTIMATION OF MODEL PARAMETERS

Goal: we want to find values of  $\beta_0$  &  $\beta_1$  such that for the data

$$y_1 = \beta_0 + \beta_1 x_1 + \epsilon_1$$

⋮

$$y_n = \beta_0 + \beta_1 x_n + \epsilon_n,$$

the sum of squares of the errors  $\sum \epsilon_i^2$  is minimized.

The values of  $\beta_0$  &  $\beta_1$  obtained by this procedure (denoted  $\hat{\beta}_0$  &  $\hat{\beta}_1$ ) are known as the "least squares estimates" of  $\beta_0$  &  $\beta_1$ .

We show that

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}.$$

Proof. we wish to minimize

$$S(\beta_0, \beta_1) = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2.$$

See that

$$\frac{\partial S}{\partial \beta_0} = -2 \sum_{i=1}^n [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]$$

$$\frac{\partial S}{\partial \beta_1} = -2 \sum_{i=1}^n x_i [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)].$$

Since we want to minimize  $S$ , we can solve

$$\left| \begin{array}{l} \frac{\partial S}{\partial \beta_0} = 0 \\ \frac{\partial S}{\partial \beta_1} = 0. \end{array} \right.$$

The resultant solutions for  $\beta_0$  &  $\beta_1$  are the desired values as required.  $\blacksquare$

In R, we can get these values via

```
> data.lsr.lm <- lm(response ~ explanatory).
```

## FITTED MODEL

For the SLR model, the fitted model is

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x,$$

where  $\hat{y}$  is the estimated mean value of  $y$  given a value of  $x$ .

## FITTED RESIDUALS

The "fitted residual" of the  $i^{th}$  observation,  $e_i$ , is defined as

$$e_i = y_i - \hat{y}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i).$$

\*  $e_i$  is a random variable in which we impose assumptions;

$e_i$  is the difference between the observed response & estimated mean response.

If we take the partial derivative wrt each parameter and set = 0 in our least squares procedure, we get that

$$\sum e_i = 0$$

$$\sum x_i e_i = 0.$$

These constraints allow us to calculate the remaining 2 residuals from  $n-2$  observations;

so we say the fitted model is associated with  $n-2$  degrees of freedom.

# LEAST SQUARES ESTIMATE OF $\sigma^2$ : $\hat{\sigma}^2$

$\textcircled{B}_1$  In the normal model, we assume

$$\varepsilon_i \sim N(0, \sigma^2).$$

$\textcircled{B}_2$  In any least squares regression model, we estimate  $\sigma^2$  by dividing the sum of squares of the residuals by the degrees of freedom.

$\textcircled{B}_3$  In particular, this means our estimate for  $\sigma^2$  is

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{n-2} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}.$$

\* note  $E[\hat{\sigma}^2] = \sigma^2$  (ie  $\hat{\sigma}^2$  is unbiased).

## RESIDUAL STANDARD ERROR: $\hat{\sigma}$

$\textcircled{B}_1$  The "residual standard error" is

$$\hat{\sigma} = \sqrt{\hat{\sigma}^2} = \sqrt{\frac{\sum e_i^2}{n-2}}$$

$\textcircled{B}_2$   $\hat{\sigma}$  can be interpreted as the estimated std dev of the errors & measures the random variation of the response given a value for  $x$ .

$\textcircled{B}_3$  The smaller  $\hat{\sigma}$  is, the more the variation in  $y$  is "explained" by  $x$ , and so the better fit the model is.

$\textcircled{B}_4$   $\hat{\sigma}$  is part of the summary R output for the fitted model:

```
> summary(audit.lm)
Call:
lm(formula = overhead ~ size)
Residuals:
    Min     1Q Median     3Q    Max 
-36639 -12874 -1997   8642  56686 
Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -27877.06    14172.00   -1.967   0.0619 .  
size         126.33      10.88    11.610 7.47e-11 *** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
Residual standard error: 23480 on 22 degrees of freedom
Multiple R-squared:  0.8597, Adjusted R-squared:  0.8533 
F-statistic: 134.8 on 1 and 22 DF,  p-value: 7.472e-11
```

# INTERPRETATION OF PARAMETER ESTIMATES

$\textcircled{B}_1$  We may interpret  $\hat{\beta}_0$  as the estimated mean change in the response  $y$  associated with a change of one unit in  $x$ .

$\textcircled{B}_2$  we may interpret  $\hat{\beta}_0$  as the estimated mean value of  $y$  at  $x=0$  only if  $x=0$  is a relevant value and is in the range of values we used to fit the model.

\* never extrapolate to values of  $x$  outside the range used to fit the model.

$\textcircled{B}_3$  Lastly, we can interpret  $\hat{\sigma}$  as a measure of the variability of the response about the fitted line.

## INFERENCE FOR $\beta_1$

$\textcircled{B}_1$  To investigate whether there is a linear relationship between  $y$  &  $x$  in the population, we can test the hypothesis  $\beta_1 = 0$ .

$\textcircled{B}_2$  We can then either use confidence intervals or hypothesis tests to test this.

$\textcircled{B}_3$  To do this, we need the least squares estimator of  $\beta_1$ :

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

# DISTRIBUTION OF $\hat{\beta}_1$

Q: we can show for the SLR model that

$$\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{S_{xx}})$$

Proof. First, note

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2} \\ &= \frac{\sum(x_i - \bar{x})y_i - \bar{y}\sum(x_i - \bar{x})}{\sum(x_i - \bar{x})^2} \\ &= \frac{\sum(x_i - \bar{x})y_i}{\sum(x_i - \bar{x})^2} \quad \because \sum(x_i - \bar{x}) = 0 \\ &= \sum c_i y_i, \quad c_i = \frac{x_i - \bar{x}}{\sum(x_i - \bar{x})^2}.\end{aligned}$$

Then, for the SLR model,  $e_i \sim N(0, \sigma^2)$ .

Since  $y_i = \beta_0 + \beta_1 x_i + e_i$ , thus

$$y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2) \quad \& \quad y_i \text{ are ind.}$$

and so

$$\hat{\beta}_1 = \sum c_i y_i \sim \text{Normal} \quad (\text{by properties of normal}).$$

Then

$$\begin{aligned}E(\hat{\beta}_1) &= E(\sum c_i y_i) = \sum c_i E(y_i) \\ &= \sum \frac{x_i - \bar{x}}{\sum(x_i - \bar{x})^2} \cdot (\beta_0 + \beta_1 x_i) \\ &= \frac{\beta_0 \sum(x_i - \bar{x}) + \beta_1 \sum x_i(x_i - \bar{x})}{\sum(x_i - \bar{x})^2} \\ &= \frac{\beta_1 \sum x_i(x_i - \bar{x})}{\sum(x_i - \bar{x})^2} \\ &= \frac{\beta_1 \sum x_i(x_i - \bar{x}) - \beta_1 \sum \underbrace{x_i}_{0}(x_i - \bar{x})}{\sum(x_i - \bar{x})^2} \\ &= \frac{\beta_1 \sum(x_i - \bar{x})^2}{\sum(x_i - \bar{x})^2} = \beta_1.\end{aligned}$$

Similarly,

$$\begin{aligned}\text{Var}(\hat{\beta}_1) &= \text{Var}(\sum c_i y_i) \\ &= \sum c_i^2 \text{Var}(y_i) \quad \because y_i \text{ is ind.} \\ &= \sum \frac{(x_i - \bar{x})^2}{(\sum(x_i - \bar{x})^2)^2} \cdot \sigma^2 \\ &= \frac{\sigma^2}{\sum(x_i - \bar{x})^2} = \frac{\sigma^2}{S_{xx}}.\end{aligned}$$

Hence  $\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{S_{xx}})$  as required.

Q: It follows that

$$\frac{\hat{\beta}_1 - \beta_1}{\text{SE}(\hat{\beta}_1)} = \frac{\hat{\beta}_1 - \beta_1}{(\hat{\sigma}/\sqrt{S_{xx}})} \sim t_{n-2}.$$

(from STAT231/330 result).

This can be used to get  $t$ -based CIs & hypothesis tests for  $\beta_1$ .

# DISTRIBUTION OF $\hat{\beta}_0$

Q: Similarly, we can show in a SLR model,

$$\hat{\beta}_0 \sim N(\beta_0, \sigma^2(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}))$$

$$\frac{\hat{\beta}_0 - \beta_0}{\text{SE}(\hat{\beta}_0)} = \frac{\hat{\beta}_0 - \beta_0}{\hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}} \sim t_{n-2}.$$

## CI FOR $\beta_1$

Q: A  $(1-\alpha)100\%$  confidence interval for  $\beta_1$  is

$$\hat{\beta}_1 \pm t_{n-2, 1-\alpha/2} \text{ SE}(\hat{\beta}_1), \quad \text{SE}(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{S_{xx}}}.$$

$t_{n-2, 1-\alpha/2}$  := the critical value from a  $t_{n-2}$  distribution corresponding to a confidence level of  $(1-\alpha)100\%$ .

Q: " $t_{n-2, 1-\alpha/2} \text{ SE}(\hat{\beta}_1)$ " is called the "margin of error" of the interval.

Q: We can use R to calculate this:

> summary(data.lm)

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-27877.06	14172.00	-1.967	0.0619
size	126.33	10.88	11.610	7.47e-11

>  $t \leftarrow qt(1-\frac{\alpha}{2}, n-2)$

↳ The CI is then  $126.33 - t(0.88)$ ,  $126.33 + t(0.88)$ .

Q: We may interpret the CI as that we are  $(1-\alpha)100\%$  confident that for every additional increase of a unit of  $x$ , the mean increase of  $y$  is between (start of CI) & (end of CI).

Q: If " $\beta_1 = 0$ " is not in the interval, then we say there is a significant relationship between  $x$  &  $y$  at the  $(1-\alpha)100\%$  confidence level.

Q: Hypothesis test for  $\beta_1$ :

- ①  $H_0: \beta_1 = 0$ ;  $H_A: \beta_1 \neq 0$
- ② (Assuming  $H_0$ ) our test statistic is

$$t = \frac{\hat{\beta}_1 - \beta_1}{\text{SE}(\hat{\beta}_1)} = \frac{\hat{\beta}_1}{\text{SE}(\hat{\beta}_1)}$$

- ③ p-value is  $p = 2P(T > t)$ ,  $T \sim t_{n-2}$ 
  - In R:  
>  $p \leftarrow 2 * (1 - pt(t, n-2))$
- ④ Check if  $p < 0.05$ ; if yes, reject  $H_0$ .

# Chapter 2:

## Multiple Regression

### MULTIPLE REGRESSION MODEL

If we expand the SLR model to  $p$  explanatory variables, we obtain the multiple linear regression model:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i, \quad i=1, \dots, n$$

This can be expressed as

$$y = X\beta + \epsilon$$

where  $y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$ ,  $X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} \in \mathbb{R}^{n \times (p+1)}$ ,

 $\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \in \mathbb{R}^{p+1}$  &  $\epsilon = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix} \in \mathbb{R}^n$ .

### NORMAL MODEL

For the normal model, where we assume  $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ , we write

$$y = X\beta + \epsilon, \quad \epsilon \sim MVN(0, \sigma^2 I),$$

where  $\text{Var}(\epsilon) = \sigma^2 I$  is the covariance matrix of the error random vector  $\epsilon$ .

### LEAST SQUARES ESTIMATION OF $\beta$

We wish to minimize

$$S(\beta_0, \dots, \beta_p) = \sum_{i=1}^n \epsilon_i^2 = \sum [y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip})]^2$$

over  $\beta_0, \dots, \beta_p$ . Taking partial derivatives and setting to 0:

$$\frac{\partial S}{\partial \beta_0} = -2 \sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip})) = 0$$

$$\frac{\partial S}{\partial \beta_1} = -2 \sum x_{i1} (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip})) = 0$$

⋮

$$\frac{\partial S}{\partial \beta_p} = -2 \sum x_{ip} (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip})) = 0$$

This yields the normal equations

$$n(\hat{\beta}_0) + \hat{\beta}_1 \sum x_{i1} + \dots + \hat{\beta}_p \sum x_{ip} = \sum y_i$$

$$\hat{\beta}_0 \sum x_{i1} + \hat{\beta}_1 \sum x_{i1}^2 + \dots + \hat{\beta}_p \sum x_{i1} x_{ip} = \sum x_{i1} y_i$$

⋮

$$\hat{\beta}_0 \sum x_{ip} + \hat{\beta}_1 \sum x_{i1} x_{ip} + \dots + \hat{\beta}_p \sum x_{ip}^2 = \sum x_{ip} y_i$$

We can write this as

$$(X^T X) \hat{\beta} = X^T y$$

and so

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

- note this needs  $X^T X$  to be invertible;  
ie full rank / all columns are linearly independent.

Note:

① The fitted line is given by

$$\hat{y} = X \hat{\beta}$$

② The vector of fitted values is

$$\hat{y} = X \hat{\beta}$$

③ The residual vector is

$$e = y - \hat{y}$$

\* sum of squares of residuals is  $\sum e_i^2 = e^T e$ .

## THE HAT MATRIX: $\hat{H}$

$\textcircled{1}$  we can express  $\hat{\mu}$  by

$$\hat{\mu} = X\hat{\beta} = X(X^T X)^{-1} X^T y = Hy$$

where

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

is the "hat" matrix which maps the vector of response variables to the vector of fitted values.

Note that

$\textcircled{1}$   $H$  is symmetric (ie  $H^T = H$ ); &

$\textcircled{2}$   $H$  is idempotent (ie  $H^2 = H$ ).

$\textcircled{3}$  We can express our residual vector  $e$  as

$$e = y - \hat{\mu} = y - Hy = (I - H)y$$

## LEAST SQUARES ESTIMATION OF $\sigma^2$

$\textcircled{1}$  The least squares estimate of  $\sigma^2$  for a p explanatory variable multiple regression model with (p+1) parameters is

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{n-(p+1)}$$

where  $df = n-(p+1)$ .

## RESIDUAL STANDARD ERROR

$\textcircled{2}$  The residual standard error is thus

$$\hat{\sigma} = \sqrt{\frac{\sum e_i^2}{n-(p+1)}}$$

## MLE FOR $\beta$

$\textcircled{1}$  The MLE for  $\beta$  is equivalent to the least squares estimate; ie the likelihood function

$$\begin{aligned} L(\beta_0, \dots, \beta_n | y_1, \dots, y_n) &= \prod_{i=1}^n f(y_i) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mu_i)^2}{2\sigma^2}}, \quad \mu_i = \beta_0 + \sum_j \beta_j x_{ij} \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{\sum (y_i - \mu_i)^2}{2\sigma^2}\right) \end{aligned}$$

or equivalently the log likelihood function

$$l(\beta_0, \dots, \beta_n | y_1, \dots, y_n) = c - \frac{\sum (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}))^2}{2\sigma^2}$$

is maximized at  $\beta = (\beta_0, \dots, \beta_p)$  that minimizes

$$\sum \varepsilon_i^2 = \sum (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}))^2$$

## GAUSS-MARKOV THEOREM & BLUE

$\textcircled{1}$  The least squares estimator

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

is the "best linear unbiased estimator" (BLUE) of  $\beta$ .

$\textcircled{2}$  More formally, if we consider the model given by

$$y = X\beta + \varepsilon, \quad E(\varepsilon) = 0, \quad \text{Var}(\varepsilon) = \sigma^2 I$$

then amongst all unbiased linear estimators

$\hat{\beta}^* = M^* y$ , the least squares estimator

$\hat{\beta} = My$  has the "smallest" variance;

i.e.

$$\text{Var}(\hat{\beta}^*) = \text{Var}(\hat{\beta}) + \sigma^2 (M^* - M)(M^* - M)^T$$

where  $(M^* - M)(M^* - M)^T$  is positive semidefinite.

- A is "positive definite" if  $a^T A a > 0$  for any vector  $a$ .

# DISTRIBUTION OF $\hat{\beta}$

We show that

$$\hat{\beta} \sim MVN(\beta, \sigma^2(X^T X)^{-1})$$

where MVN is the multivariate normal distribution.

Proof: First, we have

$$Y = X\beta + \epsilon, \quad \epsilon \sim MVN(0, \sigma^2 I).$$

Thus, by properties of MVN,

$$Y \sim MVN(X\beta, \sigma^2 I).$$

Hence  $\hat{\beta} = (X^T X)^{-1} X^T Y$  also follows a MVN distribution.

Next, see that

$$\begin{aligned} E(\hat{\beta}) &= E((X^T X)^{-1} X^T Y) \\ &= (X^T X)^{-1} X^T E[Y] \\ &= (X^T X)^{-1} X^T (X\beta) \\ &= \beta. \end{aligned}$$

Then

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \text{Var}((X^T X)^{-1} X^T Y) \\ &= (X^T X)^{-1} X^T \text{Var}(Y) [(X^T X)^{-1} X^T]^T \\ &\quad (\text{Var}(AY) = A \text{Var}(Y) A^T) \\ &= \sigma^2 (X^T X)^{-1} X^T [(X^T X)^{-1} X^T]^T \\ &= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1}, \end{aligned}$$

which gives us the desired result.

Q2: The marginal distribution of  $\hat{\beta}_j$  is thus

$$\hat{\beta}_j \sim N(\beta_j, \sigma^2 (X^T X)_{jj}^{-1}) \quad \forall j=0, \dots, p$$

Q3: We also have that

$$\frac{\hat{\beta}_j - \beta_j}{SE(\hat{\beta}_j)} \sim t_{n-(p+1)}, \quad SE(\hat{\beta}_j) = \hat{\sigma} \sqrt{(X^T X)_{jj}^{-1}}$$

## INTERPRETATION OF $\hat{\beta}_j$

Q:  $\hat{\beta}_j$  is the estimated mean change in the response associated with a change of one unit of  $x_j$  whilst holding all other variables constant.

## CIS FOR $\beta_j$

Q1: A  $(1-\alpha) 100\%$  CI for  $\beta_j$  is

$$\hat{\beta}_j \pm t_{n-(p+1), 1-\alpha/2} SE(\hat{\beta}_j)$$

Q2: If  $\beta_j = 0$  is not in the CI, then there is a significant linear relationship between  $y$  &  $x_j$ .

## HYPOTHESIS TESTS FOR $\beta_j$

Q: Hypothesis test for  $\beta_j$ :

①  $H_0: \beta_j = 0; H_A: \beta_j \neq 0$

②  $t = \frac{\hat{\beta}_j - \beta_j}{SE(\hat{\beta}_j)} = \frac{\hat{\beta}_j}{SE(\hat{\beta}_j)}$

③ p-value =  $2P(T > t)$ ,  $T \sim t_{n-(p+1)}$

④ Reject  $H_0$  if  $p < 0.05$ .

## MULTI-COLLINEARITY

We say 2 or more explanatory variables exhibit "multicollinearity" if there exist strong linear relationships between them.

This

- ① increases the variances (and thus std. errors) of the associated parameter estimators;
- ② leads to wide/imprecise CIs & inaccurate conclusions from hypothesis tests.

## VARIANCE INFLATION FACTOR / VIF

The "variance inflation factor" is a measure of multicollinearity associated with some explanatory variable  $x_j$ .

How to calculate VIF for  $x_j$ :

- ① Regress  $x_j$  onto all other  $x_i$ ; ie fit models for  $x_j$  against each other  $x_i$ ;

② Then

$$VIF_j = \frac{1}{1 - R_j^2}.$$

where  $R_j^2$  is the coefficient of determination of the model fit with  $x_j$  as the response.

- in R,  $R_j^2$  is the "multiple R squared" parameter.

Generally, we remove  $x_j$  from the model if  $VIF_j > 10 \Leftrightarrow R_j^2 > .90$ .

In R, we can do

$$> lm(x ~ \underbrace{x_1 + x_2 + \dots + x_n}_{\substack{\text{variable we} \\ \text{are testing}}} + \underbrace{\dots}_{\text{other exp. variables}})$$

and check the multiple R-squared value.

## CI FOR $\mu_{\text{new}}$

Idea: We may want to use our fitted model to estimate the mean response of a new unit in the population.

In particular,

$$\hat{\mu}_{\text{new}} = \mathbf{x}_{\text{new}}^T \hat{\beta}.$$

Then, we show that

$$\hat{\mu}_{\text{new}} \sim N(\mathbf{x}_{\text{new}}^T \hat{\beta}, \sigma^2 \mathbf{x}_{\text{new}}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_{\text{new}})$$

Proof. Recall

$$\hat{\beta} \sim MVN(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}).$$

Then  $\hat{\mu}_{\text{new}} = \mathbf{x}_{\text{new}}^T \hat{\beta}$  must also follow a normal distribution.  
See that

$$\begin{aligned} E(\hat{\mu}_{\text{new}}) &= E(\mathbf{x}_{\text{new}}^T \hat{\beta}) \\ &= \mathbf{x}_{\text{new}}^T E(\hat{\beta}) \\ &= \mathbf{x}_{\text{new}}^T \beta \end{aligned}$$

Then

$$\begin{aligned} \text{Var}(\hat{\mu}_{\text{new}}) &= \text{Var}(\mathbf{x}_{\text{new}}^T \hat{\beta}) \\ &= \mathbf{x}_{\text{new}}^T \text{Var}(\hat{\beta}) \mathbf{x}_{\text{new}} \\ &= \sigma^2 \mathbf{x}_{\text{new}}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_{\text{new}}. \end{aligned}$$

This gives the result, so we're done.  $\square$

Thus, a  $(1-\alpha)\%$  confidence interval for  $\mu_{\text{new}}$  is

$$\hat{\mu}_{\text{new}} \pm t_{n-(p+1), 1-\frac{\alpha}{2}} \hat{\sigma} \sqrt{\mathbf{x}_{\text{new}}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_{\text{new}}}$$

## PREDICTION INTERVAL FOR $y_{\text{new}}$

Q<sub>1</sub> Idea: We may also wish to use our fitted model to predict the value of the response of a new unit of the population.

Q<sub>2</sub> Then, note the variance of  $\hat{y}_{\text{new}}$  is composed of 2 sources of variation:  
① the variation associated with the parameter estimators; &  
② the variance  $\sigma^2$  associated with a random response.

Q<sub>3</sub> Thus our total variation is

$$\sigma^2 + \sigma^2 x_{\text{new}}^T (X^T X)^{-1} x_{\text{new}}.$$

Q<sub>4</sub> Hence, a  $(1-\alpha) 100\%$  prediction interval for  $y_{\text{new}}$  is

$$\hat{y}_{\text{new}} \pm t_{n-(p+1), 1-\frac{\alpha}{2}} \hat{\sigma} \sqrt{1 + x_{\text{new}}^T (X^T X)^{-1} x_{\text{new}}}$$

## CONFIDENCE & PREDICTION BANDS FOR THE SLR MODEL

Q<sub>1</sub> For the SLR model,

$$x_{\text{new}}^T (X^T X)^{-1} x_{\text{new}} = \frac{1}{n} + \frac{(x_{\text{new}} - \bar{x})^2}{s_{xx}}$$

Q<sub>2</sub> Thus, the closer  $x_{\text{new}}$  is to  $\bar{x}$ , the narrower the CIs & PIs.

Q<sub>3</sub> In general, the closer  $\{x_1, \dots, x_p\}$  is to  $\{\bar{x}_1, \dots, \bar{x}_p\}$  in the multiple regression model, the narrower the interval.