

# STA258

## University of Toronto Mississauga

### Inference for Simple Linear Regression (Intro)

AI Nosedal and Omid Jazi

Winter 2023

## Inference on Regression

estimate with residuals

$$Y = \beta_0 + \beta_1 X + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

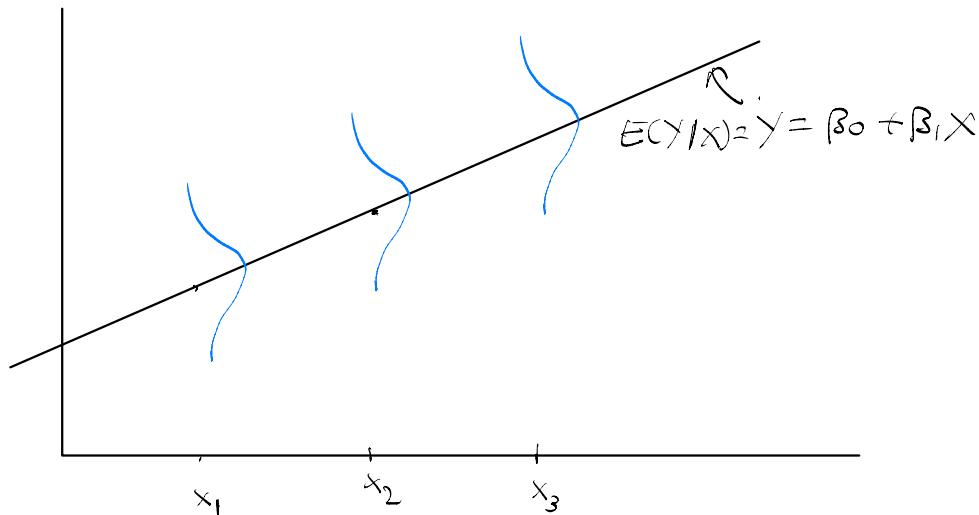
↑  
parameters

Source

can perform inference on  $\beta_0$  and  $\beta_1$ , however we are usually more interested in  $\beta_1$ .

what does the error term  $\varepsilon \sim N(0, \sigma^2)$  mean?

↳ At each value of  $X$ , the errors are distributed normally with a mean of zero and a constant variance.



Can verify with residual plots (assumptions)

We estimate  $\sigma^2$  with a value we call  $s^2$  and use  $s^2$  for inference.

$$Y = \beta_0 + \beta_1 X + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

Estimate  $\sigma^2$  with  $s^2$

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

$\uparrow$  estimate of variance

(residual variance / mean square error (MSE))

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

notice similarity

$$\text{Sample variance } s_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

$$s = \pm \sqrt{s^2}$$

$\uparrow$  estimate of standard deviation

(residual standard deviation /

root mean square error (RMSE))

In calculating  $s^2$ , why do we divide by  $n-2$ ?

Since we estimate 2 unknown parameters in the model (both  $\hat{\beta}_0$  and  $\hat{\beta}_1$ ) which are used in the calculation of  $s^2$

## Example

Suppose an appliance store conducts a 5-month experiment to determine the effect of advertising on sales revenue. The results are shown in a table below. The relationship between sales revenue,  $y$ , and advertising expenditure,  $x$ , is hypothesized to follow a first-order linear model, that is,

$$y = \beta_0 + \beta_1 x + \epsilon$$

where

$y$  = dependent variable

$x$  = independent variable

$\beta_0$  =  $y$ -intercept

$\beta_1$  = slope of the line

$\epsilon$  = error variable

## Example

Month	Advertising Expenditure $x$ (\$ hundreds)	Sales Revenue $y$ (\$ thousands)
1	1	1
2	2	1
3	3	2
4	4	2
5	5	4

The question is this: How can we best use the information in the sample of five observations in our Table to estimate the unknown  $y$ -intercept  $\beta_0$  and slope  $\beta_1$ ?

# Equation of the Least-Squares Regression Line

We have data on an explanatory variable  $x$  and a response variable  $y$  for  $n$  individuals. From the data, calculate the means  $\bar{x}$  and  $\bar{y}$  and the standard deviations  $S_x$  and  $S_y$  of the two variables, and their correlation  $r$ . The least-squares regression line is the line

$$\hat{y} = b_0 + b_1 x$$

with *slope*

$$b_1 = r \frac{S_y}{S_x}$$

and *intercept*

$$b_0 = \bar{y} - b_1 \bar{x}$$

# Least-Squares Regression Line

The **least-squares regression line** of  $y$  on  $x$  is the line that makes the sum of the squares of the vertical distances of the data points from the line as small as possible.

## Our example (again...)

$$\bar{x} = 3, \bar{y} = 2, S_x = 1.5811, S_y = 1.2247, S_{xy} = 1.75.$$

Then, the slope of the least squares line is

$$b_1 = r \frac{S_y}{S_x} = (0.9037) \left( \frac{1.2247}{1.5811} \right) = 0.7$$

and

$$b_0 = \bar{y} - b_1 \bar{x} = 2 - (0.7)(3) = -0.1.$$

The least squares line is thus

$$\hat{y} = -0.1 + 0.7x$$

# **INFERENCE FOR REGRESSION.**

# The Regression Model

We have  $n$  observations on an explanatory variable  $x$  and a response variable  $y$ . Our goal is to study or predict the behavior of  $y$  for given values of  $x$ .

- For any fixed value of  $x$ , the response  $y$  varies according to a Normal distribution. Repeated measures  $y$  are independent of each other.
- The mean response  $\mu_y$  has a straight-line relationship with  $x$ :  
$$\mu_y = \beta_0 + \beta_1 x$$
 The slope  $\beta_1$  and intercept  $\beta_0$  are **unknown** parameters.
- The standard deviation of  $y$  (call it  $\sigma$ ) is the same for all values of  $x$ . The value of  $\sigma$  is **unknown**.

The regression model has three parameters,  $\beta_0$ ,  $\beta_1$ , and  $\sigma$ .

Thus, if

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

is the predicted value of the  $i$ th  $y$  value, then the deviation of the observed value  $y_i$  from  $\hat{y}_i$  is the difference  $y_i - \hat{y}_i$  and the sum of squares of deviations to be minimized is

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2.$$

The quantity  $SSE$  is also called the **sum of squares for error**.

# Fitted Values and Residuals

Fitted Value:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

Residual:

$$\hat{\epsilon} = y_i - \hat{y}_i$$

# Regression Standard Error

The **regression standard error** is

$$s = \sqrt{\frac{1}{n-2} \sum \text{residual}^2} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y})^2} = \sqrt{\frac{SSE}{n-2}}.$$

Use  $s$  to estimate the **unknown**  $\sigma$  in the regression model.

$$\text{Var } x \quad s_x^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

## Standard Error of $\hat{\beta}_1$

The standard error of  $\hat{\beta}_1$  is the standard deviation of the sampling distribution of  $\hat{\beta}_1$  (estimate of slope  $\beta_1$ )

$$SE(\hat{\beta}_1) = \frac{s}{\sqrt{\sum_{i=1}^{n-1} (x_i - \bar{x})^2}} = \frac{s}{\sqrt{(n-1) s_x^2}}$$

Variance of  $x$

Confidence Interval for the Slope

[ Aside  $\text{statistic} \pm (\text{value from } z \text{ or } f) (\text{standard error})$  ]

$$\hat{\beta}_1 \pm t_{(n-2, \alpha/2)} \cdot SE(\hat{\beta}_1)$$

$$\hat{\beta}_1 \pm t_{(n-2, \alpha/2)} \frac{s}{\sqrt{\sum_{i=1}^{n-1} (x_i - \bar{x})^2}}$$

# Estimating $\sigma$

Advertising Expenditure $x$ (\$ hundreds)	Sales Revenue				$(y - \hat{y})^2$
	$y$ (\$ thousands)	$\hat{y}$	$y - \hat{y}$		
1	1	0.6	0.4		0.16
2	1	1.3	-0.3		0.09
3	2	2	0		0
4	2	2.7	-0.7		0.49
5	4	3.4	0.6		0.36
					SSE = 1.10

Revisit Example on advertising and sales from set 15  
 and construct a 95% confidence interval on the  
 slope. provide an interpretation of CI

From earlier :  $\hat{y} = -0.1 + 0.7x$

$x$	$y$	$\hat{y}$	residuals	$(\text{residual})^2$	$x - \bar{x}$	$(x - \bar{x})^2$
1	1	0.6	0.4	0.16	-2	4
2	1	1.3	-0.3	0.09	-1	1
3	2	2	0	0	0	0
4	2	2.7	-0.7	0.49	1	1
5	4	3.4	0.6	0.36	2	4

$$\sum x_i = 15$$

$$SSE = 1.10$$

$$\sum (x_i - \bar{x})^2 = 10$$

✓ ≡

$$\textcircled{\text{X}} \quad \frac{\sum x_i}{n} = \frac{15}{5} = 3$$

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

$$= \frac{1.10}{5-3} = 0.3667$$

$$S = \sqrt{S^2} = \sqrt{0.3667} = 0.6055 \quad \checkmark$$

$$\hat{\beta} \stackrel{\checkmark}{=} t_{(n-2, \alpha/2)} \frac{s}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad \checkmark$$

Find  $t_{(n-2, \alpha/2)}$

$$n-2 = 5-2 = 3 \quad \text{For a } 95\% \text{ CI}$$

$$1 - \alpha = 0.95$$

$$\alpha/2 = 0.025$$

$$t_{(n-2, \alpha/2)} = t_{(3, 0.025)} = 3.182$$

$$\hat{\beta}_1 \pm t_{(n-2, \alpha/2)} \frac{s}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

$$= 0.7 \pm 3.182 \left( \frac{0.6055}{\sqrt{10}} \right)$$

$$= 0.7 \pm 0.6092$$

$$= (0.0908, 1.3092)$$

(+) (+)

Interpretation:

We are 95% confident the slope ( $\beta_1$ ) for this model lies between 0.0908 and 1.3092.



# Examine possibilities of CI's

Know sign, just don't know steepness

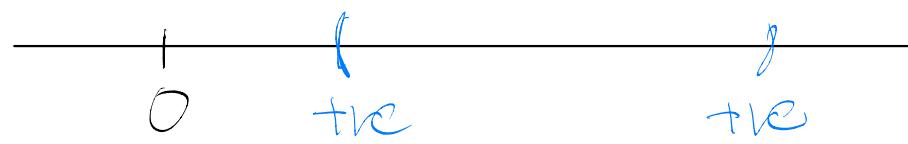
Suppose  $CI : l \text{ -ve}, -ve$



Suggests  $\beta_1$  has a -ve sign

(suggests -ve correlation, potentially good model)

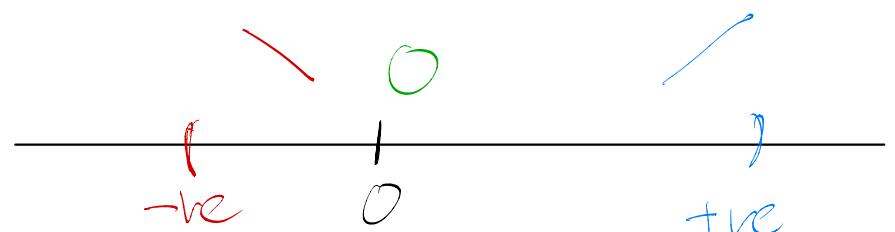
Suppose  $CI : (t\text{ve}, t\text{ve})$



Suggests  $\beta_1$  has a +ve sign

(suggests +ve correlation, potentially good model)

Suppose  $CI : l \text{ -ve}, t\text{ve}$



$\beta_1 = 0$  is plausible

suggests no LINEAR relationship between X and y

$$y = \beta_0 + \beta_1 X$$

A scatter plot with a single data point at the origin (0,0). A horizontal line passes through this point, representing a model where  $\beta_1 = 0$ .

# Regression Standard Error

The **regression standard error** is

$$s = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y})^2} = \sqrt{\frac{SSE}{n-2}} = \sqrt{\frac{1.1}{3}} = 0.6055.$$

Use  $s$  to estimate the **unknown**  $\sigma$  in the regression model.

# Confidence interval for the regression slope

A level  $C$  confidence interval for the slope  $\beta_1$  of the true regression line is

$$b_1 \pm t^* SE_{b_1}.$$

In this formula, the standard error of the least-squares slope  $b$  is

$$SE_{b_1} = \frac{s}{\sqrt{\sum(x_i - \bar{x})^2}} = \frac{s}{\sqrt{(n-1)S_X^2}}$$

and  $t^*$  is the critical value for the  $t(n-2)$  density curve with area  $C$  between  $-t^*$  and  $t^*$ .

# Estimating $\sigma$

Advertising Expenditure $x$ (\$ hundreds)	Sales Revenue				$(y - \hat{y})^2$
	$y$ (\$ thousands)	$\hat{y}$	$y - \hat{y}$		
1	1	0.6	0.4		0.16
2	1	1.3	-0.3		0.09
3	2	2	0		0
4	2	2.7	-0.7		0.49
5	4	3.4	0.6		0.36
					SSE = 1.10

# Hypothesis Test on the Slope $\beta_1$

$$H_0: \beta_1 = 0 \quad H_a: \beta_1 > 0$$

$$H_0: \beta_1 = 0 \quad H_a: \beta_1 < 0$$

$$H_0: \beta_1 = 0 \quad H_a: \beta_1 \neq 0 \quad \leftarrow \text{most common}$$

Test Stat:

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$$

[Aside

$$\text{test stat} = \frac{\text{stat} - \text{hyp value}}{SE(\text{stat})}$$

$$= \frac{\hat{\beta}_1 - 0}{\frac{s}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}} \\ \text{ssx}$$

Reference distribution:  $t$  distribution at  $n-2$  df.

For advertising example, perform a  
2 sided hyp test on slope

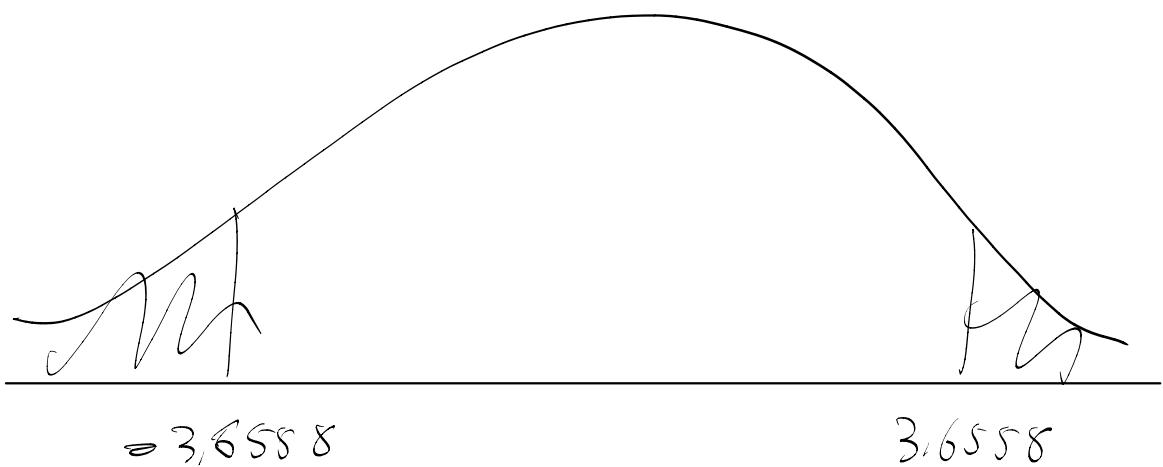
$$H_0: \beta_1 = 0 \quad H_a: \beta_1 \neq 0$$

Test Stat

$$t = -0.1 + 6.7$$

$$t^* = \frac{\hat{\beta}_1 - 0}{\sqrt{s/\sum(x_i - \bar{x})^2}} = \frac{0.7 - 0}{0.6055/\sqrt{16}} = 3.6558$$

reject  $H_0$ :  $t$  at  $n-2 = 5-2 = 3$  df



$$0.01 < p\text{-val} < 0.025$$

$$p\text{-val} < 0.05$$

reject  $H_0$ , conclude  $H_A: \beta_1 \neq 0$

Slope should be in model

## Our example

For the advertising-sales example, a 95% Confidence Interval for the slope  $\beta_1$  is

$$0.7 \pm 3.182 \left( \frac{0.6055}{\sqrt{10}} \right)$$

$$0.7 \pm 0.6092$$

Thus, we estimate with 95% confidence that the interval from 0.0908 and 1.3092 includes the parameter  $\beta_1$

## Testing the hypothesis of no linear relationship

We can also test hypotheses about the slope  $\beta_1$ . The most common hypothesis is

$$H_0 : \beta_1 = 0.$$

A regression line with slope 0 is horizontal. That is, the mean of  $y$  does not change at all when  $x$  changes. So this  $H_0$  says that there is no true linear relationship between  $x$  and  $y$ .

# Significance test for regression slope

To test the hypothesis  $H_0 : \beta_1 = 0$ , compute the  $t$  statistic

$$t = \frac{b_1}{SE_{b_1}}.$$

In terms of a random variable  $T$  having the  $t(n - 2)$  distribution, the P-value for a test of  $H_0$  against:

$H_a : \beta_1 \neq 0$  is  $2P(T > |t|)$ .

$H_a : \beta_1 > 0$  is  $P(T > t)$ .

$H_a : \beta_1 < 0$  is  $P(T < t)$ .

# Our example

$$\alpha = 0.05$$

- 1)  $H_0 : \beta_1 = 0$  vs  $H_a : \beta_1 > 0$
- 2)  $t_* = \frac{b_1}{SE_{b_1}} = \frac{0.7}{0.1914} = 3.6572$

$$3) P\text{-value} = P(T > t) = P(T > 3.6572) \text{ d.f.} = n-2 = 5 - 2 = 3.$$

Using t-distribution table,  $0.01 < P\text{ value} < 0.025$

- 4) Since  $P\text{-value} < \alpha = 0.05$ , we reject  $H_0$ .

## Our example (different $H_a$ )

$$\alpha = 0.05$$

- 1)  $H_0 : \beta_1 = 0$  vs  $H_a : \beta_1 \neq 0$
- 2)  $t_* = \frac{b_1}{SE_{b_1}} = \frac{0.7}{0.1914} = 3.6572$

$$3) P\text{-value} = 2P(T > |t|) = 2P(T > 3.6572) \text{ d.f.} = n-2 = 5 - 2 = 3.$$

Using Table 3,  $0.02 < P\text{ value} < 0.05$

- 4) Since  $P\text{-value} < \alpha = 0.05$ , we reject  $H_0$ .

# Our example (different $H_a$ )

```
x=c(1,2,3,4,5);  
y=c(1,1,2,2,4);  
mod=lm(y~x);  
summary(mod);
```

## Review

$$Y = \beta_0 + \beta_1 X + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

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

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

r : coeff of correlation (strength)

$r^2$ : |||| determination (% variability)

$$S^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2} = \overline{SS_E}$$

$$S = \sqrt{S^2}$$

$$SE(\hat{\beta}_1) = \frac{S}{\sqrt{s_{xx}}}$$

$$(I : \hat{\beta}_1 \pm t_{(n-2, \alpha/2)} \cdot SE(\hat{\beta}_1))$$

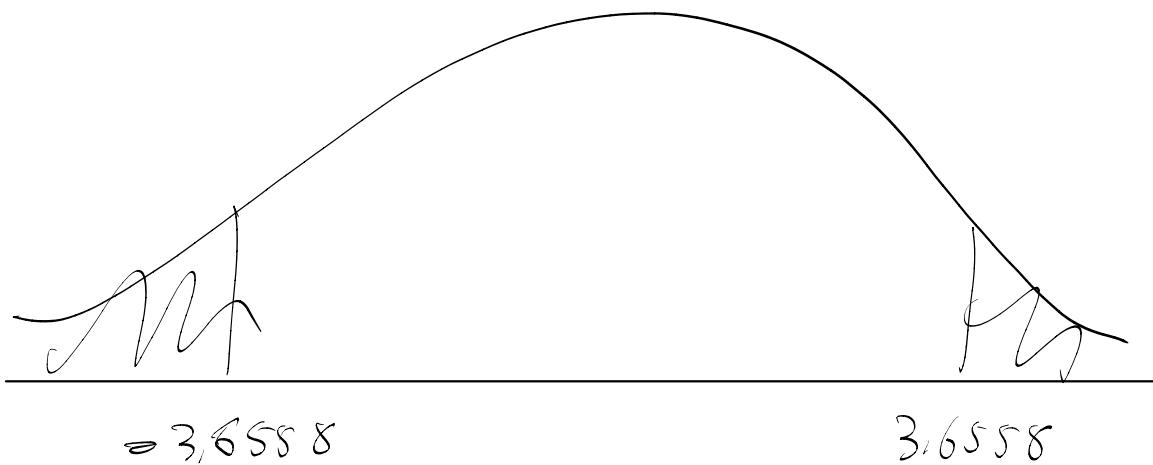
$$H_0: \beta_1 = 0$$

$$\text{hyp test : test stat } t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$$

$$t^* = \frac{\hat{\beta}_1 - 0}{\frac{s/\sqrt{\sum(x_i - \bar{x})^2}}{0.6055/\sqrt{16}}} = \frac{0.7 - 0}{0.6055/\sqrt{16}} = 3.6558$$

$SE(\hat{\beta}_1)$

referenze des t: t ab  $n-2 = 5-2 = 3$  df



$$0.01 < p-value < 0.028$$

$$0.02 < p-value < 0.05$$

reject  $H_0$ , conclude  $H_A: \beta_1 \neq 0$

Slope should be in model

# Our example (different $H_a$ )

```
##  
## Call:  
## lm(formula = y ~ x)  
##  
## Residuals:  
##           1            2            3            4            5  
## 4.000e-01 -3.000e-01 -3.886e-16 -7.000e-01 6.000e-01  
##  
## Coefficients:  
##               Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -0.1000   0.6351  -0.157  0.8849  
## x            0.7000   0.1915   3.656  0.0354 *## ---  
## Signif. codes:  
## 0 *** 0.001 ** 0.01 * 0.05 . 0.1 ' ' 1  
##  
## Residual standard error: 0.6055 on 3 degrees of freedom
```

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

$\hat{\beta}_0$

$\hat{\beta}_1$

$SE(\hat{\beta}_1)$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x = -0.10 + 0.70x$$

$$SE(\hat{\beta}_1) = 0.1918$$

By default, R conducts the following test for each coefficients

$$H_0: \beta_j = 0 \quad H_a: \beta_j \neq 0 \quad (\text{two sided})$$

test stat

$$t^* = \frac{\hat{\beta}_j - 0}{SE(\hat{\beta}_j)}$$

---

For advertising and sales data

$$H_0: \beta_1 = 0 \quad H_a: \beta_1 \neq 0$$

$$t^* = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} = \frac{0.70}{0.1918} = 3.656$$

```
> summary(model);
```

Call:

```
lm(formula = camrys$Price ~ Odometer, data = camrys)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.68679	-0.27263	0.00521	0.23210	0.70071

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )							
(Intercept)	17.248727	0.182093	94.72	<2e-16	***						
Odometer	-0.066861	0.004975	-13.44	<2e-16	***						
---											
Signif. codes:	0	'***'	0.001	'**'	0.01	'*'	0.05	'.'	0.1	' '	1

Residual standard error: 0.3265 on 98 degrees of freedom

Multiple R-squared: 0.6483, Adjusted R-squared: 0.6447

F-statistic: 180.6 on 1 and 98 DF, p-value: < 2.2e-16

## The three deviations associated with a data point

We square all three deviations for each one of our data points, and sum over all n points. Here, cross terms drop out, and we are left with the following equation:

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

$$SST = SSE + SSR$$

Total sum of squares = Sum of squares for error + Sum of squares for regression.

From Exercise 11.15 (HW?) we have the following identity:

$$\begin{aligned} SSE &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n (y_i - \bar{y})^2 - \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= S_{YY} - \hat{\beta}_1 S_{XY} \end{aligned}$$

Notice that this provides an easier computational method of finding SSE.

# ANOVA TABLE

Analysis of Variance

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

$$MS = \frac{SS}{df}$$

Source of Variation	Sum of Squares	Degrees of freedom	Mean square	Computed F
Regression	SSR	1	$MSR = SSR/1$	$SSR/s^2$
Error	SSE	n-2	$s^2 = SSE/n-2$	X
Total	SST	n-1	MSE	X

$$SSR + SSE = SS_T$$

$$(n-2) = n-1$$

For the general multivariate regression model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

$\underbrace{\beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}_p \text{ predictors}$

ANOVA can be used for testing

$$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0 \quad \begin{array}{l} \text{all coeffs are zero} \\ (\text{none of the variables are good predictors}) \end{array}$$

$$H_a: \text{At least one } \beta_j \neq 0$$

$$\exists \beta_j \neq 0 \quad j=1, \dots, p$$

Test statistic

$$F = \frac{\text{MSR}}{\text{MSE}} = \frac{\text{SSR}/p}{\text{SSE}/(n-p-1)} \sim F_{(p, n-p-1)}$$

Reference dist: F with num = p, denom = n - p - 1 df

# Our example

```
x=c(1,2,3,4,5);  
y=c(1,1,2,2,4);  
mod=lm(y~x);  
anova(mod);
```

# Our example

```
## Analysis of Variance Table
##
## Response: y
##           Df Sum Sq Mean Sq F value    Pr(>F)
## x          1   4.9   4.9000  13.364 0.03535 *
## Residuals  3   1.1   0.3667
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

$$H_0: \beta_1 = 0$$

# Testing the Coefficient of Correlation

We can use the coefficient of correlation to test for a linear relationship between two variables. When there is no linear relationship between the two variables,  $\rho = 0$ . To determine whether we can infer that  $\rho = 0$ , we test the hypothesis

$$H_0 : \rho = 0.$$

# Significance test for regression slope

To test the hypothesis  $H_0 : \rho = 0$ , compute the  $t$  statistic

$$t = \frac{\sqrt{n-2}}{\sqrt{1-r^2}} r.$$

In terms of a random variable  $T$  having the  $t(n-2)$  distribution, the P-value for a test of  $H_0$  against:

$H_a : \rho \neq 0$  is  $2P(T > |t|)$ .

$H_a : \rho > 0$  is  $P(T > t)$ .

$H_a : \rho < 0$  is  $P(T < t)$ .

# STA258: Statistics with Applied Probability

## Simple Linear Regression

---

Nishan Mudalige

Winter 2025

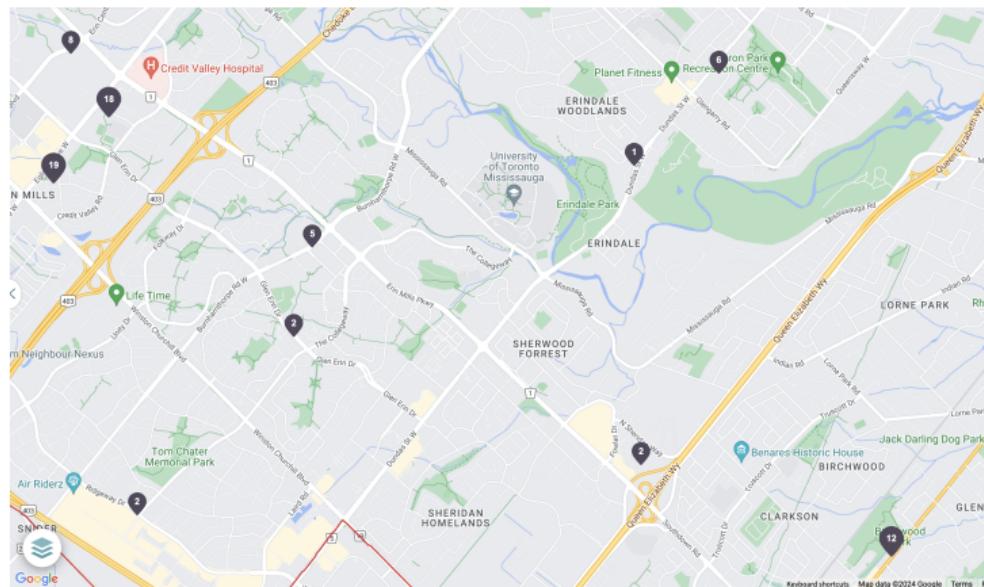


## Outline

# 1 SLR Example



## Example - Apartments Around UTM







## Example Ctd...

Price (\$1000)	Area (sq ft)	Beds	Baths
620	11.0	2	2
590	6.5	2	1
620	10.0	2	2
700	8.4	2	2
680	8.0	2	2
500	5.7	1	1
760	12.0	2	2
800	14.0	3	1
660	7.3	2	1



## Example Ctd...

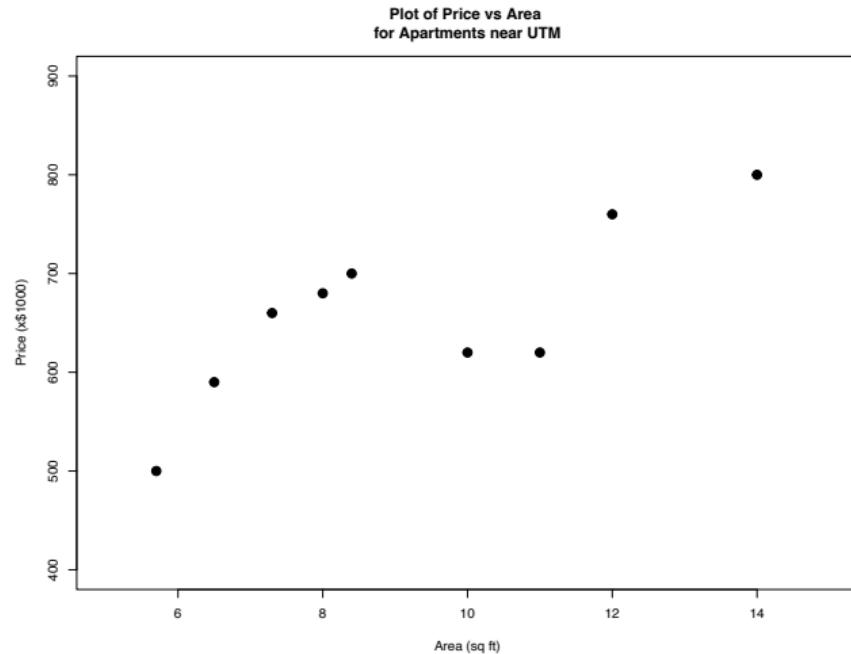
Y

X

Price (\$1000)	Area (sq ft)	Beds	Baths
620	11.0	2	2
590	6.5	2	1
620	10.0	2	2
700	8.4	2	2
680	8.0	2	2
500	5.7	1	1
760	12.0	2	2
800	14.0	3	1
660	7.3	2	1

## Example Ctd...

Price (×\$1000)	Area (×100 sq ft)
620	11.0
590	6.5
620	10.0
700	8.4
680	8.0
500	5.7
760	12.0
800	14.0
660	7.3



Surf

$y$	$x$	$(x - \bar{x})$	$(y - \bar{y})$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$
620	11.0	1.8	-38.9	-70.02	3.24
590	6.5	-2.7	-68.9	186.03	7.29
620	10.0	0.8	-38.9	-31.12	0.64
700	8.4	-0.8	41.1	-32.88	0.64
680	8.0	-1.2	21.1	-25.32	1.44
500	5.7	-3.5	-158.9	556.15	12.25
760	12.0	2.8	101.1	283.08	7.84
800	14.0	4.8	141.1	677.28	23.04
660	7.3	-1.9	1.1	-2.09	3.61

$$\Sigma y = 5,930 \quad \bar{x} = 82.9$$

$$S_{xy} = 1,541.11 \quad S_{xx} = 60.00$$

$$\bar{y} = \frac{\sum y}{n} = \frac{5930}{9} = 658.89$$

$$\bar{x} = \frac{\sum x}{n} = \frac{82.9}{9} = 9.21$$

---

## Example Ctd - Find $\hat{\beta}_1$ and $\hat{\beta}_0$

$$\hat{\beta}_1 = \frac{s_{xy}}{s_{xx}} = \frac{1541.11}{60} = 25.69$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 658.87 - (25.69)(9.21) = 422.28$$

## Example Ctd - Equation of regression line

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

$$\hat{y} = 422.28 + 25.69 x$$

predicted price (\$/sq ft)

estimate of intercept ( $\hat{\beta}_0$ )

estimate of coefficient of area (100 sq ft)

area

Interpretations of  $\hat{\beta}_0$ ,  $\hat{\beta}_1$

$\hat{\beta}_1 = 25.69$

unit of  $x$

1

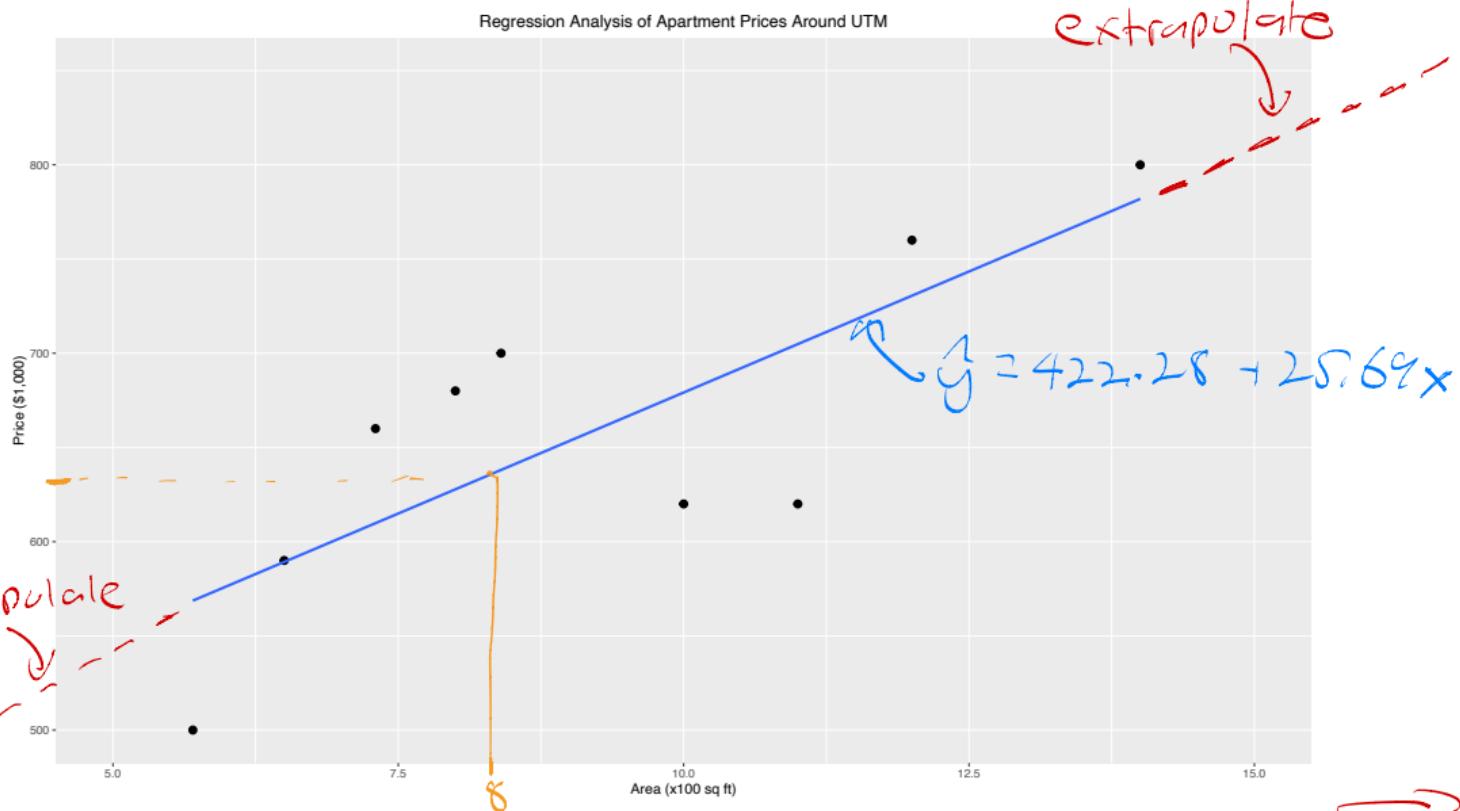
For an increase in area ( $x$ ) by 1 ( $x/100 \text{ sq ft}$ ), we expect the price ( $y$ ) to increase by  $25.69$  (\$1,000) on average  
 $(\hat{\beta}_1 > 0)$

$\hookrightarrow \$ 25,690$

$\hat{\beta}_0 = 422.28$  ( $\times \$100 = \$422,280$ )

The predicted price when area is zero

$\hookrightarrow$  limitation of model



## Example Ctd

$x=8$

Use the model to predict the price of an apartment which is 800 sq ft.

$$x=8, \quad \hat{y} = 422.28 + 25.69(8) \\ = 627.8 \quad (\$1000)$$

interpolation

## Example Ctd

$$x=25$$

Use the model to predict the price of an apartment which is 2,500 sq ft.

$$\begin{aligned}x &= 25 \quad \hat{y} = 422.28 + 25.69(25) \\&= 1064.53 \quad (\$1000)\end{aligned}$$

extrapolation

## R Code

- We can create a simple linear regression model in R using the `lm` command.

### R Code

```
lm(y ~ x, data = data_source)
```

- The data is available in the `apt_around_utm.csv` file.

### R Code

```
apt = read.csv(file.choose())
# apt = read.csv("~/PATH_TO_FILE/apt_around_utm.csv")

apt_model = lm(price ~ area, data = apt)
```

## R Output

### R Output

```
> apt_model

Call:
lm(formula = price ~ area, data = apt)
```

### Coefficients:

(Intercept)

422.26

area

25.69

422.28

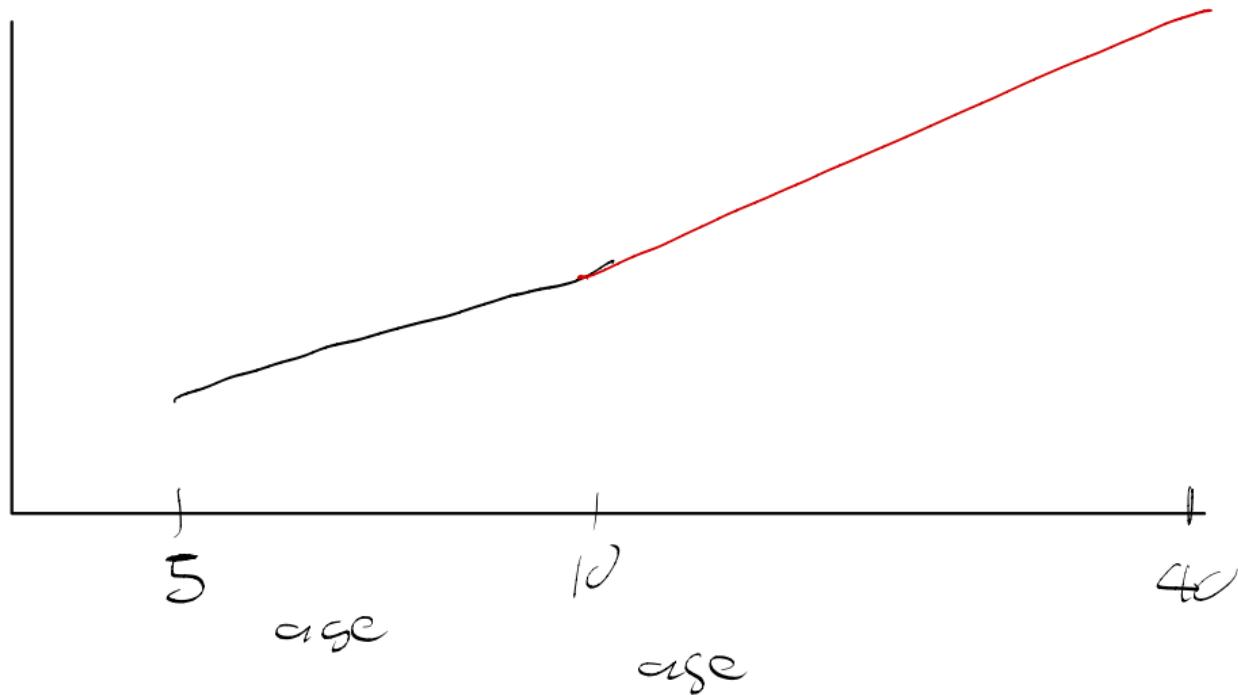
---

## Interpolation and Extrapolation

- **Interpolation** is calculating predicted values of  $y$  using our linear model while working within the range of  $x$  in which data was available to construct our model.
- **Extrapolation** is calculating predicted values of  $y$  using our linear model outside the range of  $x$  used to obtain the linear model.
- Interpolation is usually safe if we have a good linear model.
- Extrapolation must be performed carefully since extrapolations that are done without any foresight can be very inaccurate.

extrapolation can be dangerous

height



## Example Ctd...

Recall our model:  $\hat{y} = 25.69x + 422.26$

$y$	$x$	$\hat{y}$	$y - \hat{y}$	$(y - \hat{y})^2$
620	11.0	704.85	-84.85	7,198.73
590	6.5	589.24	0.76	0.58
620	10	679.16	-59.16	3,499.36
700	8.4	638.05	61.95	3,837.62
680	8.0	627.78	52.22	2,727.4
500	5.7	568.69	-68.69	4,718.13
760	12.0	730.54	29.46	868.17
800	14.0	781.92	18.08	327.06
660	7.3	609.79	50.21	2,520.79

$\text{SSE} = 25,697.83$

## Example Ctd - Hypothesis Test

Conduct a hypothesis test on the slope at the 5% level of significance.

$$H_0 : \beta_1 = 0 \quad \text{vs.} \quad H_a : \beta_1 \neq 0$$

$$S^2 = \frac{\text{SSE}}{n-2} = \frac{25697.83}{7} = 3670.2618$$

$$s = \sqrt{s^2} = 7.28$$

$$SE(\hat{\beta}_1) = \frac{s}{\sqrt{SS_{xx}}} = \frac{7.28}{\sqrt{60}} = 7.82$$

*X earlier*

test stat

$$t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} = \frac{28.69 - 0}{7.82} = 3.284$$

t-dist at 7 df

$$0.005 < p\text{-val} < 0.01$$

Suff evidence to reject  $H_0$ , conclude  
model should contain  $\beta_1$ , suggesting  
there is a relationship between price and  
area

## Example Ctd...

*f model*

Recall our model:  $\hat{y} = 25.69x + 422.26$

*P*

<i>y</i>	<i>x</i>	$\hat{y}$	$(y - \hat{y})^2$	$(\hat{y} - \bar{y})^2$	$(y_i - \bar{y})^2$
620	11.0	704.85	7,198.73	2,112.00	1,512.35
590	6.5	589.24	0.58	4,850.88	4,745.68
620	10	679.16	3,499.36	410.73	1,512.35
700	8.4	638.05	3,837.62	434.20	1690.12
680	8.0	627.78	2,727.4	968.04	445.68
500	5.7	568.69	4,718.13	8,136.08	2,5245.68
760	12.0	730.54	868.17	5,133.21	10,223.46
800	14.0	781.92	327.06	1,5135.47	19,912.35
660	7.3	609.79	2,520.79	2,410.45	1.23
			25,697.83	39,591.06	65,288.90

*SSE* + *SSR* = *SSTotal*

## Example Ctd...

coeff of determination



coeff of correlation

Use the information to Calculate  $r^2$  and  $r$ .

$$r^2 = \frac{SSR}{SS_{\text{Total}}} = \frac{39,591.06}{65,288.90} = 0.6063 = 60.63\%$$

In Interpretation

Approx 60.63% of variability in price ( $y$ )

is explained by the regression model

$$r = \pm \sqrt{r^2}$$
$$= \pm \sqrt{0.6063}$$
$$= \pm 0.779$$

$$\hat{y} = 422.26 + 25.69 \times \hat{\beta}_1 > 0$$

Interpretation

Strong positive correlation between area ( $x$ )  
and price ( $y$ )

## R Code

- R can do a lot of the work for us

### R Code

```
model = lm(y ~ x, data = data_source)  
summary(model)
```

- For our example we can type

### R Code

```
apt_model = lm(price ~ area, data = apt)  
summary(apt_model)
```

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

## Example - R Output

	$\hat{\beta}_0$	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	422.256	74.834	5.643	0.00078	***
area	25.690	7.823	3.284	0.01341	*
---				$SE(\hat{\beta}_1)$	
Signif. codes:	0 ***	0.001 **	0.01 *	0.05 .	0.1 ' 1

Residual standard error: 60.59 on 7 degrees of freedom

Multiple R-squared: 0.6064, Adjusted R-squared: 0.5502

F-statistic: 10.78 on 1 and 7 DF, p-value: 0.01341

By default, R performs 2-sided test

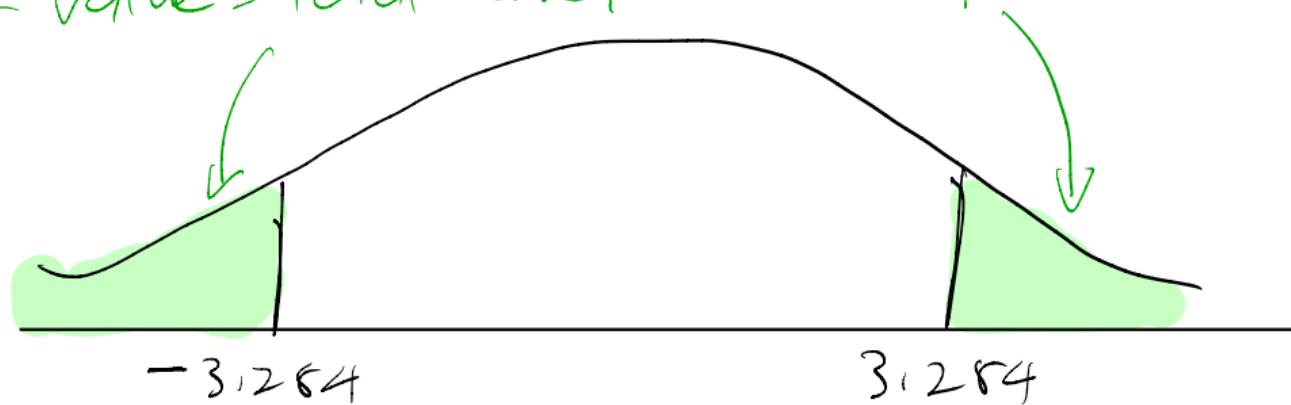
$$H_0: \beta_1 = 0 \quad H_a: \beta_1 \neq 0$$

Test stat

$$t^* = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} = \frac{28.690 - 0}{7.823} = 3.284 \sim t_{(n-2)}$$

df = 9 - 2 = 7

$$p\text{-value} = \text{total area} = 0.0134$$



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

## Example - R Output

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	422.256	74.834	<del>5.618</del>	<del>&lt; 0.001 ***</del>
area	25.690	7.823	<del>3.271</del>	<del>&lt; 0.051 *</del>
---				

Signif. codes: 0 ‘\*\*\*’ 0.001 ‘\*\*’ 0.01 ‘\*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 60.59 on 7 degrees of freedom

Multiple R-squared: 0.6064, Adjusted R-squared: 0.5502

F-statistic: 10.78 on 1 and 7 DF, p-value: 0.01341

## Residual Plots (not in slides)

$$Y = \beta_0 + \beta_1 X + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

Residual plots are used to verify the assumption related to the error terms.  $\varepsilon \sim N(0, \sigma^2)$

$\xrightarrow{\text{mean=0}}$   $\curvearrowright \text{constant variance}$

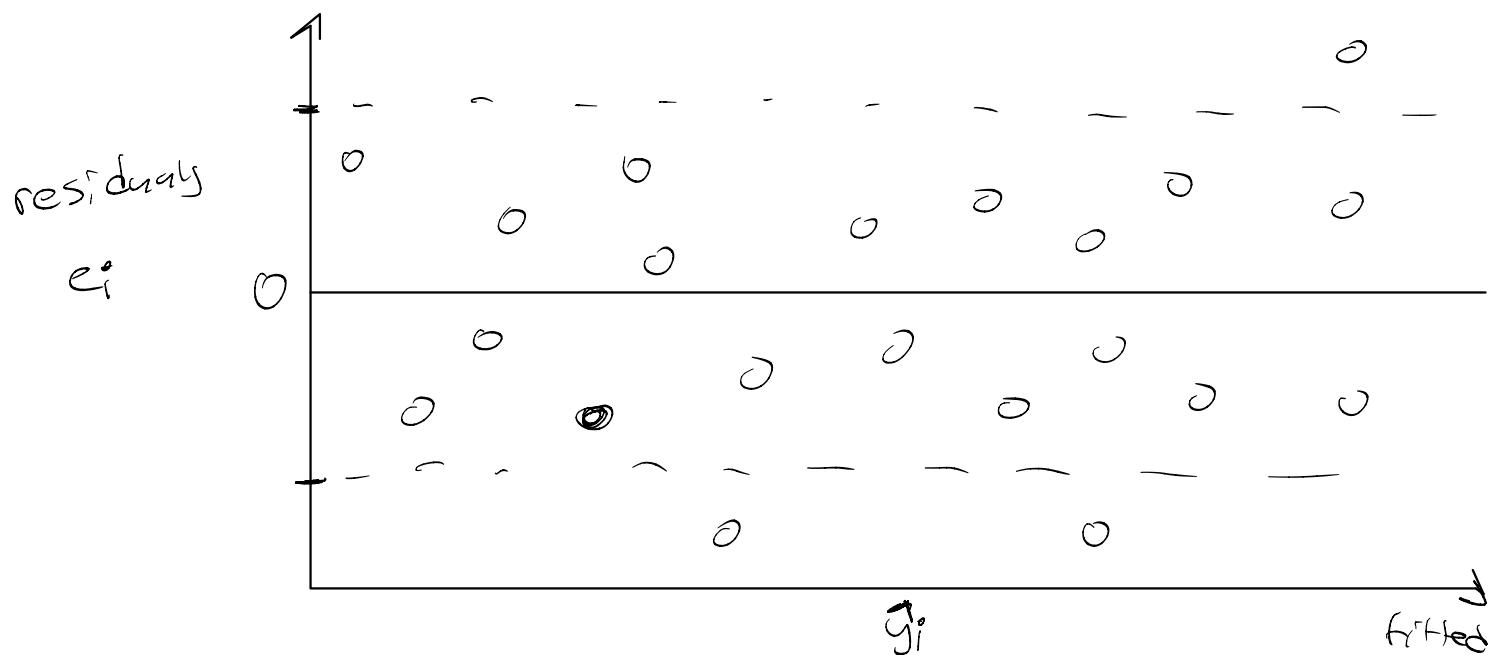
Plot residuals ( $e_i = y_i - \hat{y}_i$ ) against fitted ( $\hat{y}_i$ ) values

$y\text{-axis}$                                      $x\text{-axis}$

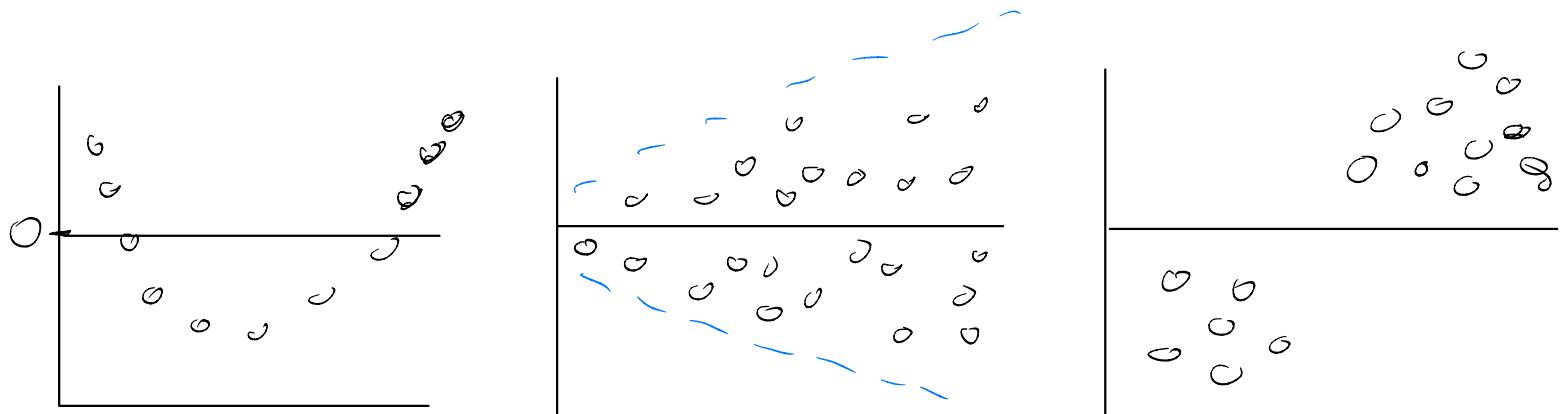
If assumption  $\varepsilon \sim N(0, \sigma^2)$  is satisfied, residual plot should have the following features

- ✓ Random Scattering (No obvious pattern)  
pattern (eg curve) suggests non-linear relationship  
Random Scattering also suggest indep
- ✗ Majority of residuals within a horizontal band  
with approx half residuals above zero and half below)  
Suggests constant variance (homoskedastic)
- ✗ No influential points / clustering

Example of a plot satisfying assumptions



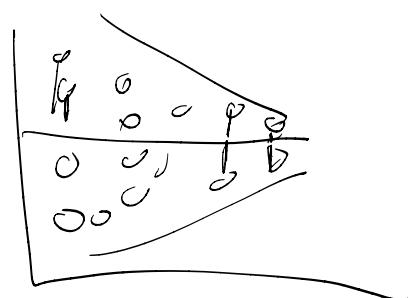
Examples of plots violating assumptions



Suggests  
non-linear  
relationship

non-constant  
variance  
(heteroskedastic)

clustering



## Assumptions of SLR

$$Y = \beta_0 + \beta_1 X + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

1) Relationship between  $X$  and  $Y$  is linear

2) Residuals are STA 302: regression

Independent  
Have constant variance  
Normal

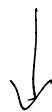
} Verify with resid plots

---

Inference on One Sample (Z or t)



Inference on two samples (Z or t)



Inference on 3 or more groups

(F-tests with ANOVA)