

COMM215
First the Foundation, then Innovation

LESSON 10 & 11 SIMPLE LINEAR REGRESSION

SAMIE L.S. LY

-  1. Simple Linear Regression Model
- 2. Least Square Method
- 3. Coefficient of Determination
- 4. Model Assumptions
- 5. Testing for Significance
- 6. Covariance & Coefficient of Correlation
- 7. Using the Estimated Regression –Equation for Estimation and Prediction

SIMPLE LINEAR REGRESSION MODEL

REGRESSION MODEL

REGRESSION EQUATION

ESTIMATED REGRESSION EQUATION

THE OBJECTIVE OF SLR

Let's say, I would like to create the ultimate equation to figure out my Final COMM 215 grade.

**What influences my Final COMM 215 grade?
Think of a few examples.**

FINAL COMM 215 Grade (y) =

Study Time (x_1) +

Work Time (x_2) +

of hours spent on Facebook (x_3) +

Family Time (x_4) +

anything else you can think of (x_{\dots}) ...

THE OBJECTIVE OF SLR

Let's say, I would like to create the ultimate equation to figure out my Final COMM 215 grade.

Let's say,

I study 5 hours a week,

work part-time for 25 hours a week,

I spend 3 hours on facebook,

And I have 2 kids.

Here is 1 scenario.

THE OBJECTIVE OF SLR

Let's say, I would like to create the ultimate equation to figure out my Final COMM 215 grade.

What if it was someone else?
With a different profile?

Do I have to make my calculations all over again?

If I set up a regression line, I just need to plug in values and get an estimation of my Final COMM 215 grade.

Voila!

Then you might ask... how?

THE OBJECTIVE OF SLR

Let's say, I would like to create the ultimate equation to figure out my Final COMM 215 grade.

How do I create this regression line?

Answer: By gathering data from history.

I am going to take a sample of individuals who took COMM 215 before and write down their profile.

Hours per week	Study Time	Work Time	Family Time	Facebook Time
Bob	10	25	5	0
Sally	12	0	15	15
Eric	3	40	10	0

THE OBJECTIVE OF SLR

Let's say, I would like to create the ultimate equation to figure out my Final COMM 215 grade.

Since we are only considering 1 independent variable, let's take just 1, Study Time as the main indicator of your Final COMM 215 grade.

Hours per week	Study Time(x)	Final Grade(y)
Bob	10	89
Sally	12	67
Eric	3	45

THE OBJECTIVE OF SLR

Let's say, I would like to create the ultimate equation to figure out my Final COMM 215 grade.

We've generated this equation!

Now, if Michelle asks, if I study 8 hours a week, what would be my estimated Final COMM 215 grade?

$$y = 3.4478x + 38.269$$

SIMPLE LINEAR REGRESSION

1 VARIABLE

FINAL COMM 215 GRADE (y) = STUDY TIME (x)

MULTIPLE LINEAR REGRESSION

MORE THAN 1 VARIABLE

FINAL COMM 215 GRADE (y) =

STUDY TIME (x_1) + WORK TIME(x_2) +

As a way of predicting sales

Managerial decisions are often made based on the relationship between two or more variables.

The statistical process is called regression analysis used to develop an equation showing how the variables are related.

**The variable
being predicted is
called the dependent variable.**

**The variable or variables
being used to predict the value of the
dependence variable are
called independent variables.**

FINAL COMM 215 Grade (y) =

Study Time (x_1) +

Work Time (x_2) +

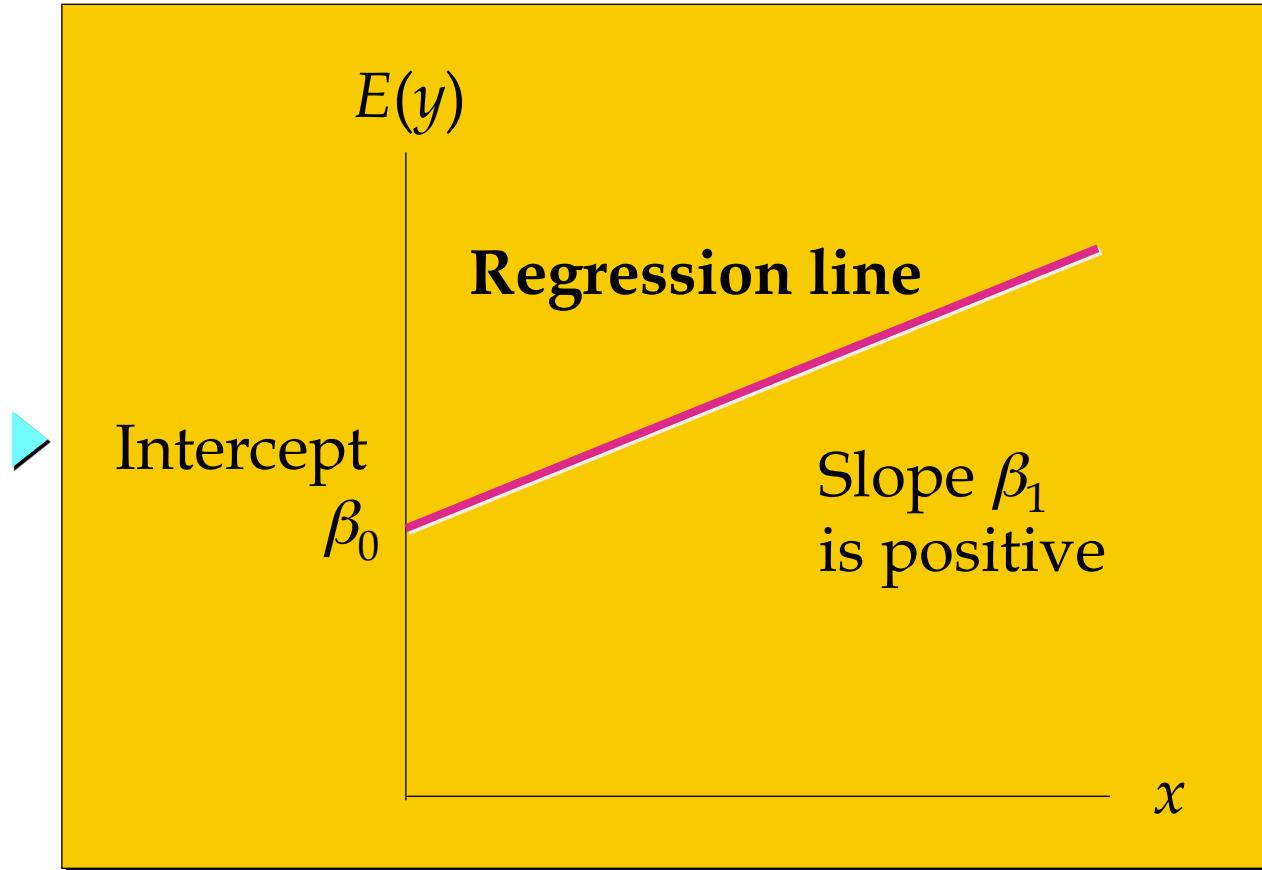
of hours spent on Facebook (x_3) +

Family Time (x_4) +

anything else you can think of (x_{\dots}) ...

SIMPLE LINEAR REGRESSION EQUATION

- Positive Linear Relationship

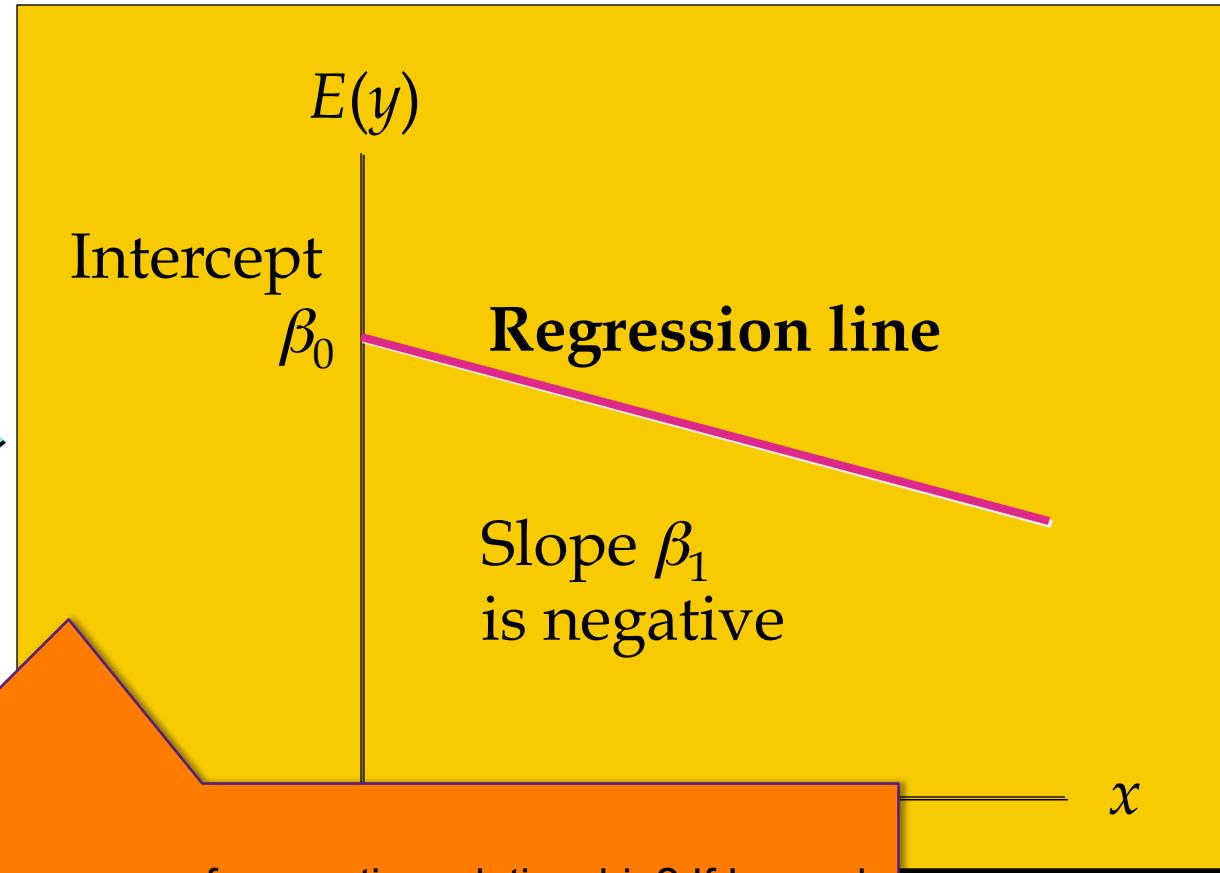


- The relationship between the two variables is approximated by a straight line.

Bowerman, et al. (2017) pp. 534

SIMPLE LINEAR REGRESSION EQUATION

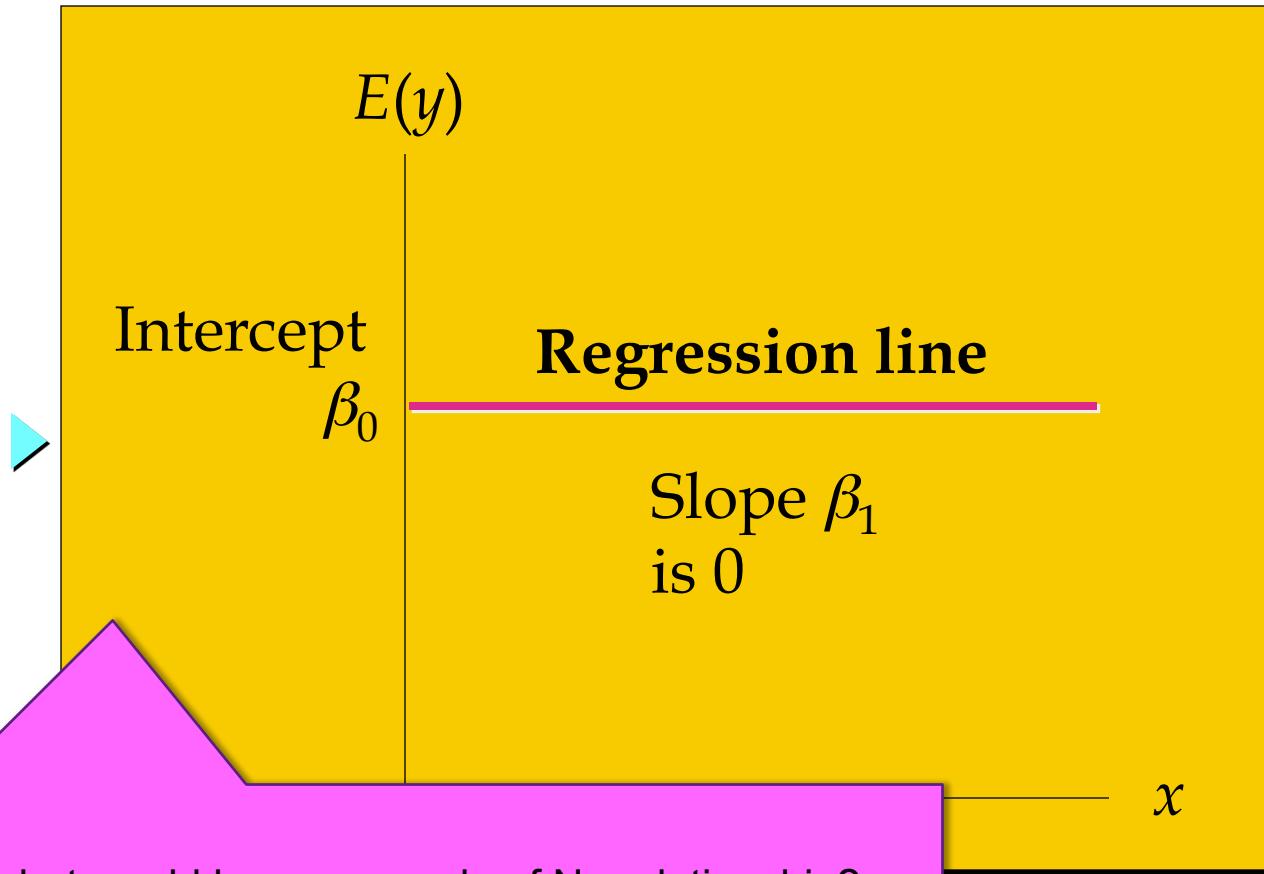
■ Negative Linear Relationship



When would be a case of a negative relationship? If I spend all my time watching movies instead of studying, would it increase my grade? Or decrease my grade?

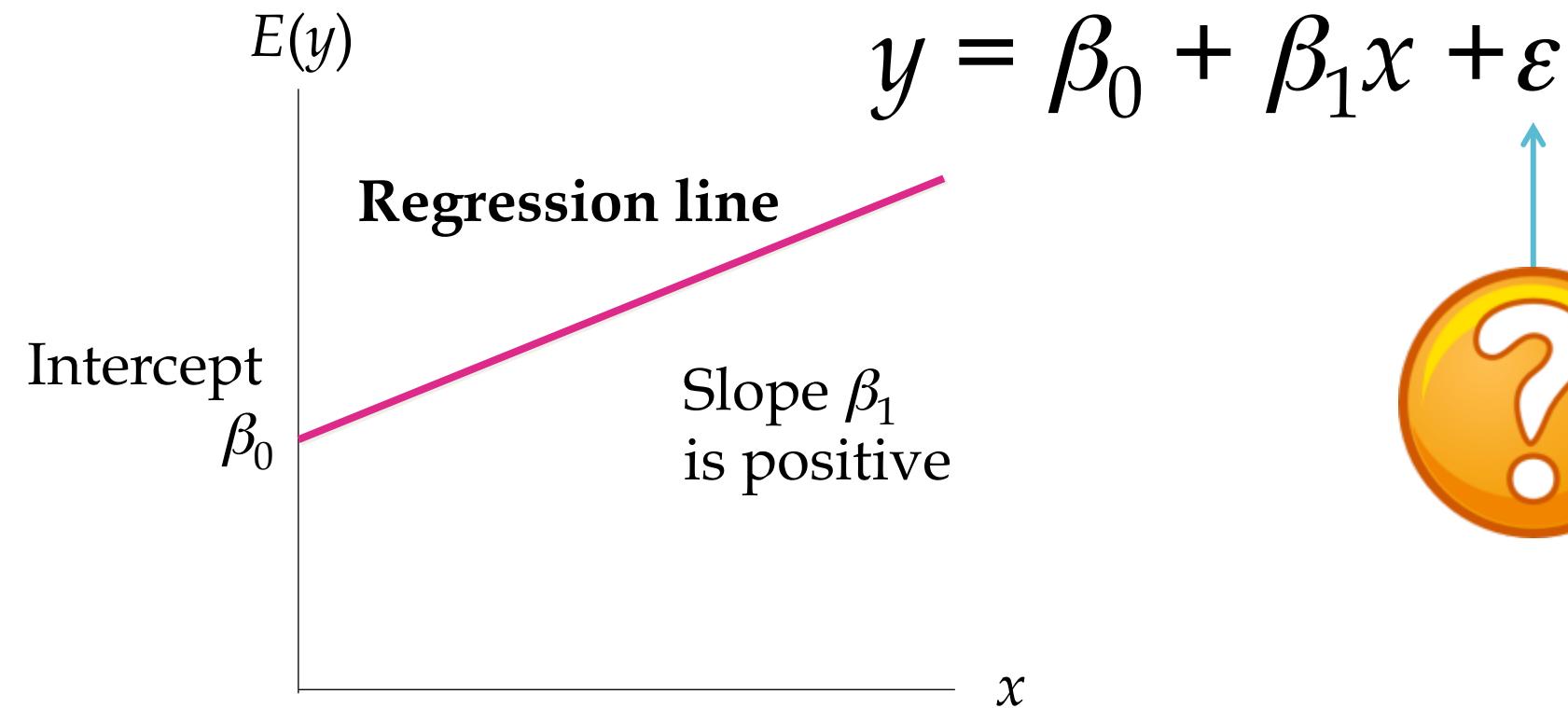
SIMPLE LINEAR REGRESSION EQUATION

■ No Relationship



Hmm... what would be an example of No relationship?
If my friend eats 5 ice creams everyday, does it have any
relationship with my grade? Not really right?

SIMPLE LINEAR REGRESSION MODEL



β_0 and β_1 – parameters of the model

ϵ is a random variable referred to as the error term.

ϵ – variability in y that cannot be explained

CHARACTERISTICS OF THE ERROR TERM

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

The more variables you add
The small ϵ will become
because now you know more!

$$\begin{aligned} y = & \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 \\ & + \beta_5 x_5 + \beta_6 x_6 + \beta_7 x_7 + \beta_8 x_8 \dots + \beta_i x_i \end{aligned}$$



CHARACTERISTICS OF THE ERROR TERM

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

Cannot be explained!

Cannot be calculated!

The more variables you add
The small ϵ will become
because now you know more!

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5 + \beta_6 x_6 + \beta_7 x_7 + \beta_8 x_8 \dots + \beta_i x_i$$

Cannot.. Just don't ask sigh...

REGRESSION EQUATION



Describes how the expected value of y , $E(y)$ is related to x

$$E(y) = \beta_0 + \beta_1 x$$

- Graph of the regression equation is a straight line.
- β_0 is the y intercept of the regression line.
- β_1 is the slope of the regression line.
- $E(y)$ is the expected value of y for a given x value.

ESTIMATED SIMPLE LINEAR REGRESSION EQUATION

$$E(y) = \beta_0 + \beta_1 x$$

Sample statistics are

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

- The graph is called the estimated regression line.
- b_0 is the y intercept of the line.
- b_1 is the slope of the line.
- \hat{y} is the estimated value of y for a given x value.

Least Square Method

Coefficient of Determination

Model Assumptions

Testing for Significance

Covariance & Coefficient of Correlation

Using the Estimated Regression –Equation for Estimation and Prediction

2. LEAST SQUARE METHOD

Is a procedure for using **sample data** to find the estimated regression equation.

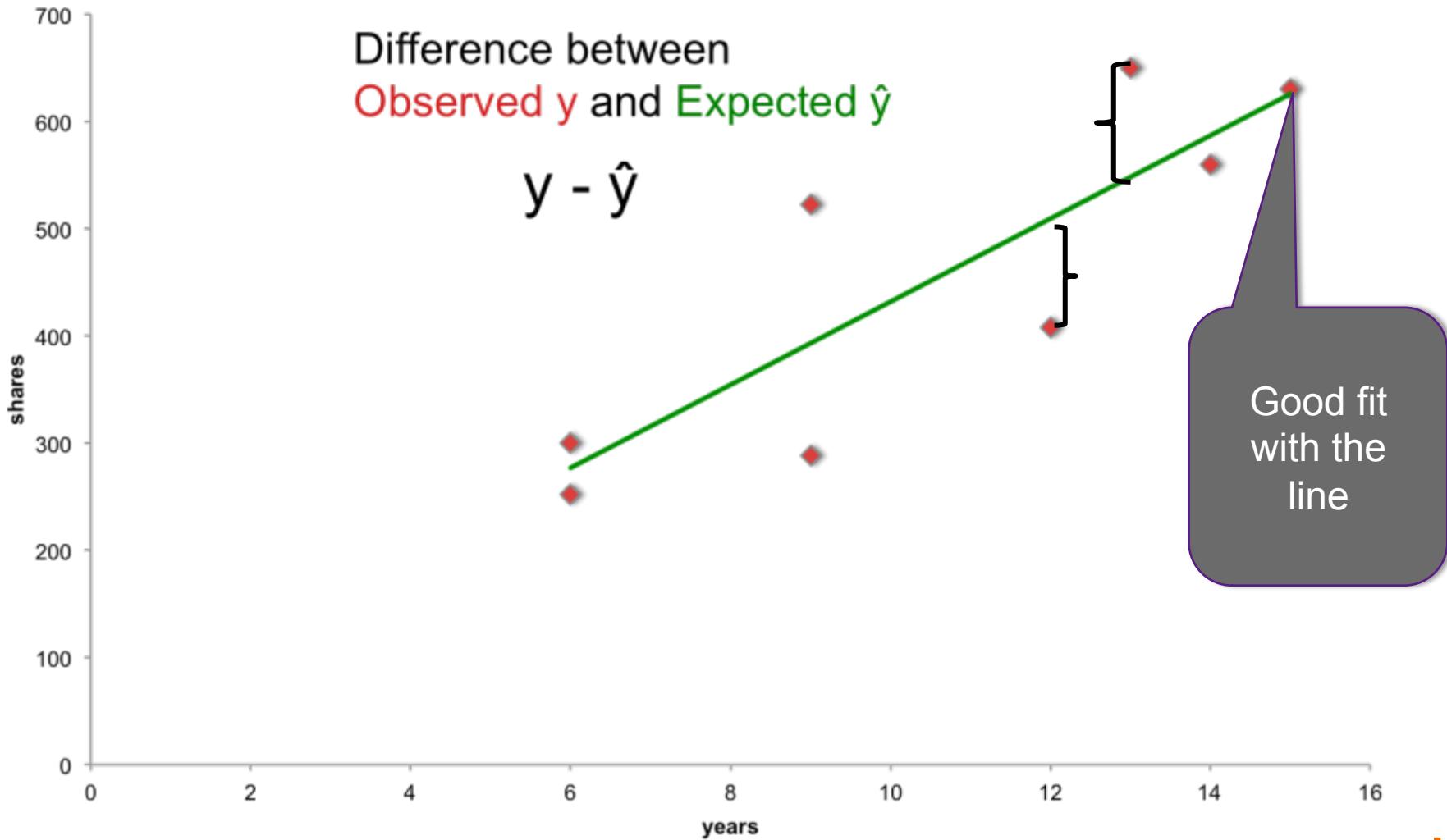
$$\hat{y} = b_0 + b_1 x$$

The goal to using the least square method

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y}_i)}{\sum (x_i - \bar{x})^2} = \frac{\sum x_i y_i - \sum x_i \sum y_i / n}{\sum x_i^2 - (\sum x_i)^2 / n}$$

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

Ownership of company stock vs years with the firm



LEAST SQUARES METHOD

Least Squares Criterion

$$\min \sum (y_i - \hat{y}_i)^2$$

where:

y_i = observed value of the dependent variable
for the i th observation

\hat{y}_i = estimated value of the dependent variable
for the i th observation

LEAST SQUARES METHOD

Slope for the Estimated Regression Equation

►
$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

where:

x_i = value of independent variable for i th observation

y_i = value of dependent variable for i th observation

\bar{x} = mean value for independent variable

\bar{y} = mean value for dependent variable

LEAST SQUARES METHOD

- y -Intercept for the Estimated Regression Equation

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

Coefficient of Determination

Model Assumptions

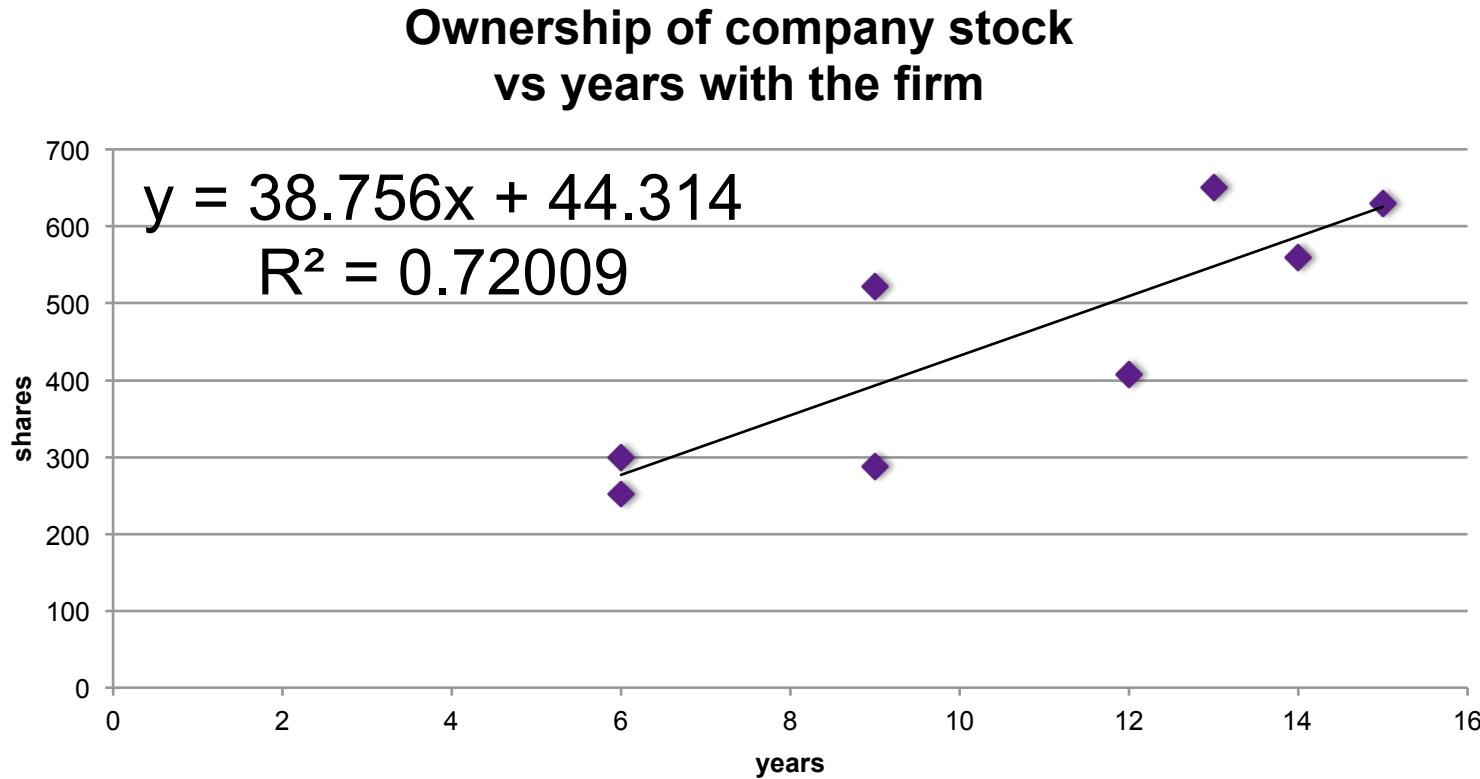
Testing for Significance

Covariance & Coefficient of Correlation

Using the Estimated Regression –Equation for Estimation and Prediction

3. COEFFICIENT OF DETERMINATION

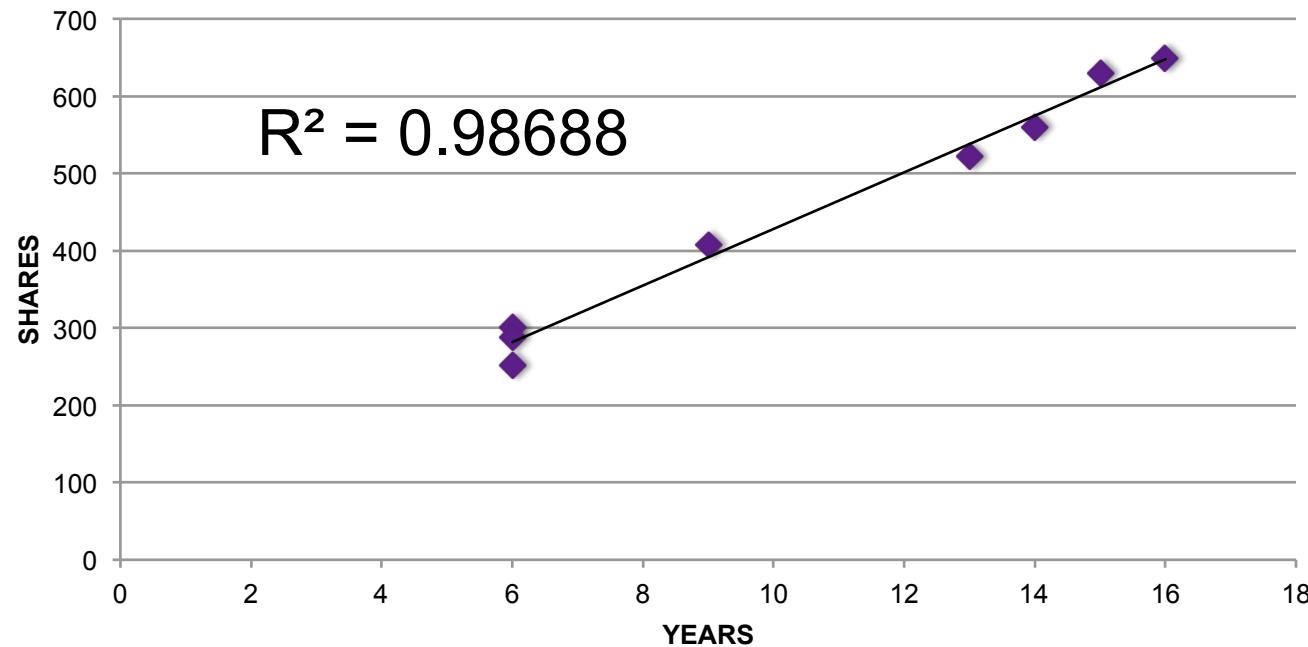
CORRELATION COEFFICIENT



COEFFICIENT OF DETERMINATION

x	y
6	300
9	408
14	560
6	252
6	288
16	650
15	630
13	522

Ownership of company stock
vs years with the firm



COEFFICIENT OF DETERMINATION

PROVIDES A MEASURE OF THE GOODNESS OF FIT FOR THE ESTIMATED REGRESSION EQUATION

**IN OTHER WORDS
DOES THE INDEPENDENT VARIABLE
EXPLAIN
THE DEPENDENT VARIABLE WELL?**

Explain 40% of
the grade

FINAL COMM 215 Grade

(y) =

Study Time (x_1) +

Work Time (x_2) +

**# of hours spent on Facebook
(x_3) +**

Family Time (x_4) +

**anything else you can think of
($x_{...}$) ...**

Explain 15% of
the grade

Explain 9% of
the grade

Explain 15% of
the grade

Explain 40% of
the grade

FINAL COMM 215 Grade

(y) =

These 4 Variables
explain
79% of your
grade!

Explain 15% of
the grade

Study Time (x_1) +

Work Time (x_2) +

**# of hours spent on Facebook
(x_3) +**

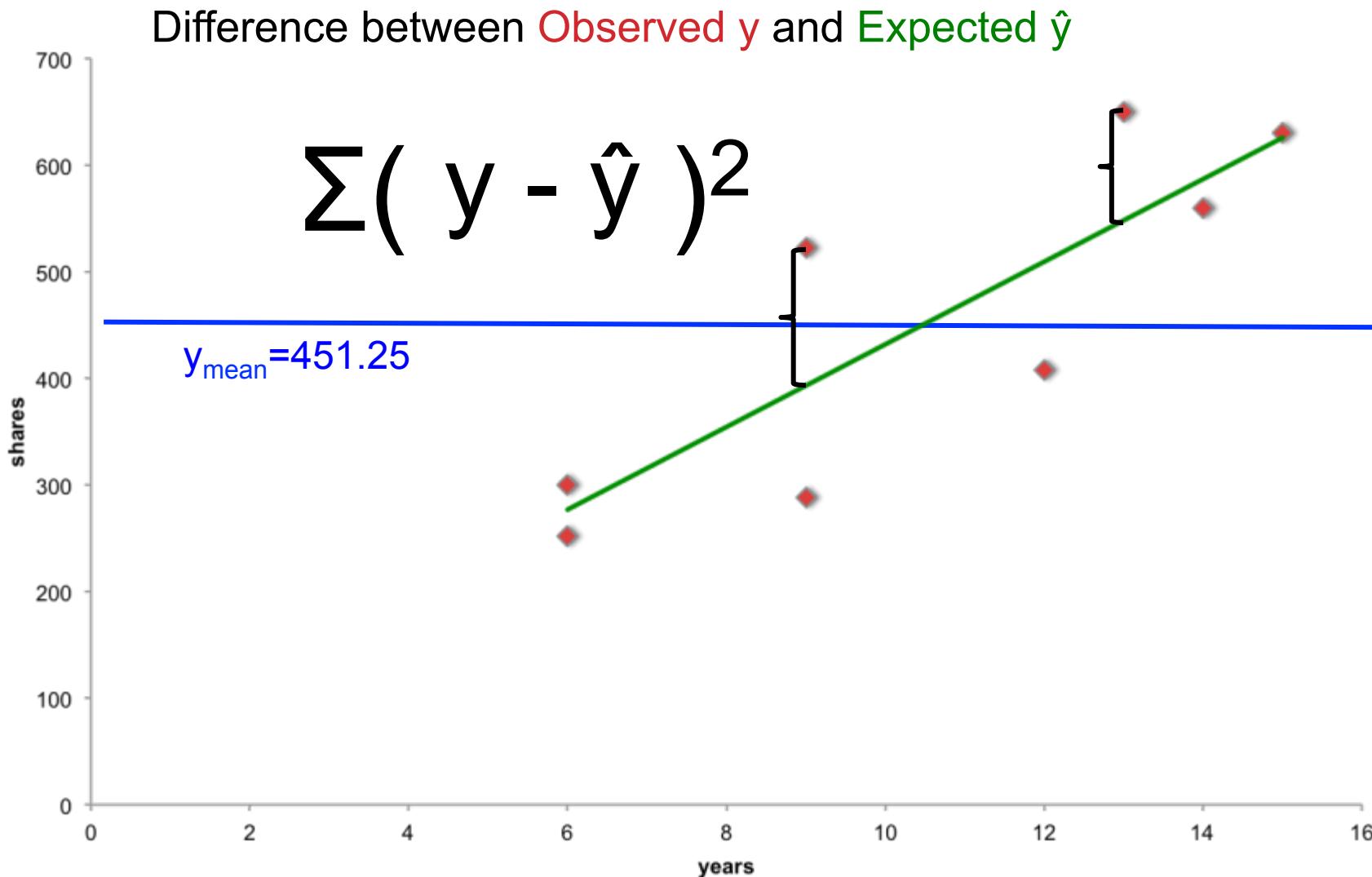
Family Time (x_4) +

Explain 15% of
the grade

Explain 9% of
the grade

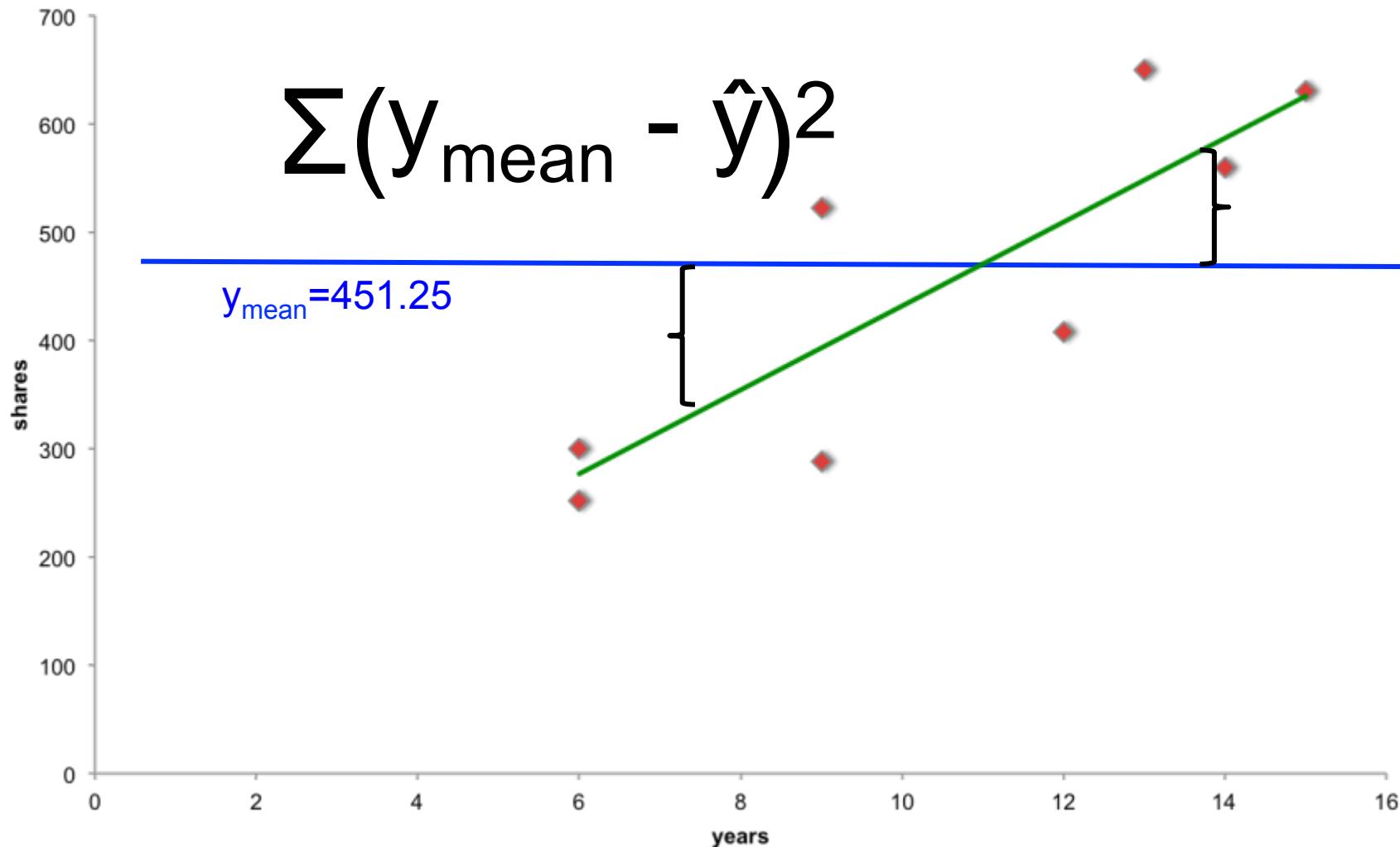
**anything else you can think of
($x_{...}$) ...**

Sum of Squares due to Error (SSE)



Sum of Squares due to Regression (SSR)

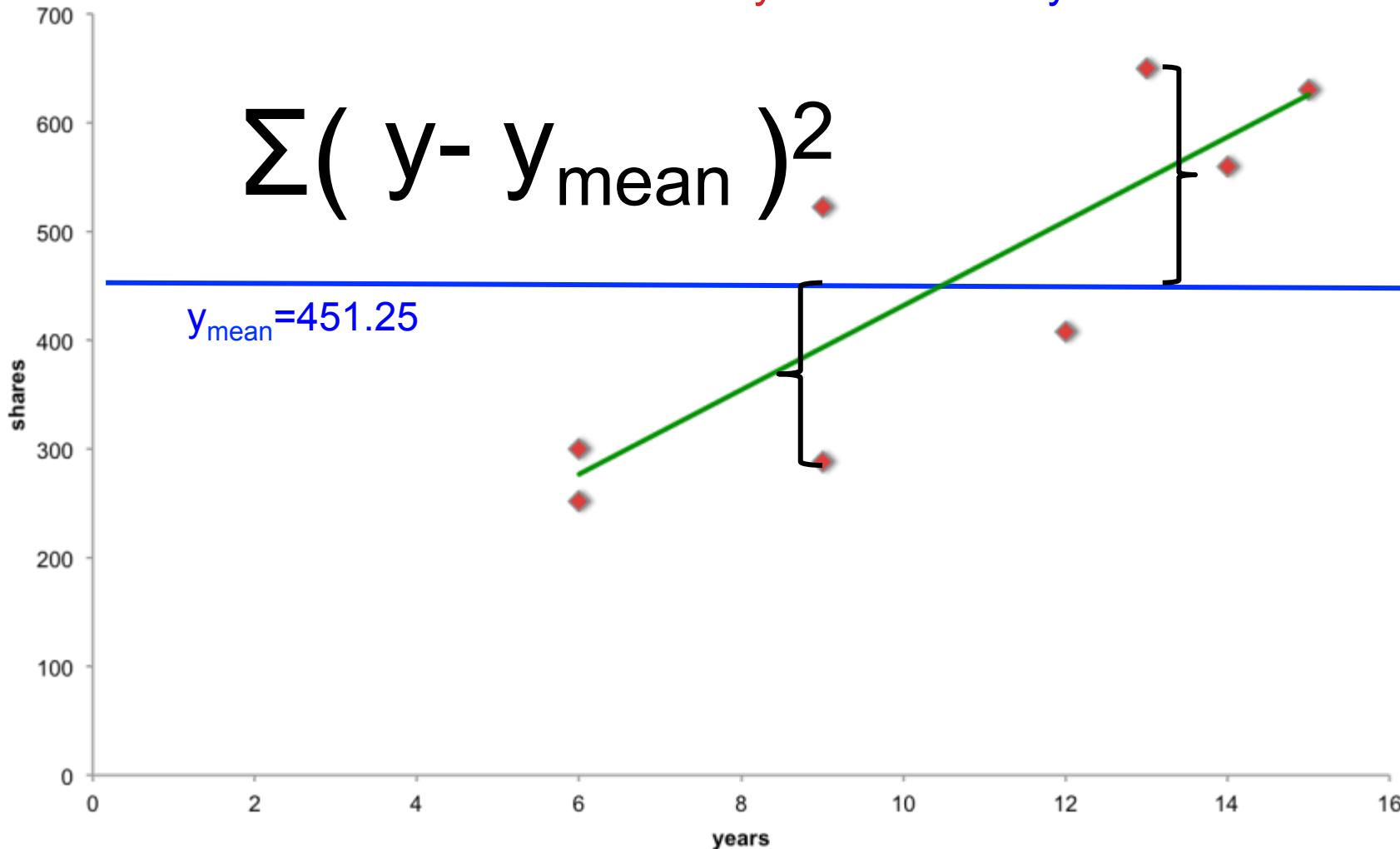
Difference between Mean of y and Expected \hat{y}



Bowerman, et al. (2017) pp. 546

Total Sum of Squares (SST)

Difference between Observed y and Mean of y



COEFFICIENT OF DETERMINATION

Relationship Among SST, SSR, SSE



$$\text{SST} = \text{SSR} + \text{SSE}$$

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$$

where:

SST = total sum of squares

SSR = sum of squares due to regression

SSE = sum of squares due to error

COEFFICIENT OF DETERMINATION

- The coefficient of determination is:


$$r^2 = \text{SSR}/\text{SST}$$

where:

SSR = sum of squares due to regression
= explained variation

SST = total sum of squares
= total variation

COEFFICIENT OF DETERMINATION (R^2)

Expresses proportion of the variation in the dependent variable (y) that is explained by the regression line:

$$\hat{y} = b_0 + b_1 x_1$$

COEFFICIENT OF CORRELATION (R)

Describes both the direction and the strength of the linear relationship between two variables

$$r = (\text{sign of } b_1) \sqrt{r^2}$$

Model Assumptions

Testing for Significance

Covariance & Coefficient of Correlation

**Using the Estimated Regression –Equation for Estimation
and Prediction**

4. MODEL ASSUMPTIONS

Before conducting regression analysis...

What determines an appropriate model for the relationship between the dependent and the independent variable(s).

Even if the data fits well , the estimated regression equation should not be used until further analysis of

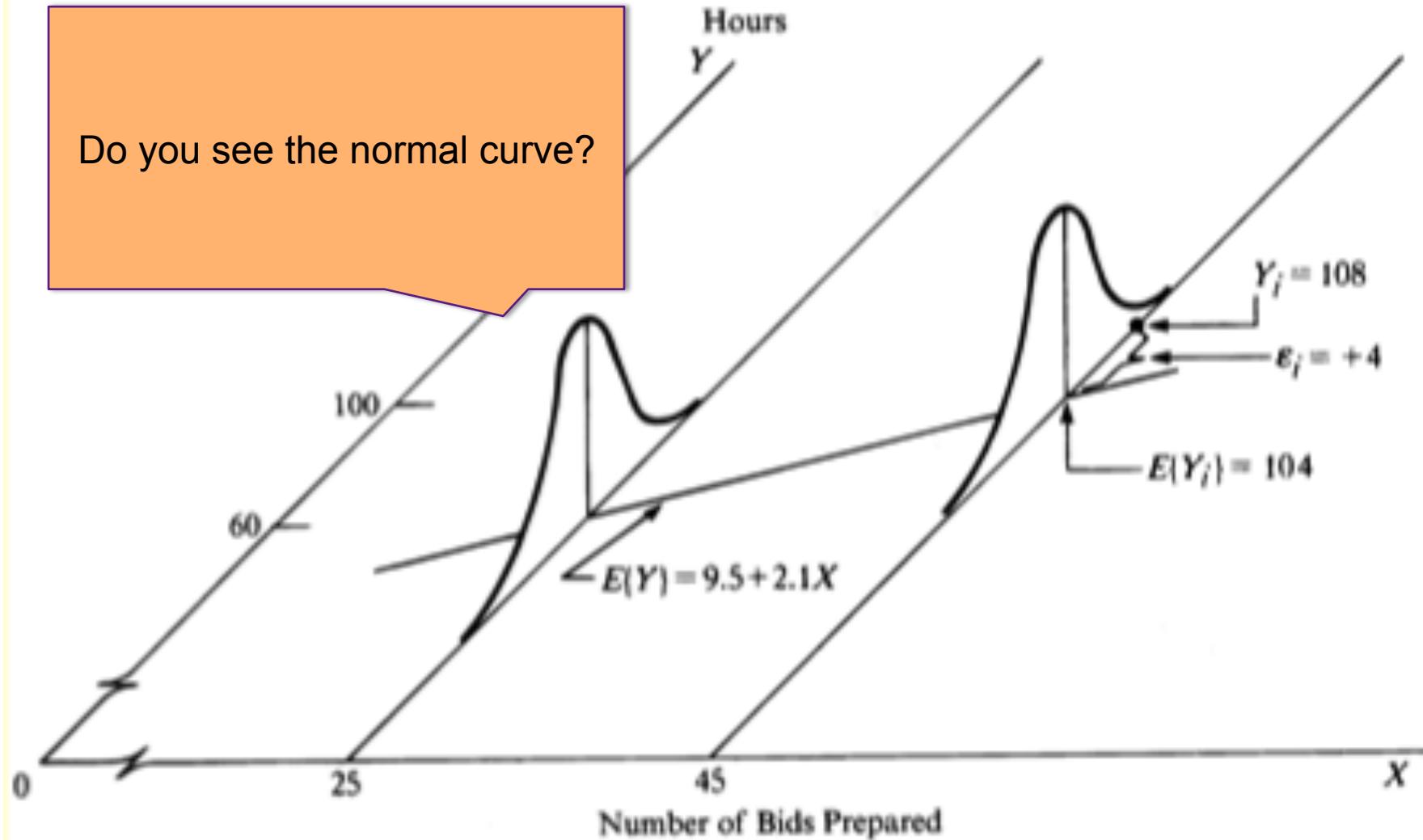
how appropriate the assumed model is.

One way to determine is to test for significance.

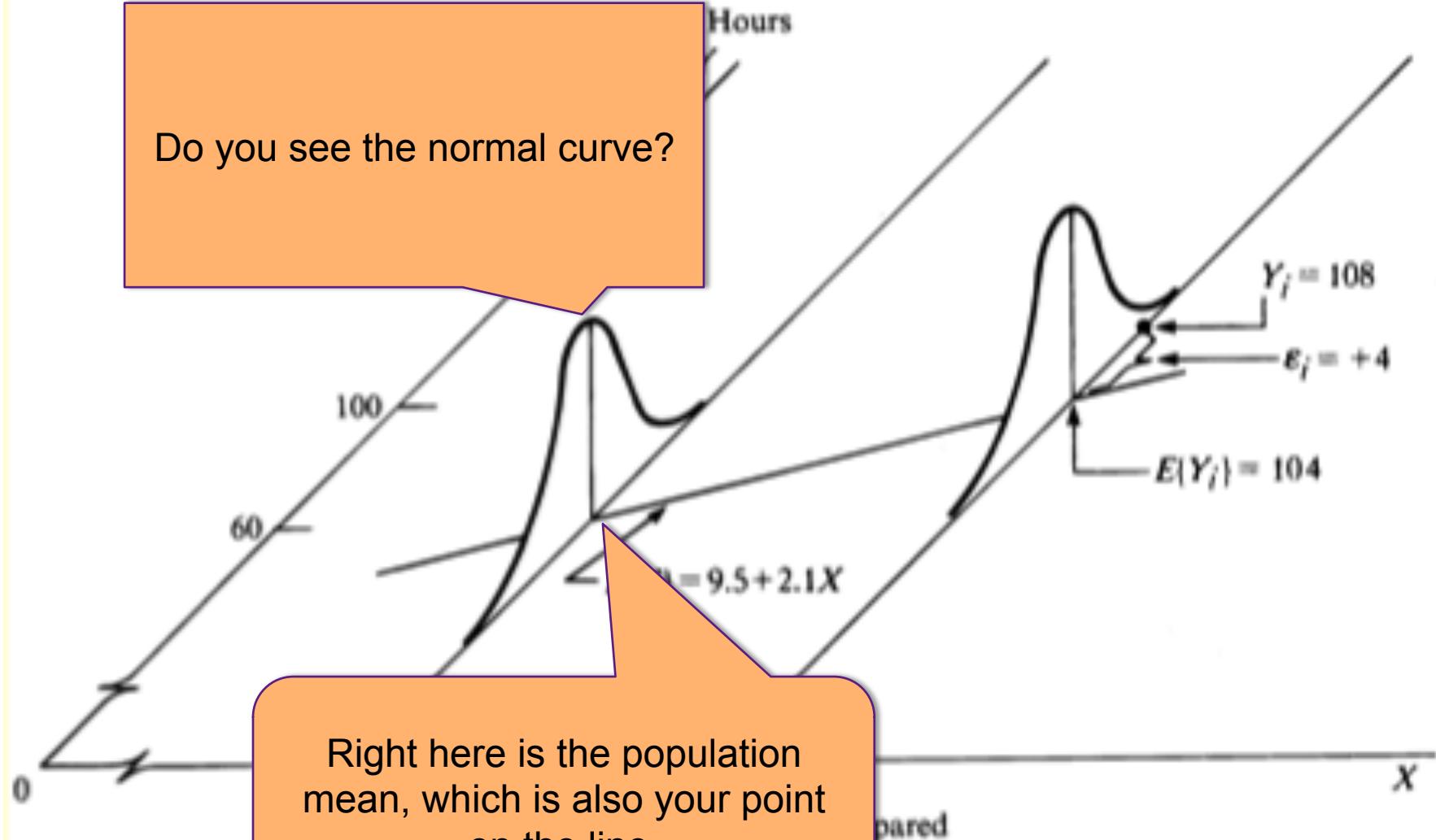
ASSUMPTIONS- ERROR TERM

1. The error term ε is a random variable with an expected value of zero. In estimating an element that is unpredictable, it is best to assume it to be zero.
2. The variance of ε , denoted by σ^2 is the same for all values of x.
3. The values of error are independent.
4. The error term is normally distributed.

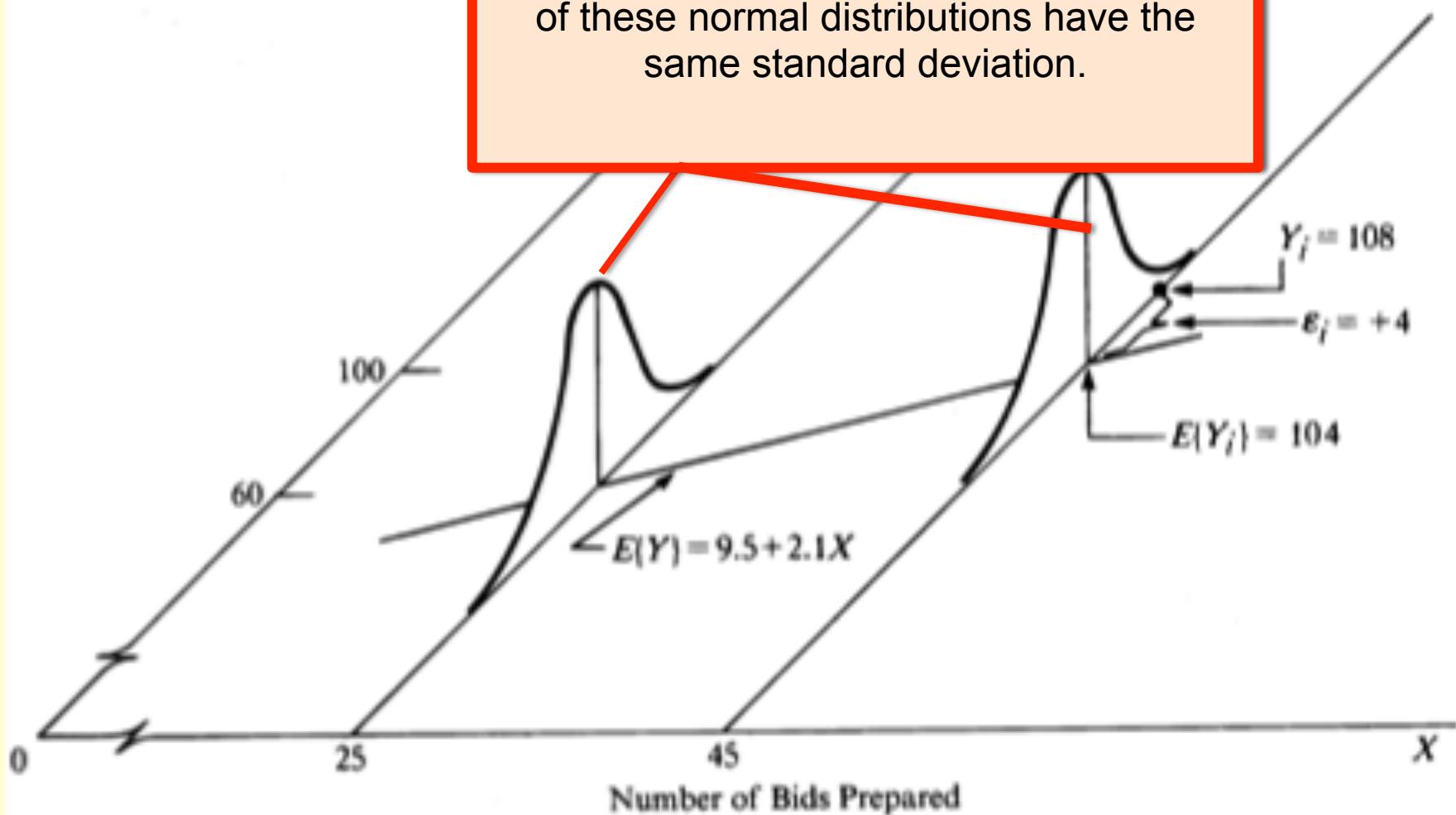
Do you see the normal curve?



Do you see the normal curve?

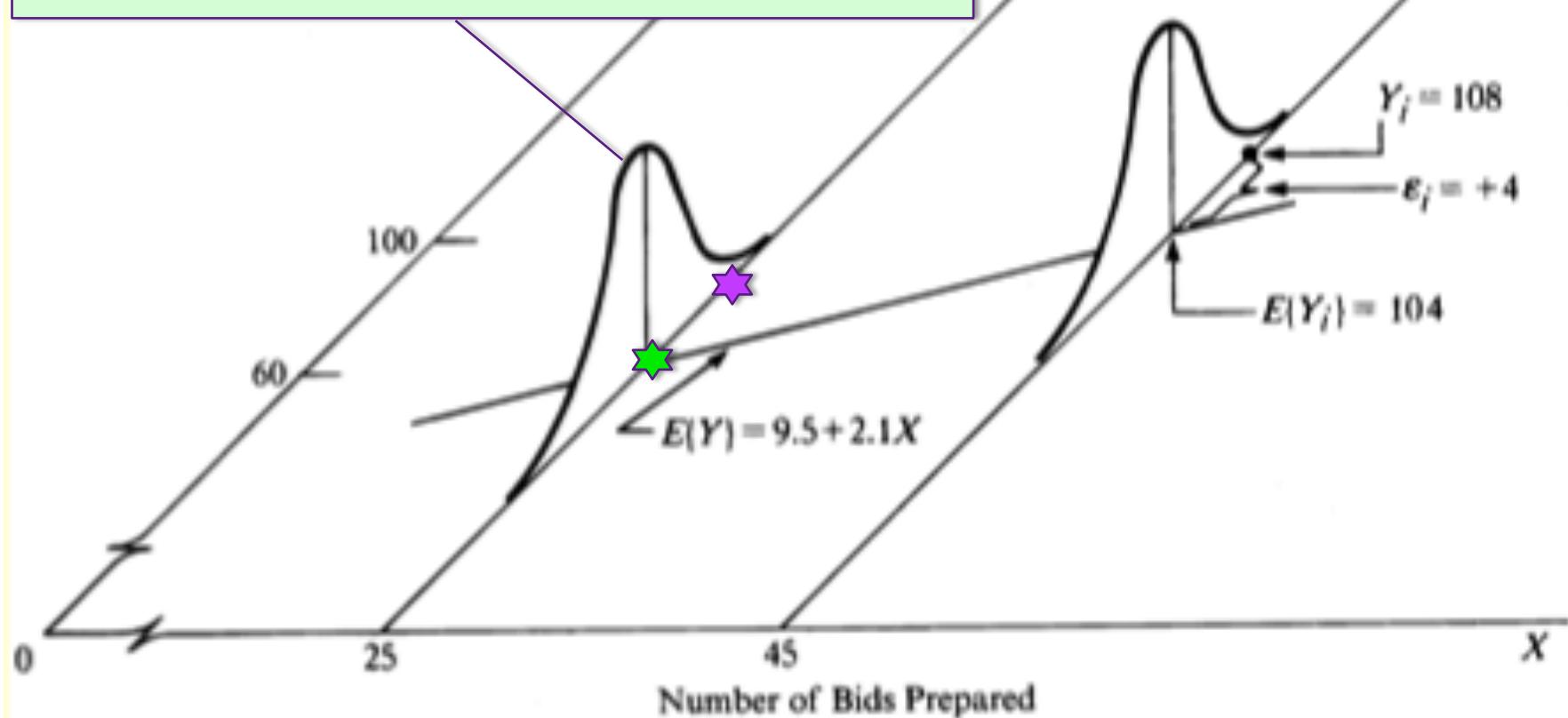


Assumption # 2: The variance/ standard deviation for each value x is the same. All of these normal distributions have the same standard deviation.



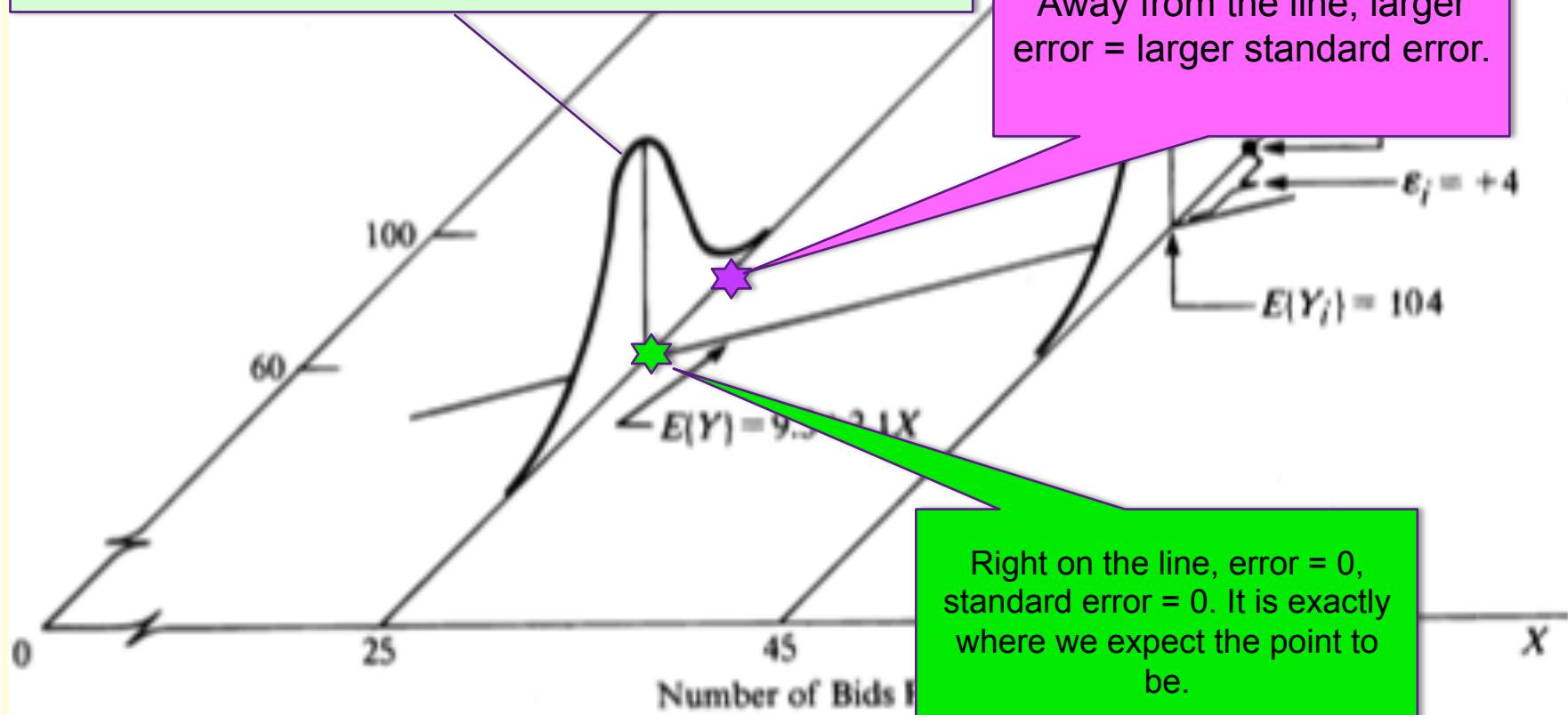
Assumptions # 3 & 4: the error terms are independent and are normally distributed.

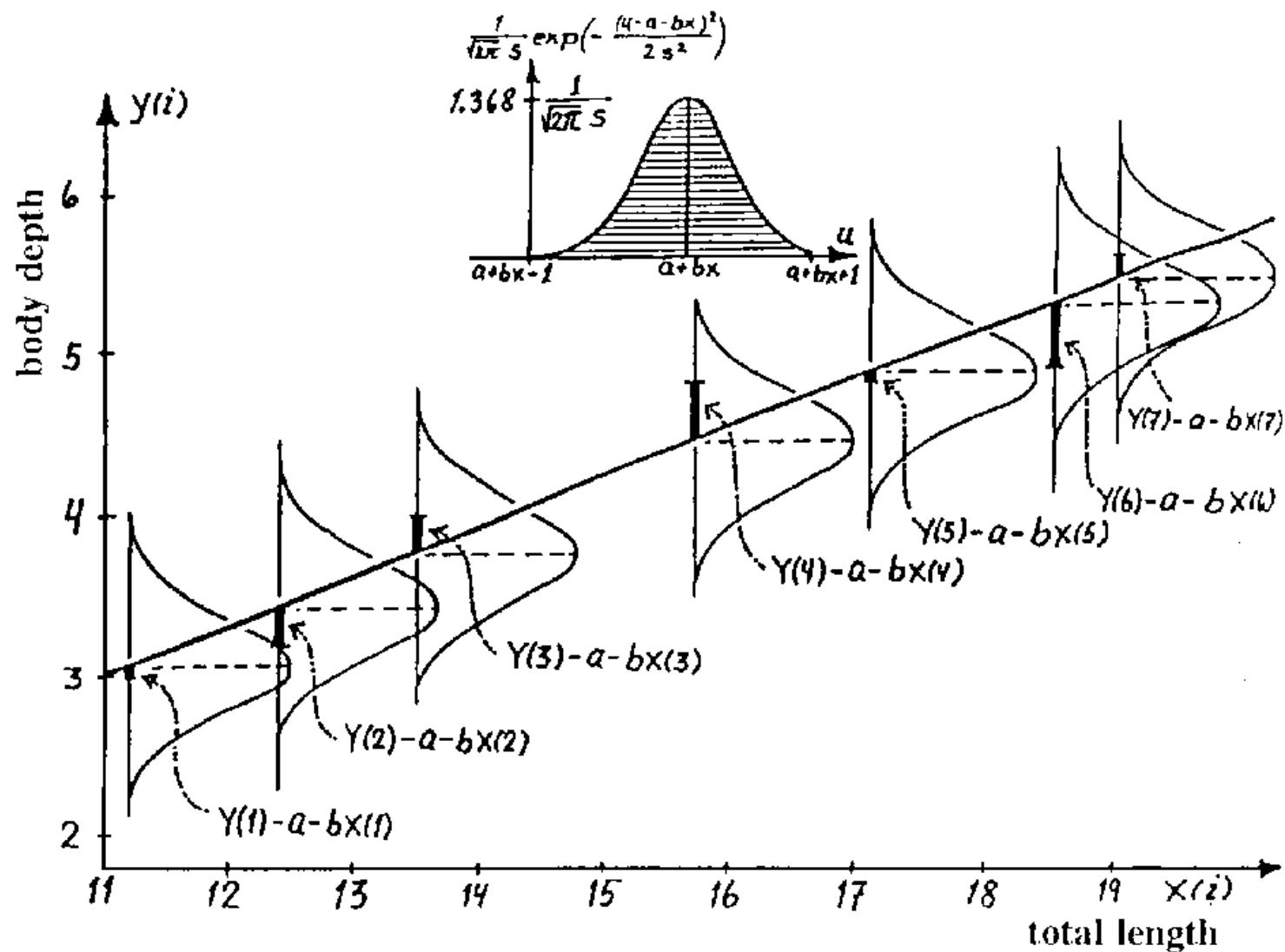
If the point is right on the line, then you have an error of 0. If your point is away from the line, the further it is away, the larger the error term is.



Assumptions # 3 & 4: the error terms are independent and are normally distributed.

If the point is right on the line, then you have an error of 0. If your point is away from the line, the further it is away, the larger the error term is.





Testing for Significance

Covariance & Coefficient of Correlation

Using the Estimated Regression –Equation for Estimation and Prediction

5. TESTING FOR SIGNIFICANCE

ESTIMATE σ^2

TESTING FOR SIGNIFICANCE

CONFIDENCE INTERVAL FOR B_1

TESTING FOR SIGNIFICANCE

An Estimate of σ^2

- ▶ The mean square error (MSE) provides the estimate of σ^2 , and the notation s^2 is also used.

$$s^2 = \text{MSE} = \text{SSE}/(n - 2)$$

- ▶ where:

$$\text{SSE} = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - b_0 - b_1 x_i)^2$$

$$\text{SSE} = \sum y_i^2 - b_0 \sum y_i - b_1 \sum x_i y_i$$

TESTING FOR SIGNIFICANCE

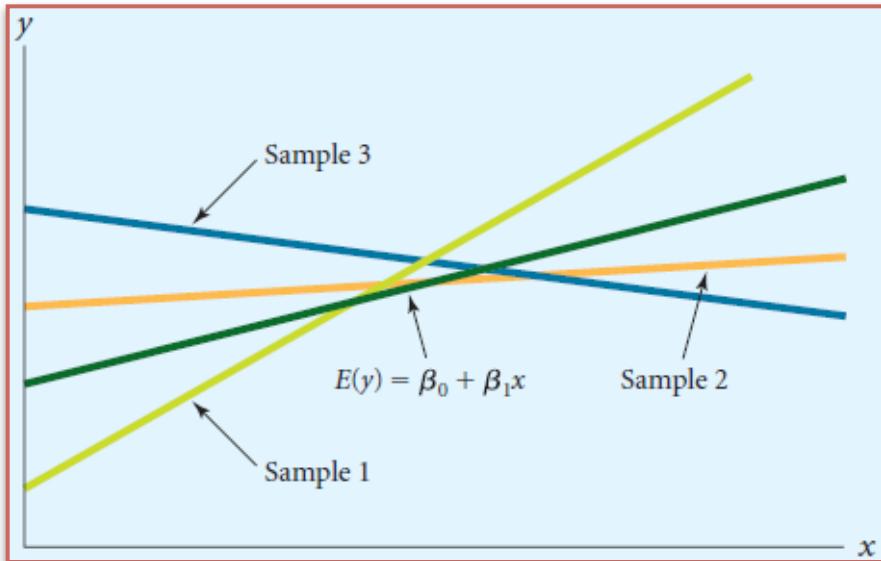
An Estimate of σ

- ▶ • To estimate σ we take the square root of σ^2 .
- ▶ • The resulting s is called the standard error of the estimate.

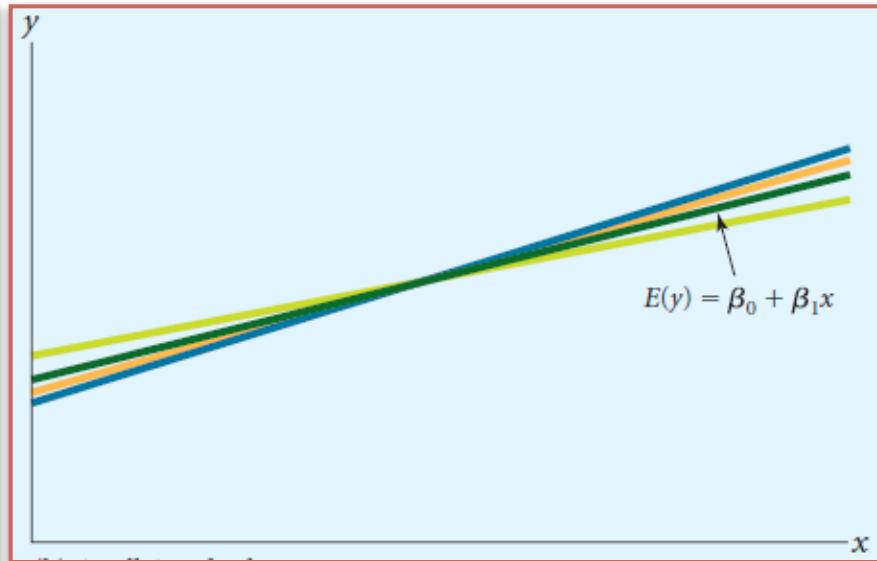
$$s = \sqrt{\text{MSE}} = \sqrt{\frac{\text{SSE}}{n - 2}}$$

Why is this $n-2$? it actually is $n-k-1$.
In a simple linear regression, since you always have 1 independent variable ($k=1$), therefore automatically it becomes $n-1-1$.

STANDARD ERROR OF THE ESTIMATE



Large Standard Error



Small Standard Error

TESTING FOR SIGNIFICANCE

- ▶ To test for a significant regression relationship, we must conduct a hypothesis test to determine whether the value of β_1 is zero.
- ▶ Two tests are commonly used:
 t Test and F Test
- ▶ Both the t test and F test require an estimate of σ^2 , the variance of ε in the regression model.

TESTING FOR SIGNIFICANCE: T TEST

Hypotheses


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

Test Statistic



$$t = \frac{b_1}{s_{b_1}}$$

where

$$s_{b_1} = \frac{s}{\sqrt{\sum(x_i - \bar{x})^2}}$$

TESTING FOR SIGNIFICANCE: *T* TEST

1. Set up Hypotheses. $H_0: \beta_1 = 0$ $H_a: \beta_1 \neq 0$
2. What is the appropriate test statistic to use?. $t = \frac{b_1}{S_{b_1}}$
3. Calculate the test statistic value.
4. Find the critical value for the test statistic. $\alpha = .05$
5. Define the decision rule
6. Make your decision
7. Interpret the conclusion in context

CONFIDENCE INTERVAL FOR B_1

The form of a confidence interval for β_1 is:

$$b_1 \pm t_{\alpha/2} s_{b_1}$$

b_1 is the point estimator

$t_{\alpha/2} s_{b_1}$ is the margin of error

where $t_{\alpha/2}$ is the t value providing an area of $\alpha/2$ in the upper tail of a t distribution with $n - 2$ degrees of freedom

CONFIDENCE INTERVAL FOR B_1

- ▶ ■ We can use a 95% confidence interval for β_1 to test the hypotheses just used in the t test.
- ▶ ■ H_0 is rejected if the hypothesized value of β_1 is not included in the confidence interval for β_1 .

CONFIDENCE INTERVAL FOR B_1

► ■ Rejection Rule

Reject H_0 if 0 is not included in the confidence interval for β_1 .

► ■ 95% Confidence Interval for β_1

$$b_1 \pm t_{\alpha/2} s_{b_1} = 5.0 \pm 2.048(2.25) = 5.0 \pm 4.608$$

or 0.392 to 9.608

► ■ Conclusion

0 is not included in the confidence interval.

Reject H_0

SOME CAUTIONS ABOUT THE INTERPRETATION OF SIGNIFICANCE TESTS

- ▶ ■ Rejecting $H_0: \beta_1 = 0$ and concluding that the relationship between x and y is significant does not enable us to conclude that a cause-and-effect relationship is present between x and y .
- ▶ ■ Just because we are able to reject $H_0: \beta_1 = 0$ and demonstrate statistical significance does not enable us to conclude that there is a linear relationship between x and y .

What is a standard error of the estimate?

What is a standard error of the slope?

$$s = \sqrt{\text{MSE}} = \sqrt{\frac{\text{SSE}}{n - 2}}$$

$$s_{b_1} = \frac{s}{\sqrt{\sum(x_i - \bar{x})^2}}$$

Testing for Significance: t Test

1. Set up Hypotheses.
2. What is the appropriate test statistic to use?.
$$t = \frac{b_1}{S_{b_1}}$$
3. Calculate the test statistic value.
4. Find the critical value for the test statistic. $\alpha = .05$
5. Define the decision rule
6. Make your decision
7. Interpret the conclusion in context

3. Calculate the test statistic value.

5

$$t = \frac{b_1}{s_{b_1}}$$

4

$$s_{b_1} = \frac{s}{\sqrt{\sum(x_i - \bar{x})^2}}$$

3

$$\sum(x_i - \bar{x})^2 = (\sum x_i^2) - \frac{(\sum x_i)^2}{n} = SS_{xx}$$

2

$$s_{\varepsilon} = \sqrt{\frac{SSE}{n-k-1}} \quad \text{Standard Error of Estimate (} s_{\varepsilon} \text{)}$$

1

SSE

PROBLEM # 10.5

At 5% level significance. The manager of Colonial Furniture has been reviewing weekly advertising expenditures. During the past 6 months, all advertisements for the store have appeared in the local newspaper. The number of ads per week has varied from one to seven. The store's sales staff has been tracking the number of customers who enter the store each week. The number of ads and the number of customers per week for the past 26 weeks were recorded.

- c. Can the manager infer that the larger the number of ads, the larger the number of customers?

$$SSE = \sum(y - \hat{y})^2 = 424281.17$$

Ads x	Customer y	$\hat{y} = 296.92 + 21.356x$	$y - \hat{y}$	$(y - \hat{y})^2$
5.00	353.00	403.70	-50.70	2570.49
6.00	319.00	425.06	-106.06	11247.88
3.00	440.00	360.99	79.01	6242.90
2.00	332.00	339.63	-7.63	58.25
4.00	172.00	382.34	-210.34	44244.60
2.00	331.00	339.63	-8.63	74.51
4.00	344.00	382.34	-38.34	1470.26
2.00	483.00	339.63	143.37	20554.38
4.00	329.00	382.34	-53.34	2845.58
2.00	532.00	339.63	192.37	37005.45
7.00	496.00	446.41	49.59	2458.97
5.00	393.00	403.70	-10.70	114.49
4.00	376.00	382.34	-6.34	40.25
7.00	372.00	446.41	-74.41	5537.15
2.00	512.00	339.63	172.37	29710.73
5.00	254.00	403.70	-149.70	22410.09
5.00	459.00	403.70	55.30	3058.09
2.00	153.00	339.63	-186.63	34831.50
1.00	426.00	318.28	107.72	11604.46
6.00	566.00	425.06	140.94	19865.21
6.00	596.00	425.06	170.94	29221.85
5.00	395.00	403.70	-8.70	75.69
6.00	676.00	425.06	250.94	62972.89
3.00	194.00	360.99	-166.99	27884.99
2.00	135.00	339.63	-204.63	41874.26
7.00	367.00	446.41	-79.41	6306.27
424281.17				

$$SSE = \sum(y - \hat{y})^2$$

PROBLEM # 10.5

At 5% level significance. The manager of Colonial Furniture has been reviewing weekly advertising expenditures. During the past 6 months, all advertisements for the store have appeared in the local newspaper. The number of ads per week has varied from one to seven. The store's sales staff has been tracking the number of customers who enter the store each week. The number of ads and the number of customers per week for the past 26 weeks were recorded.

- c. Can the manager infer that the larger the number of ads, the larger the number of customers?

$$SSE = \sum(y - \hat{y})^2 = 424281.17$$

$$s_e = \sqrt{\frac{SSE}{n-k-1}}$$

$$s_e = \sqrt{\frac{424281.17}{26-1-1}}$$

$$s_e = 132.96$$

PROBLEM # 10.5

At 5% level significance. The manager of Colonial Furniture has been reviewing weekly advertising expenditures. During the past 6 months, all advertisements for the store have appeared in the local newspaper. The number of ads per week has varied from one to seven. The store's sales staff has been tracking the number of customers who enter the store each week. The number of ads and the number of customers per week for the past 26 weeks were recorded.

- c. Can the manager infer that the larger the number of ads, the larger the number of customers?

4

$$S_{b_1} = \frac{s}{\sqrt{\sum(x_i - \bar{x})^2}} = 132.96$$

$$S_{b_1} = \frac{132.96}{\sqrt{86.6544}} = 14.28$$

3

Ads x	x-xbar	(x-xbar) ²
5.00	0.88	0.7744
6.00	1.88	3.5344
3.00	-1.12	1.2544
2.00	-2.12	4.4944
4.00	-0.12	0.0144
2.00	-2.12	4.4944
4.00	-0.12	0.0144
2.00	-2.12	4.4944
4.00	-0.12	0.0144
2.00	-2.12	4.4944
7.00	2.88	8.2944
5.00	0.88	0.7744
4.00	-0.12	0.0144
7.00	2.88	8.2944
2.00	-2.12	4.4944
5.00	0.88	0.7744
5.00	0.88	0.7744
2.00	-2.12	4.4944
1.00	-3.12	9.7344
6.00	1.88	3.5344
6.00	1.88	3.5344
5.00	0.88	0.7744
6.00	1.88	3.5344
3.00	-1.12	1.2544
2.00	-2.12	4.4944
7.00	2.88	8.2944
4.12		86.6544

$$SS_{xx} =$$

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

PROBLEM # 10.5

At 5% level significance. The manager of Colonial Furniture has been reviewing weekly advertising expenditures. During the past 6 months, all advertisements for the store have appeared in the local newspaper. The number of ads per week has varied from one to seven. The store's sales staff has been tracking the number of customers who enter the store each week. The number of ads and the number of customers per week for the past 26 weeks were recorded.

- c. Can the manager infer that the larger the number of ads, the larger the number of customers?

$$s_{b_1} = \frac{s}{\sqrt{\sum(x_i - \bar{x})^2}} = 14.28$$

$$\hat{y} = 296.92 + 21.356x$$

$$t = \frac{b_1}{s_{b_1}} = + 21.356 / 14.28$$

$$t = 1.4955$$

Covariance & Coefficient of Correlation

Using the Estimated Regression –Equation for Estimation and Prediction

6. COVARIANCE & COEFFICIENT OF CORRELATION

Covariance

Interpretation of the covariance

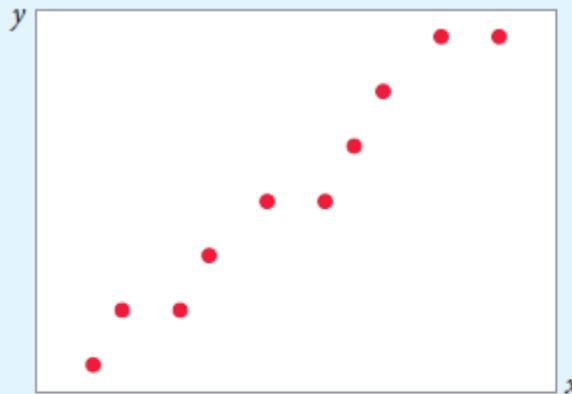
Correlation coefficient

MEASURES OF ASSOCIATION BETWEEN TWO VARIABLES

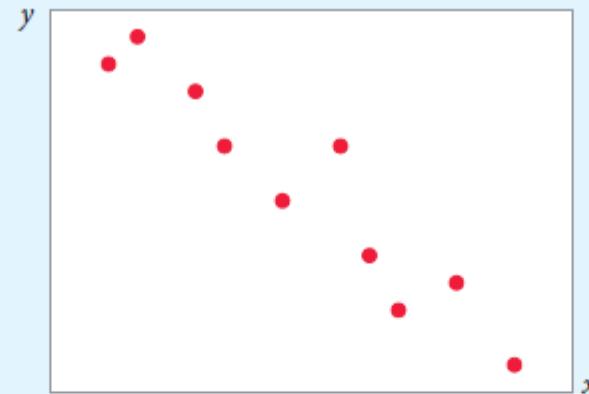
- ▶ Thus far we have examined numerical methods used to summarize the data for one variable at a time.
- ▶ Often a manager or decision maker is interested in the relationship between two variables.
- ▶ Two descriptive measures of the relationship between two variables are covariance and correlation coefficient.

2 VARIABLE RELATIONSHIPS

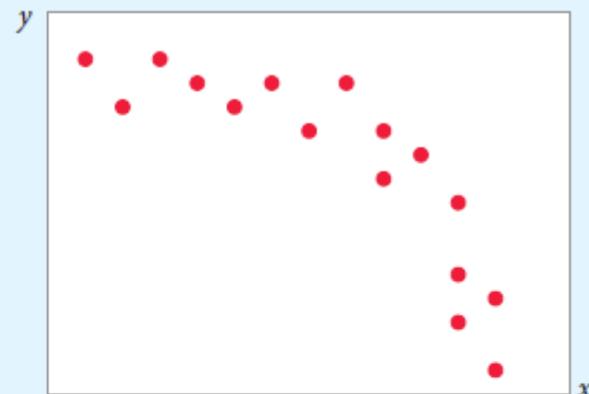
(a) Linear



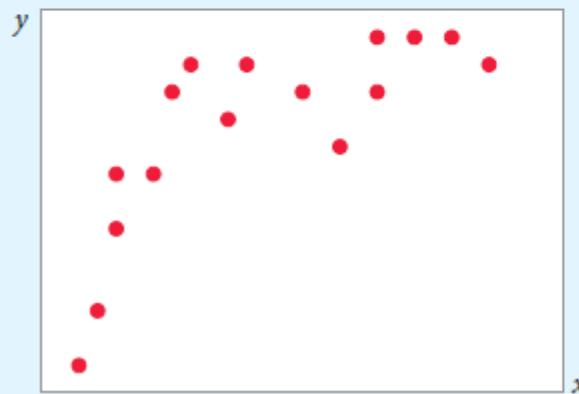
(b) Linear



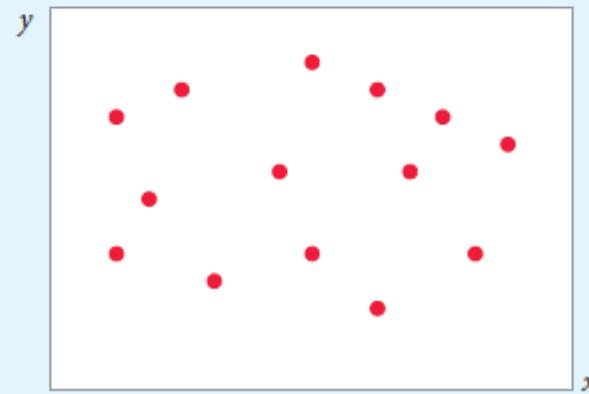
(c) Curvilinear



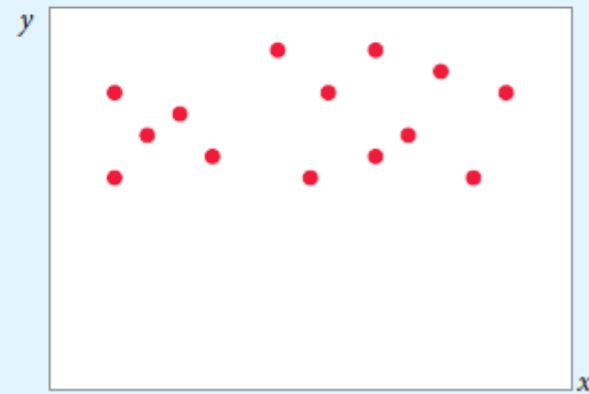
(d) Curvilinear



(e) No Relationship



(f) No Relationship



COVARIANCE

- ▶ The covariance is a measure of the linear association between two variables.
- ▶ Positive values indicate a positive relationship.
- ▶ Negative values indicate a negative relationship.

COVARIANCE

- The covariance is computed as follows:

$$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

for samples

$$\sigma_{xy} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{N}$$

for populations

CORRELATION COEFFICIENT

- ▶ Correlation is a measure of linear association and not necessarily causation.
- ▶ Just because two variables are highly correlated, it does not mean that one variable is the cause of the other.

CORRELATION COEFFICIENT

- The correlation coefficient is computed as follows:

$$r_{xy} = \frac{s_{xy}}{s_x s_y}$$

for
samples

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

for
populations



Pearson Product Moment
Correlation Coefficient.

CORRELATION COEFFICIENT

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{[\sum (x - \bar{x})^2][\sum (y - \bar{y})^2]}}$$

$$r = \frac{n \sum xy - \sum x \sum y}{\sqrt{[n(\sum x^2) - (\sum x)^2][n(\sum y^2) - (\sum y)^2]}}$$

r - Sample correlation coefficient

n - Sample size

x - Value of the independent variable

y - Value of the dependent variable

CORRELATION COEFFICIENT

- ▶ The coefficient can take on values between -1 and +1.
- ▶ Values near -1 indicate a strong negative linear relationship.
- ▶ Values near +1 indicate a strong positive linear relationship.
- ▶ The closer the correlation is to zero, the weaker the relationship.

Covariance and Correlation Coefficient

x	y	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$
277.6	69	10.65	-1.0	-10.65
259.5	71	-7.45	1.0	-7.45
269.1	70	2.15	0	0
267.0	70	0.05	0	0
255.6	71	-11.35	1.0	-11.35
272.9	69	5.95	-1.0	-5.95
Average	267.0	70.0		Total -35.40
Std. Dev.	8.2192	.8944		

Ads x	Customer y	x-xmean	(x-xmean)²	y-ymean	(y-ymean)²	(x-xmean)(y-ymean)
5.00	353.00	0.88	0.77	-31.81	1011.88	-27.993
6.00	319.00	1.88	3.53	-65.81	4330.96	-123.723
3.00	440.00	-1.12	1.25	55.19	3045.94	-61.813
2.00	332.00	-2.12	4.49	-52.81	2788.90	111.957
4.00	172.00	-0.12	0.01	-212.81	45288.10	25.537
2.00	331.00	-2.12	4.49	-53.81	2895.52	114.077
4.00	344.00	-0.12	0.01	-40.81	1665.46	4.897
2.00	483.00	-2.12	4.49	98.19	9641.28	-208.163
4.00	329.00	-0.12	0.01	-55.81	3114.76	6.697
2.00	532.00	-2.12	4.49	147.19	21664.90	-312.043
7.00	496.00	2.88	8.29	111.19	12363.22	320.227
5.00	393.00	0.88	0.77	8.19	67.08	7.207
4.00	376.00	-0.12	0.01	-8.81	77.62	1.057
7.00	372.00	2.88	8.29	-12.81	164.10	-36.893
2.00	512.00	-2.12	4.49	127.19	16177.30	-269.643
5.00	254.00	0.88	0.77	-130.81	17111.26	-115.113
5.00	459.00	0.88	0.77	74.19	5504.16	65.287
2.00	153.00	-2.12	4.49	-231.81	53735.88	491.437
1.00	426.00	-3.12	9.73	41.19	1696.62	-128.513
6.00	566.00	1.88	3.53	181.19	32829.82	340.637
6.00	596.00	1.88	3.53	211.19	44601.22	397.037
5.00	395.00	0.88	0.77	10.19	103.84	8.967
6.00	676.00	1.88	3.53	291.19	84791.62	547.437
3.00	194.00	-1.12	1.25	-190.81	36408.46	213.707
2.00	135.00	-2.12	4.49	-249.81	62405.04	529.597
7.00	367.00	2.88	8.29	-17.81	317.20	-51.293
4.12	384.81		86.65		463802.04	1850.577
			1.86		136.21	74.023
			3.466176		18552.081544	

Ads x	Customer y	x-xmean	(x-xmean)²	y-ymean	(y-ymean)²	(x-xmean) (y-ymean)
5.00	353.00	0.88	0.77	-31.81	1011.88	-27.993
...
6.00	566.00	1.88	3.53	181.19	32829.82	340.637
6.00	596.00	1.88	3.53	211.19	44601.22	397.037
5.00	395.00	0.88	0.77	10.19	103.84	8.967
6.00	676.00	1.88	3.53	291.19	84791.62	547.437
3.00	194.00	-1.12	1.25	-190.81	36408.46	213.707
2.00	135.00	-2.12	4.49	-249.81	62405.04	529.597
7.00	367.00	2.88	8.29	-17.81	317.20	-51.293
4.12	384.81		86.65		463802.04	1850.577
			1.86		136.21	74.023

Ads x	Customer y	x-xmean	(x-xmean) ²	y-ymean	(y-ymean) ²	(x-xmean)(y-ymean)
5.00	353.00	0.88	0.77	-31.81	1011.88	-27.993
6.00	310.00	1.88	3.53	-65.81	4330.96	-123.723
			1.25	55.19	3045.94	-61.813
			4.49	-52.81	2788.90	111.957
			0.01	-212.81	45288.10	25.537
			4.49	-53.81	2895.52	114.077
			0.01	77.62	4.897	
2.00	483.00	-2.12	4.49	-8.81	-208.163	
4.00	329.00	-0.12	0.01	77.62	6.697	
2.00	532.00	-2.12	4.49	1696.62	-312.043	
7.00	496.00	2.88	8.29	32829.82	320.227	
5.00	393.00	0.88	0.77	44601.22	7.207	
4.00	376.00	-0.12	0.01	103.84	1.057	
7.00	372.00	-2.12	12.81	84791.62		
2.00	512.0	-2.12	27.19	36408.46		
5.00	254.0	0.88	30.81	62405.04		
5.00	459.0	0.88	74.19	317.20		
2.00	153.0	-2.12	31.81	463802.04		
1.00	426.0	-2.12	41.19	1850.577		
6.00	566.00	1.88	3.53	74.023		
6.00	596.00	1.88	3.53			
5.00	395.00	0.88	0.77			
6.00	676.00	1.88	3.53			
3.00	194.00	-1.12	1.25			
2.00	135.00	-2.12	4.49			
7.00	367.00	2.88	8.29			
4.12	384.81		86.65			
			1.86			
				136.21		
					88	

Ads x	Customer y	x-xmean	(x-xmean) ²	y-ymean	(y-ymean) ²	(x-xmean) (y-ymean)
5.00	353.00	0.88	0.77	-31.81	1011.88	-27.993
		
			3.53	-65.81	4330.96	-123.723
			4.49	-231.81	53735.88	491.437
			9.73	41.19	1696.62	-128.513
			3.53	181.19	32829.82	340.637
			3.53	211.19	44601.22	397.037
			0.77	10.19	103.84	8.967
6.00	676.00	1.88	3.53	291.19	84791.62	547.437
3.00	194.00	-1.12	1.25	-190.81	36408.46	213.707
2.00	135.00	-2.12	4.49	-249.81	62405.04	529.597
7.00	367.00	2.88	8.29	-17.81	317.20	-51.293
4.12	384.81		86.65		463802.04	1850.577

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

1.86

$$s_y = \sqrt{\frac{\sum(y_i - \bar{y})^2}{n - 1}}$$

136.21

74.023

$$s_{xy} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

Using the Estimated Regression –Equation for Estimation and Prediction

7. ESTIMATION

POINT ESTIMATION

INTERVAL ESTIMATION

CONFIDENCE INTERVAL FOR THE MEAN VALUE OF Y

PREDICTION INTERVAL FOR AN INDIVIDUAL VALUE OF Y

POINT ESTIMATION

If 3 TV ads are run prior to a sale, we expect the mean number of cars sold to be:

► $\hat{y} = 10 + 5(3) = 25 \text{ cars}$

Using the Estimated Regression Equation for Estimation and Prediction

- Confidence Interval Estimate of $E(y_p)$

$$\hat{y}_p \pm t_{\alpha/2} s_{\hat{y}_p}$$

- Prediction Interval Estimate of y_p

$$y_p \pm t_{\alpha/2} s_{\text{ind}}$$

where:

confidence coefficient is $1 - \alpha$ and
 $t_{\alpha/2}$ is based on a t distribution
with $n - 2$ degrees of freedom

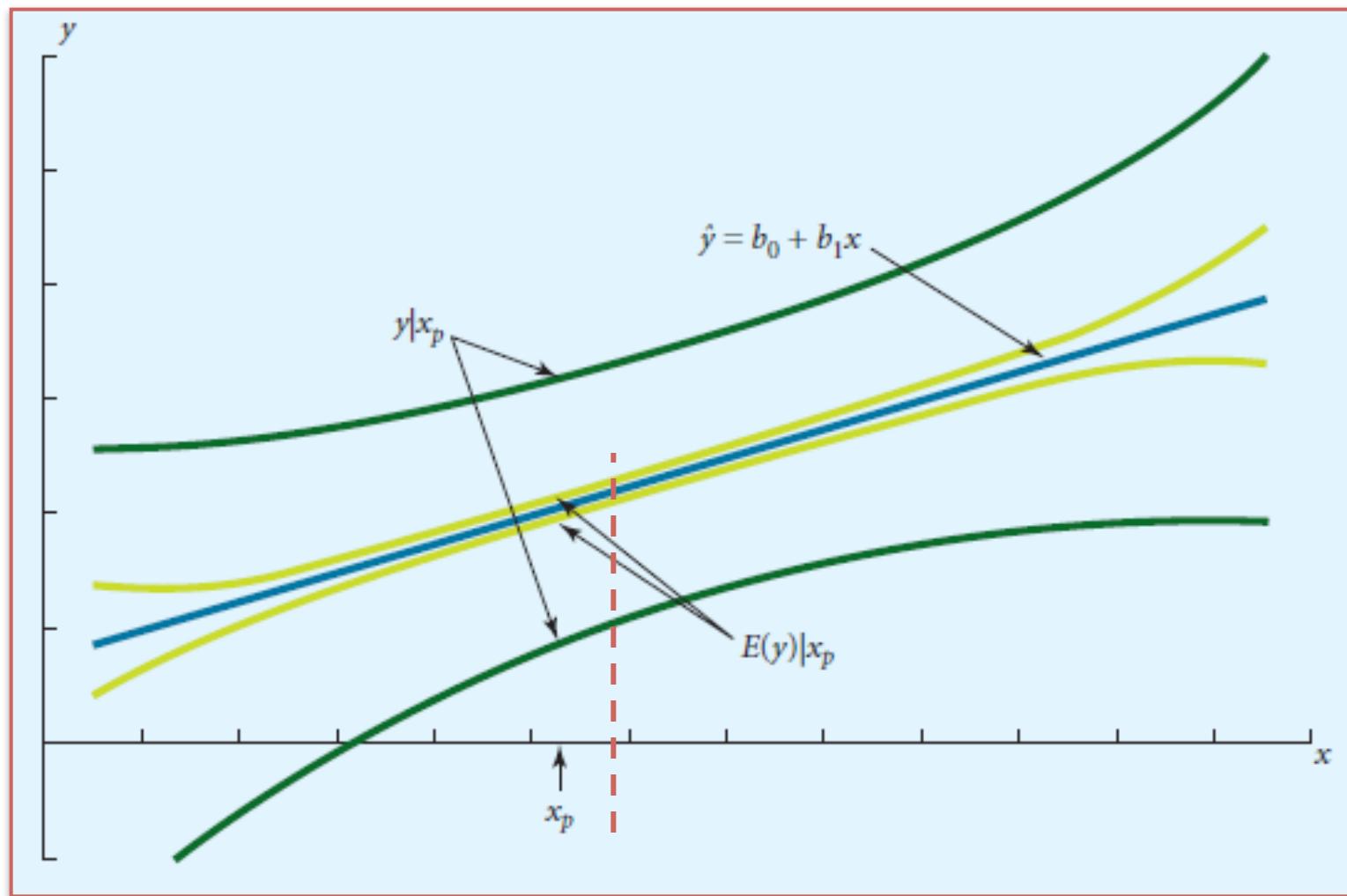
Confidence Interval for $E(y_p)$

- ▶ ■ Estimate of the Standard Deviation of \hat{y}_p

$$\hat{y}_p \pm t_{\alpha/2} s_{\hat{y}_p}$$

$$s_{\hat{y}_p} = s \sqrt{\frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

CONFIDENCE VS PREDICTION INTERVAL



Prediction Interval for y_p

- ▶ ■ Estimate of the Standard Deviation of an Individual Value of y_p

$$y_p \pm t_{\alpha/2} s_{\text{ind}}$$

$$s_{\text{ind}} = s \sqrt{1 + \frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$