

# Today's Agenda

---

Simple Regression

Multivariate Regression

## Question

---

Suppose you start up a company that has developed a drug that is supposed to increase IQ. You know that the standard deviation of IQ in the general population is 15. You test your drug on 36 patients and obtain a mean IQ of 97.65. Using an alpha value of 0.05, is this IQ significantly different than the population mean of 100?

# Question

---

$$z = \frac{97.65 - 100}{2.5} = -0.94$$

Level of Significance = 0.05, two tailed,  $0.05/2 = 0.025$ , Z value = -1., Since calculated value is less than tab null accepted.

# Question Normal Distribution

- A company's share price is normally distributed with a mean of Rs 800 and a standard deviation of Rs 300. A random sample of 16 days share price is taken. (a) What is the probability that the share price of the sample exceeds Rs 900?

$$\mu = 800, \bar{x} = 900, \text{S.E} = 300/\sqrt{16} = 75$$

$$P(\bar{x} > 900) = \frac{900 - 800}{75} = 1.33$$

$0.5 - 0.4082 = 0.0918$ , Hence there is 9.18% probability.



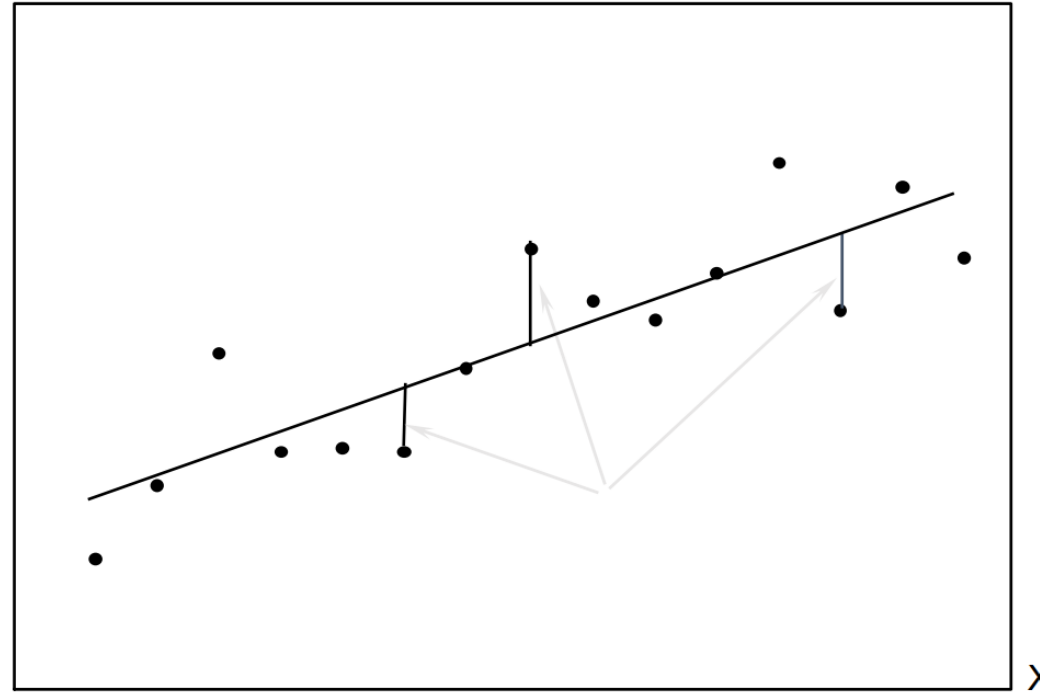
# Regression Analysis

---

The statistical technique that expresses the relationship between two or more variables in the form of an equation to estimate the value of a variable, based on the given value of another variable, is called regression analysis.

# Linear Regression model

Objective: The line that **BEST** fit the data.



*Minimize the sum of difference between actual and fitted values?*

*Or*

*Minimize the sum of squares of difference between the actual and fitted value?*

# Regression Analysis

---

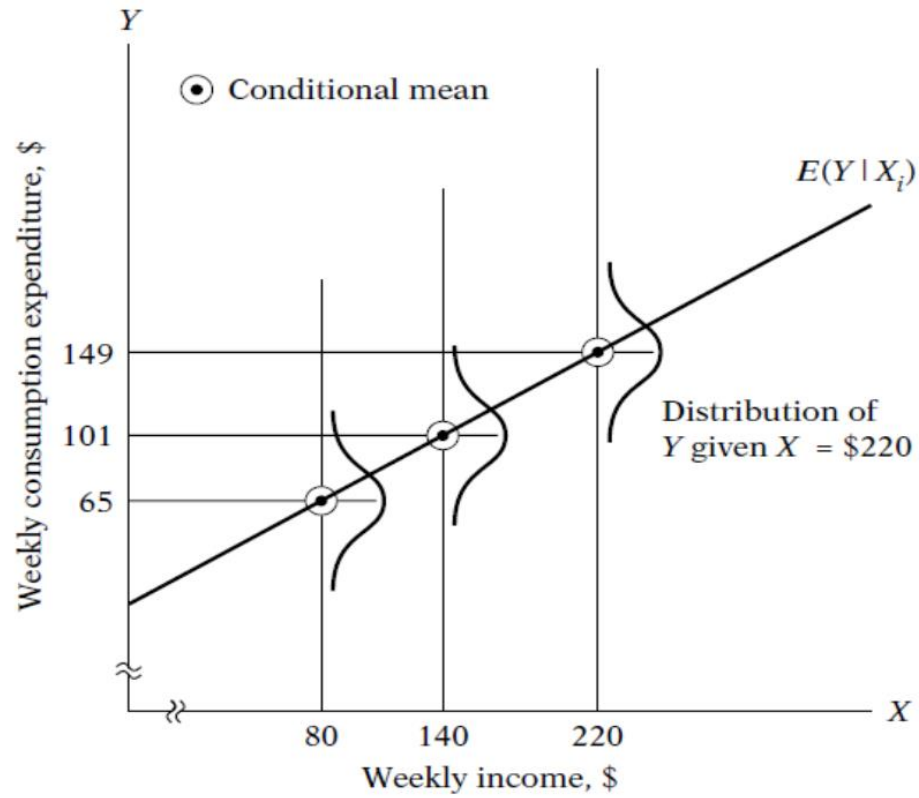
- **A time series is a set of observations  $x$ , each one being recorded at a specific time  $t$ .**
- Time series data, as the name suggests, are data that have been collected over a period of time on one or more variables.
- **Problems that can be solved with Time-Series:**
  - How the value of a country's stock index has varied with that country's macroeconomic fundamentals.
  - How the value of a company's stock price has varied when it announced the value of its dividend payment.
  - The effect on a country's exchange rate with increase in its trade deficit.
  - How interest rates are determined.
  - Finding out risk in an asset class.

# Regression Analysis

---

Dependent variable	Explanatory variable
↕	↕
Explained variable	Independent variable
↕	↕
Predictand	Predictor
↕	↕
<b>Regressand</b>	<b>Regressor</b>
↕	↕
Response	Stimulus
↕	↕
Endogenous	Exogenous
↕	↕
Outcome	Covariate
↕	↕
Controlled variable	Control variable

# Regression Analysis



# Regression Analysis

---

You are fitting a below mentioned straight -line equation.

$$y_i = \beta_1 + \beta_2 x_2$$

where  $\beta_1$  and  $\beta_2$  are unknown but fixed parameters known as the **regression coefficients**.

In regression analysis our interest is in estimating the PRFs.

# Regression Analysis

---

## Assumptions of Linear Regression Model

- 1. Linear regression model.** The regression model is **linear in the parameters**,  **$X$  values are fixed in repeated sampling**. Values taken by the regressor  $X$  are considered fixed in repeated samples. More technically,  $X$  is assumed to be *non stochastic*.
- 2. Homoscedasticity or equal variance of  $u_i$ .**
- 3. No autocorrelation between the disturbances.**
- 4. Zero covariance between  $u_i$  and  $X_i$ ,**
- 5. The number of observations  $n$  must be greater than the number of parameters to be estimated.**

# Regression Analysis

---

- **PRECISION OR STANDARD ERRORS**

- It is evident that least-squares estimates are a function of the sample data. But since the data are likely to change from sample to sample, the estimates will change ipso facto.
- In statistics the precision of an estimate is measured by its standard error.

**Standard error of estimate or the standard error of the regression (se).**

- *It is simply the standard deviation of the Y values about the estimated regression line and is often used as a summary measure of the “goodness of fit” of the estimated regression line.*

# Regression Analysis

---

## PROPERTIES OF LEAST-SQUARES ESTIMATORS

- It is **linear**, that is, a linear function of a random variable, such as the dependent variable  $Y$  in the regression model.
- It is **unbiased**, that is, its average or expected value,  $E(\hat{\beta}_2)$ , is equal to the true value.
- It has minimum variance in the class of all such linear unbiased estimators; an unbiased estimator with the least variance is known as an **efficient estimator**.

# Regression Analysis

---

## THE COEFFICIENT OF DETERMINATION- $R^2$

### A MEASURE OF “GOODNESS OF FIT”

We now consider the **goodness of fit** of the fitted regression line to a set of data; that is, we shall find out how “well” the sample regression line fits the data.

The **coefficient of determination**  $r^2$  (two-variable case) or  $R^2$  (multiple regression) is a summary measure that tells how well the sample regression line fits the data.

$$R^2 = \frac{ESS}{TSS}$$

Where **ESS** = Explained sum of square  
**TSS** = Total sum of square

# Regression Analysis

## Output Interpretation

Regression Statistics										
Multiple R	0.983559									
R Square	0.967389									
Adjusted R Square	0.964672									
Standard Error	19.00868									
Observations	14									
ANOVA										
	df	SS	MS	F	Significance F					
Regression	1	128625.3	128625.3	355.9772	2.75E-10					
Residual	12	4335.96	361.33							
Total	13	132961.2								
	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%		
Intercept	-22.5464	10.43676	-2.16029	0.051685	-45.2861	0.193368	-45.2861	0.193368		
X Variable 1	3.269721	0.1733	18.86736	2.75E-10	2.892132	3.64731	2.892132	3.64731		

# Regression Analysis

---

Null Hypothesis of Regression Co-efficient  $H_0 =$  All coefficients  $=0$

# Significance testing...

---

H0:  $\beta_1 = 0$  (no linear relationship)

H1:  $\beta_1 \neq 0$  (linear relationship does exist)

$$T_{n-2} = \frac{\hat{\beta} - 0}{s.e.(\hat{\beta})}$$

# Residual Analysis: check assumptions

---

$$e_i = Y_i - \hat{Y}_i$$

The residual for observation  $i$ ,  $e_i$ , is the difference between its observed and predicted value

Check the assumptions of regression by examining the residuals

- Examine for linearity assumption

- Examine for constant variance for all levels of  $X$  (homoscedasticity)

- Evaluate normal distribution assumption

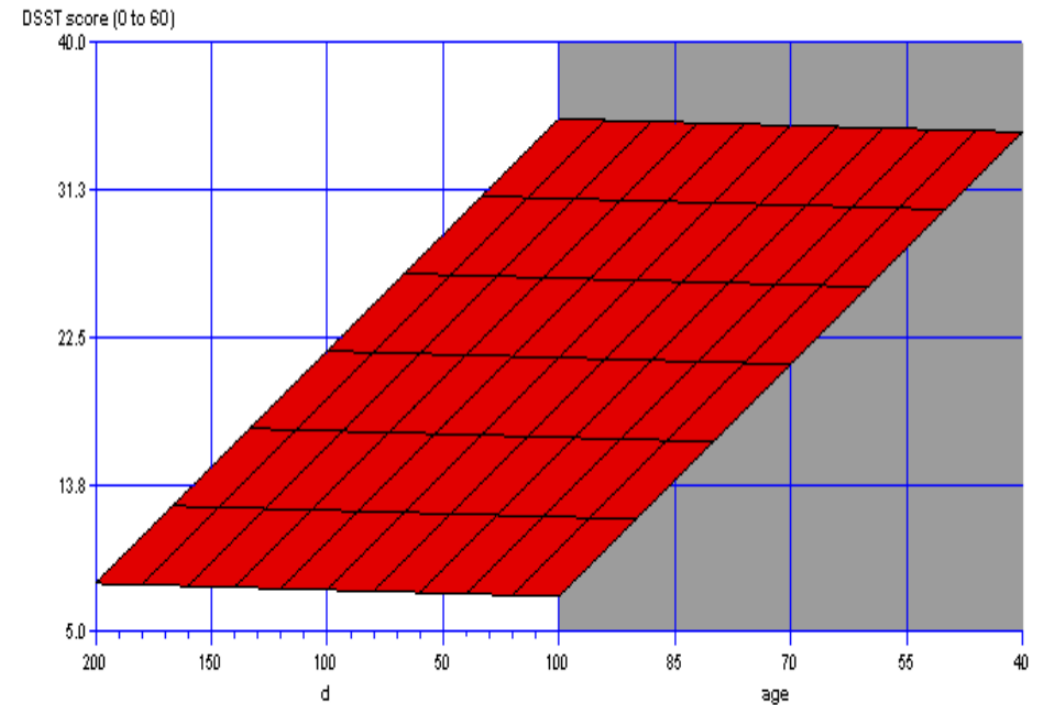
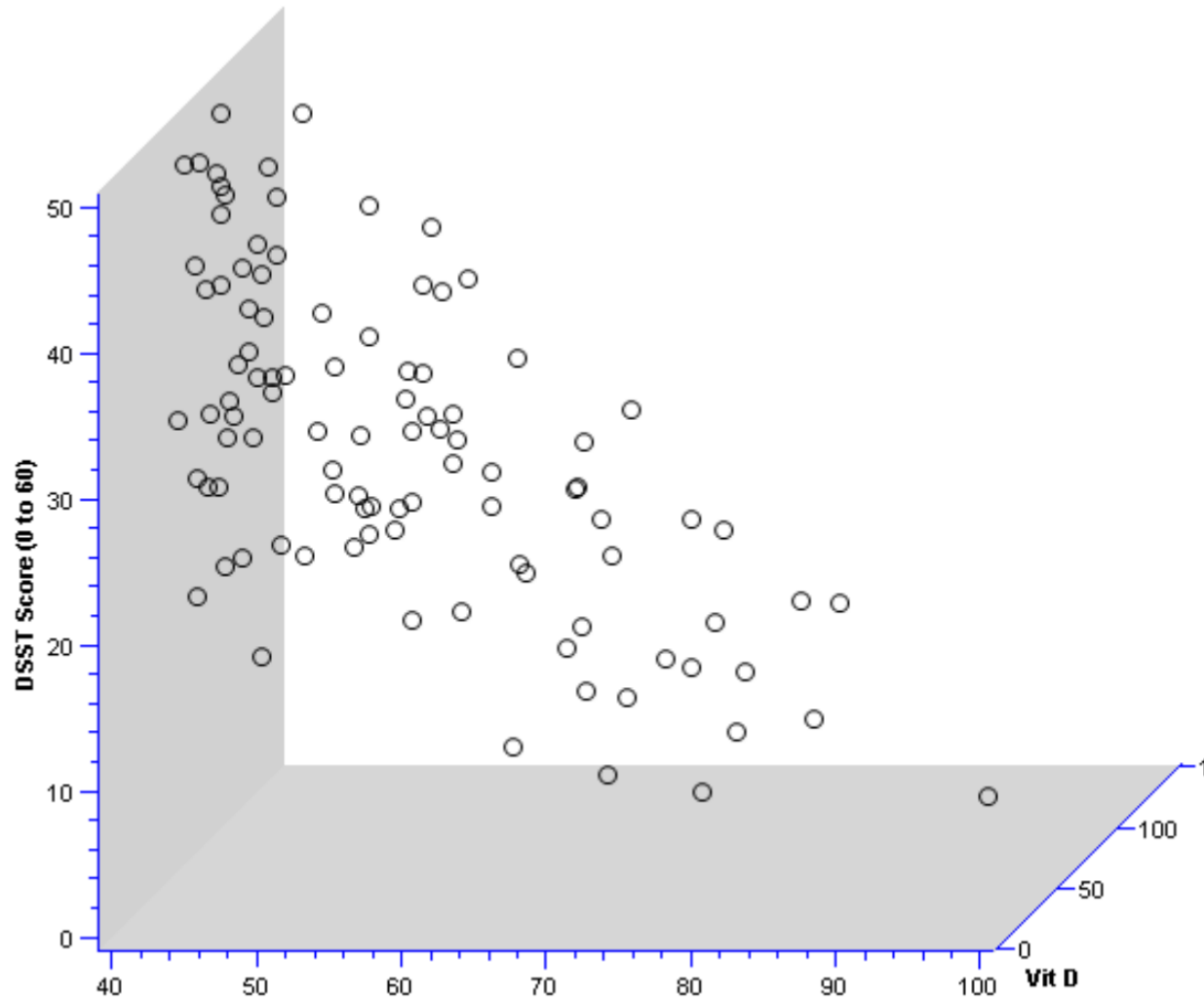
- Evaluate independence assumption

## Graphical Analysis of Residuals

- Can plot residuals vs.  $X$

# Multiple linear regression ( We fit a Plane)

---



# Functions of multivariate analysis:

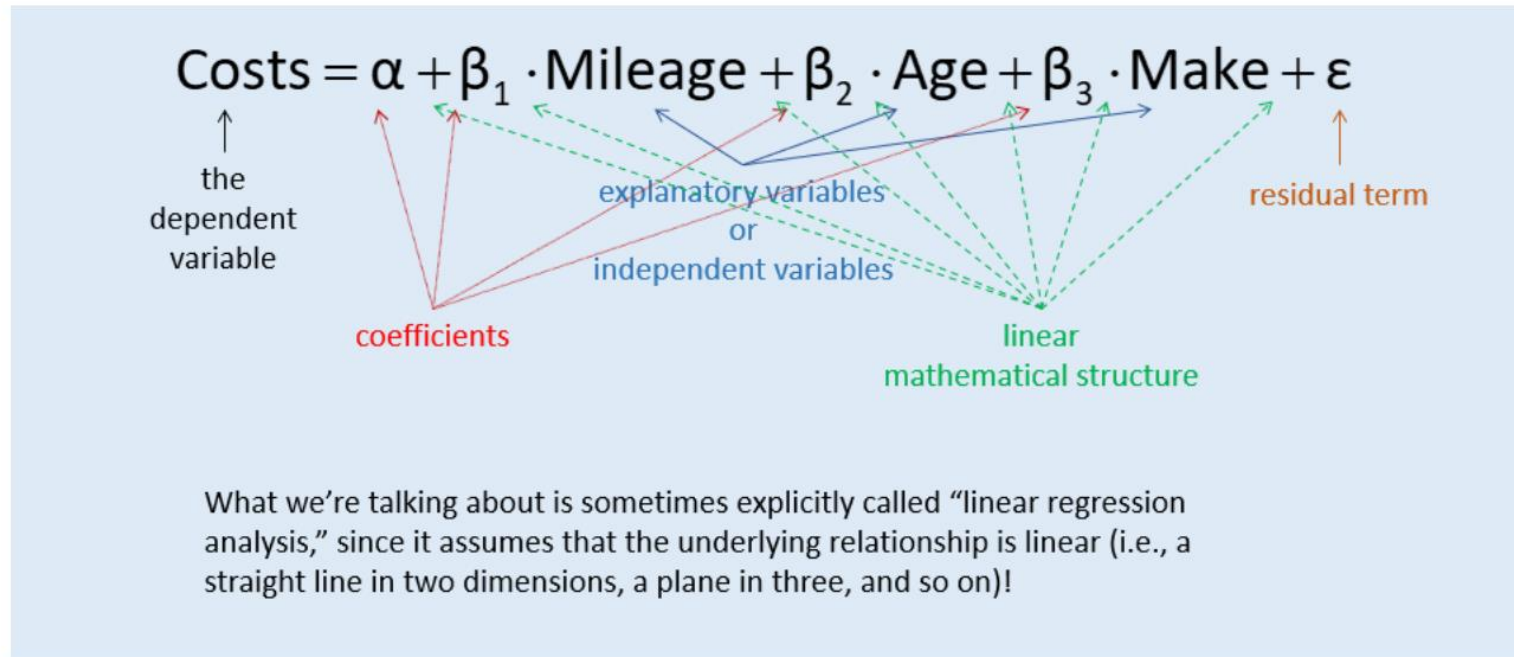
---

- Control for confounders
- Test for interactions between predictors (effect modification)
- Improve predictions

# The Regression Model

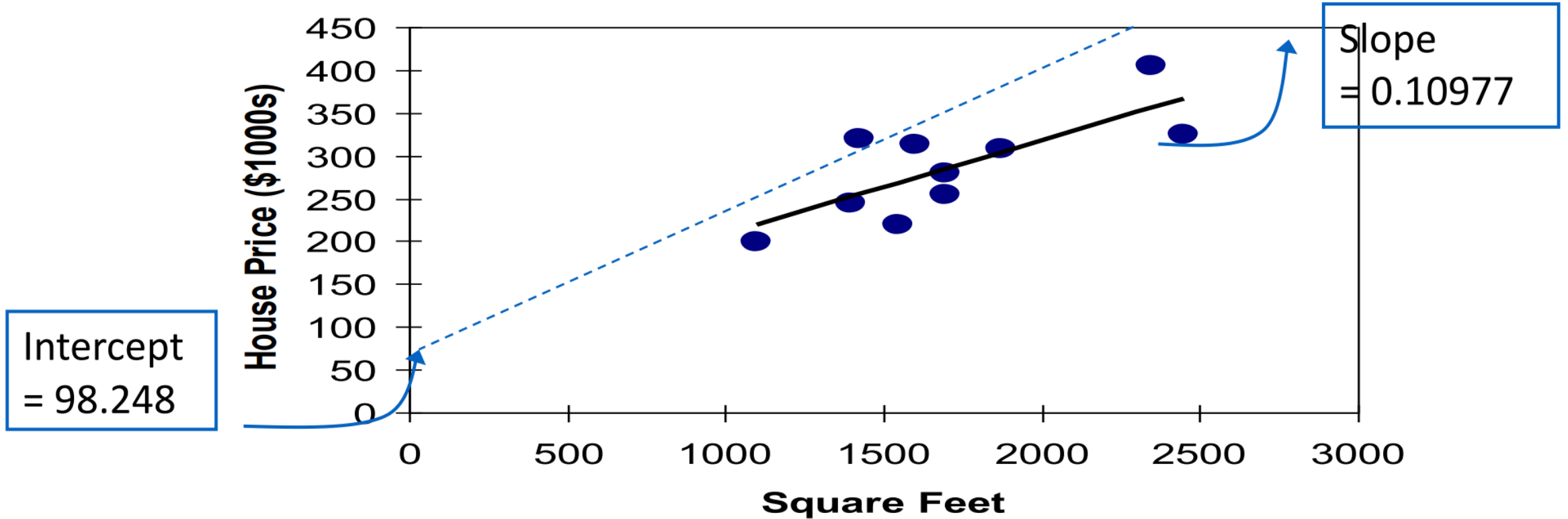
---

$$\text{Costs} = \alpha + \beta_1 \cdot \text{Mileage} + \beta_2 \cdot \text{Age} + \beta_3 \cdot \text{Make} + \varepsilon$$

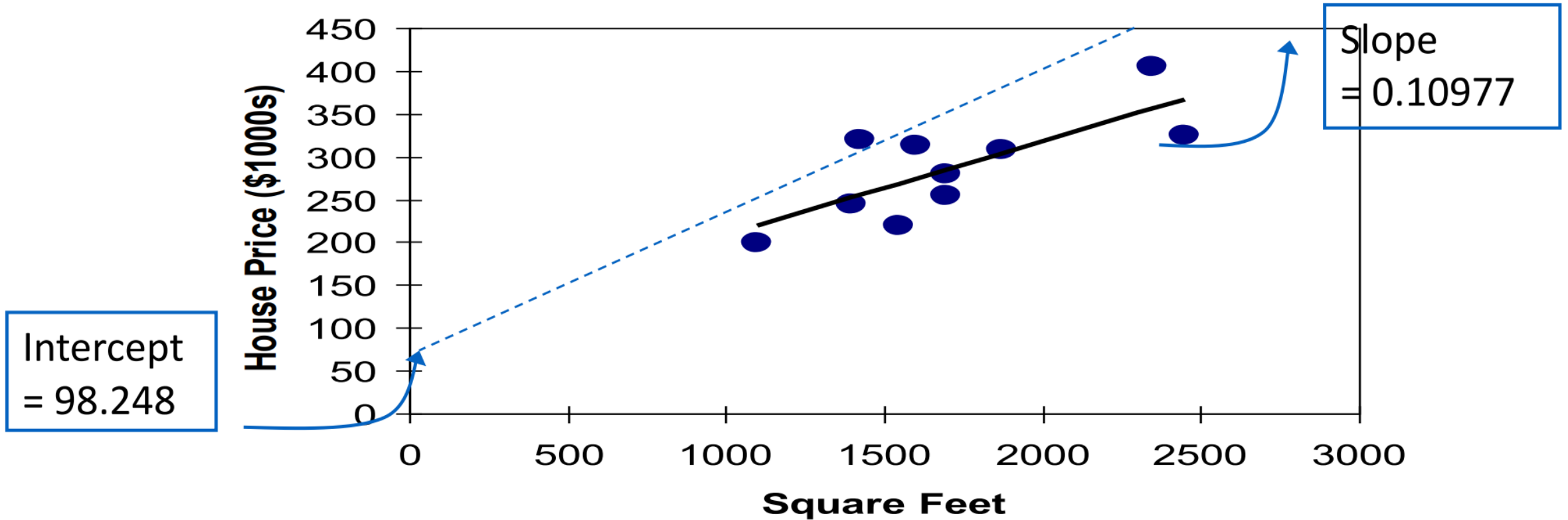


# The Regression Model

- House price model: scatter plot and regression line



- House price model: scatter plot and regression line



## Interpretation of the Intercept, $b_0$

---

$$\text{Price} = 98.24833 + 0.10977 (\text{sales})$$

$b_0$  is the estimated average value of Y when the value of X is zero (if  $X = 0$  is in the range of observed X values)

# Excel Output

## Regression Statistics

Multiple R	0.76211
R Square	0.58082
Adjusted R Square	0.52842
Standard Error	41.33032
Observations	10

$$r^2 = \frac{SSR}{SST} = \frac{18934.9348}{32600.5000} = 0.58082$$

58.08% of the variation in house prices is explained by variation in square feet

## ANOVA

	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	1	18934.9348	18934.9348	11.0848	0.01039
Residual	8	13665.5652	1708.1957		
Total	9	32600.5000			

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	98.24833	58.03348	1.69296	0.12892	-35.57720	232.07386
Square Feet	0.10977	0.03297	3.32938	0.01039	0.03374	0.18580

## Standard Error of Estimate

---

The standard deviation of the variation of observations around the regression line is estimated by

$$S_{YX} = \sqrt{\frac{SSE}{n-2}} = \sqrt{\frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n-2}}$$

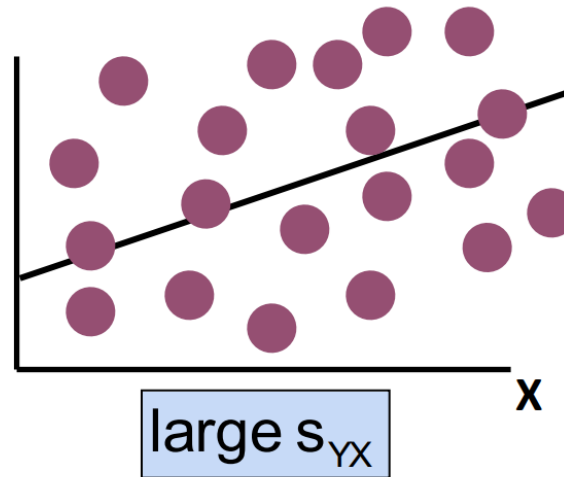
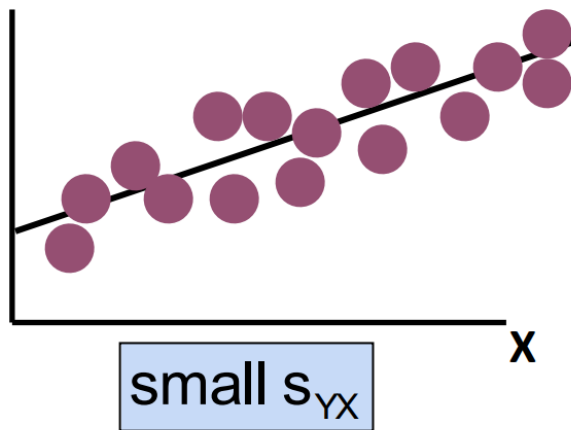
Where

SSE = error sum of squares

n = sample size

# Comparing Standard Errors

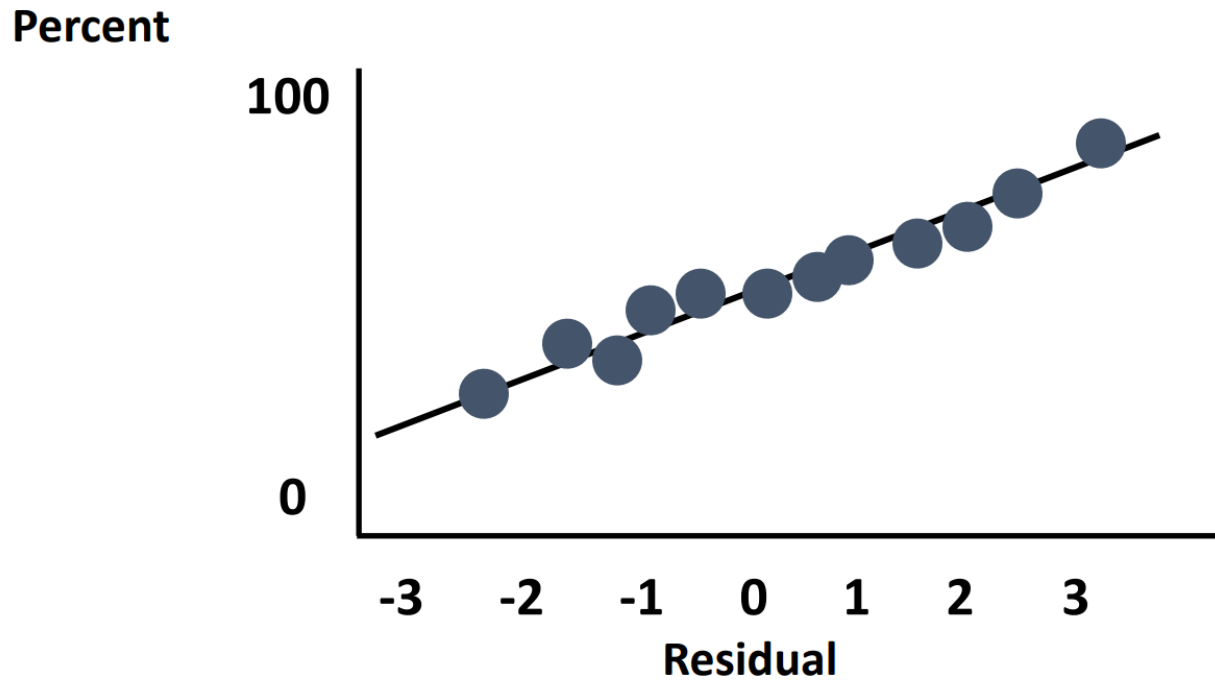
$S_{YX}$  is a measure of the variation of observed Y values from the regression line



The magnitude of  $S_{YX}$  should always be judged relative to the size of the Y values in the sample data

# Residual Analysis for Normality

A normal probability plot of the residuals can be used to check for normality:



# Regression Estimator

---

Regression estimator should be BLUE

B= Best

L= Linear

U= Unbiased

E= Estimator