

# INTRODUCTION TO DATA ANALYTICS

***Class # 16***

**Relation Analysis**

**Dr. Sreeja S R**

*Assistant Professor*

**Indian Institute of Information Technology**

**IIIT Sri City**

# THIS TOPIC INCLUDES...

- Regression Analysis
  - Simple Linear Regression
  - Multiple Linear Regression
  - Non-Linear Regression Analysis
- Auto-Regression Analysis

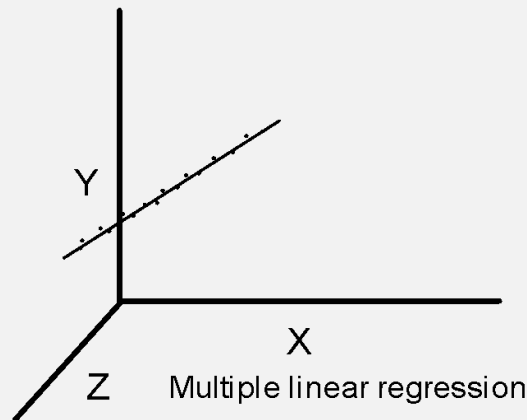
# REGRESSION ANALYSIS

# REGRESSION ANALYSIS

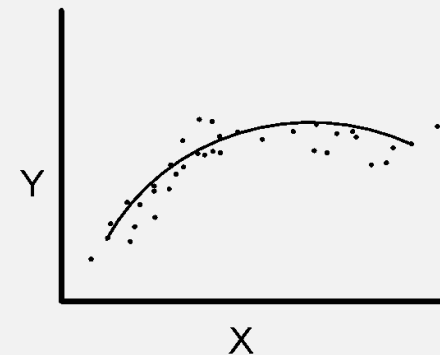
- The regression analysis is a statistical method to deal with the formulation of mathematical model depicting relationship amongst variables, which can be used for the purpose of prediction of the values of dependent variable, given the values of independent variables.
- **Classification of Regression Analysis Models**
  - Linear regression models
    1. Simple linear regression
    2. Multiple linear regression
  - Non-linear regression models



Simple linear regression



Multiple linear regression

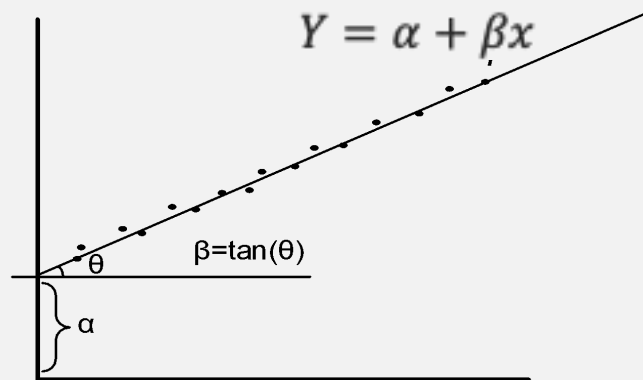


Non-linear regression

# SIMPLE LINEAR REGRESSION MODEL

In simple linear regression, we have only two variables:

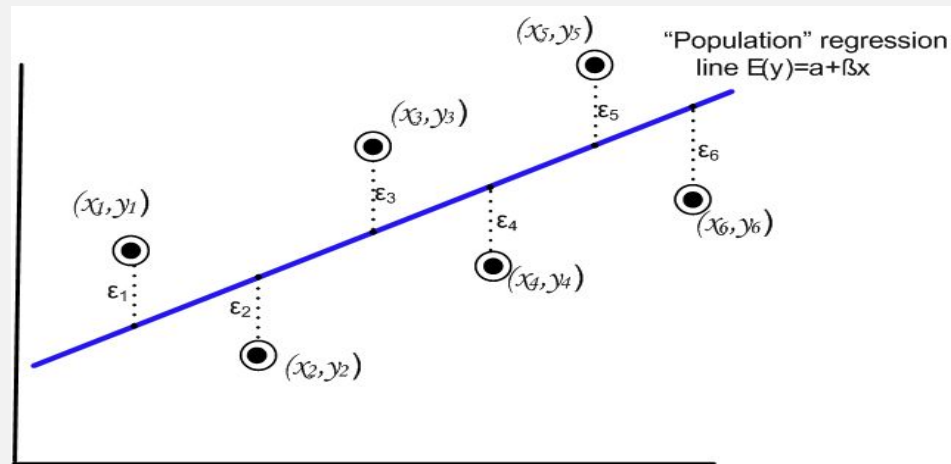
- Dependent variable (also called **Response**), usually denoted as  $Y$ .
- Independent variable (alternatively called **Regressor**), usually denoted as  $x$ .
- A reasonable form of a relationship between the Response  $Y$  and the Regressor  $x$  is the linear relationship, that is in the form  $Y = \alpha + \beta x$



## Note:

- There are infinite number of lines (and hence  $\alpha_s$  and  $\beta_s$ )
- The concept of regression analysis deal with finding the best relationship between  $Y$  and  $x$  (and hence best fitted values of  $\alpha$  and  $\beta$ ) quantifying the strength of that relationship.

# REGRESSION ANALYSIS



Given the set  $[(x_i, y_i), i = 1, 2, \dots, n]$  of data involving  $n$  pairs of  $(x, y)$  values, our objective is to find “true” or population regression line such that  $Y = \alpha + \beta x + \epsilon$

Here,  $\epsilon$  is a random variable with  $E(\epsilon) = 0$  and  $var(\epsilon) = \sigma^2$ . The quantity  $\sigma^2$  is often called the **error variance**.

## Note:

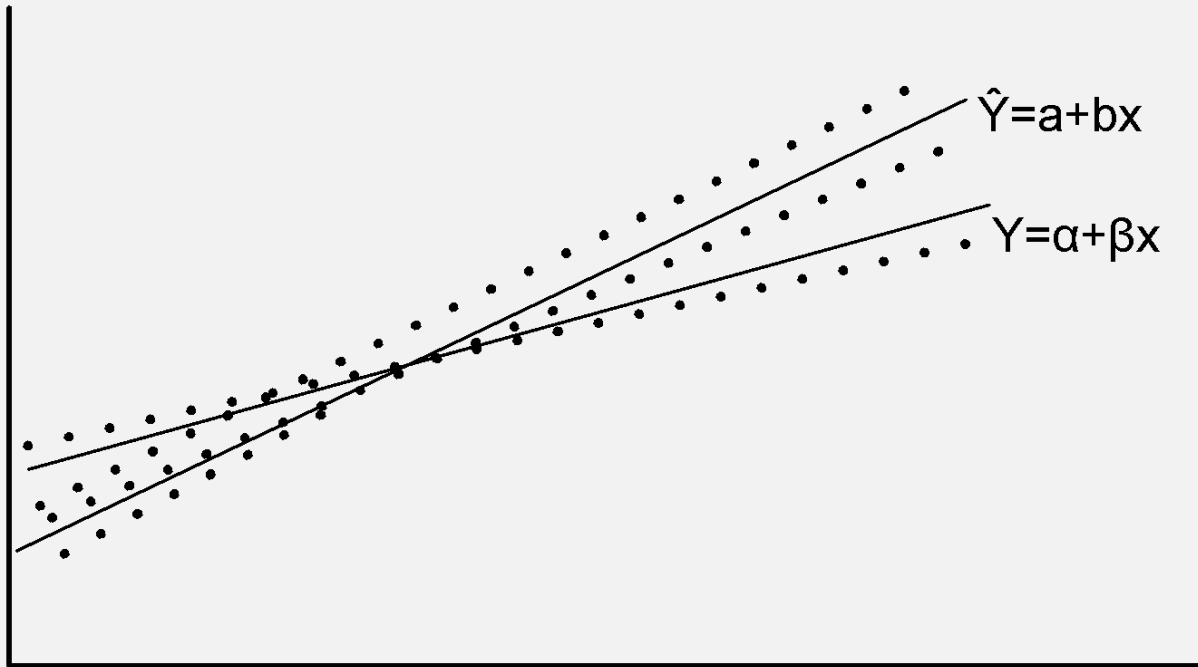
- $E(\epsilon) = 0$  implies that at a specific  $x$ , the  $y$  values are distributed around the “true” regression line  $Y = \alpha + \beta x$  (i.e., the positive and negative errors around the true line is reasonable).
- $\alpha$  and  $\beta$  are called **regression coefficients**.
- $\alpha$  and  $\beta$  values are to be estimated from the data.

# TRUE VERSUS FITTED REGRESSION LINE

- The task in regression analysis is to estimate the regression coefficients  $\alpha$  and  $\beta$ .
- Suppose, we denote the estimates  $a$  for  $\alpha$  and  $b$  for  $\beta$ . Then the fitted regression line is

$$\hat{Y} = a + bx$$

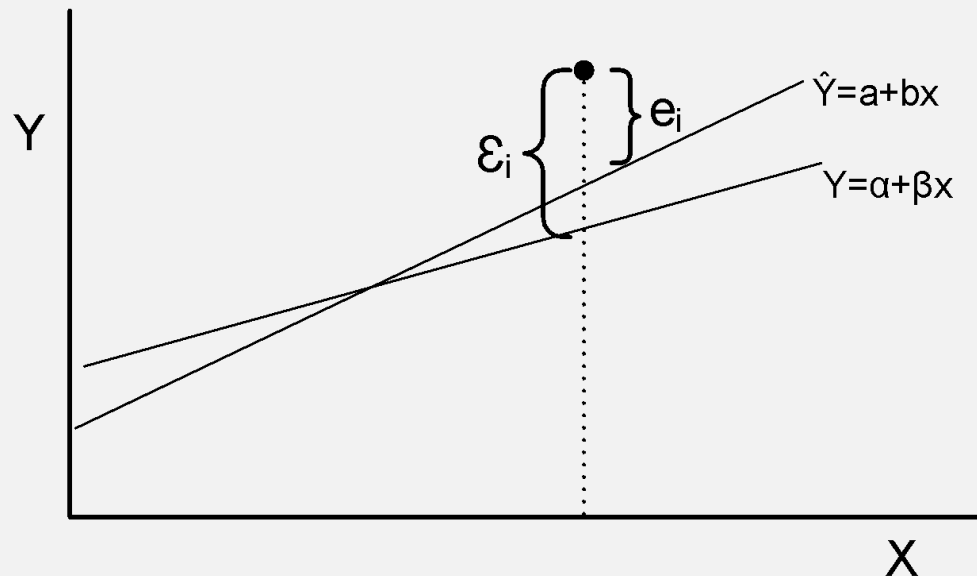
where  $\hat{Y}$  is the predicted or fitted value.



# LEAST SQUARE METHOD TO ESTIMATE $\alpha$ AND $\beta$

This method uses the concept of **residual**. A residual is essentially an error in the fit of the model  $\hat{Y} = a + bx$ . Thus,  $i^{th}$  residual is

$$e_i = Y_i - \hat{Y}_i, i = 1, 2, 3, \dots, n$$





# LEAST SQUARE METHOD

- The **residual sum of squares** is often called **the sum of squares of the errors** about the fitted line and is denoted as SSE

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

- We are to minimize the value of SSE and hence to determine the parameters of  $a$  and  $b$ .
- Differentiating SSE with respect to  $a$  and  $b$ , we have

$$\frac{\partial(SSE)}{\partial a} = -2 \sum_{i=1}^n (y_i - a - bx_i)$$

$$\frac{\partial(SSE)}{\partial b} = -2 \sum_{i=1}^n (y_i - a - bx_i) \cdot x_i$$

For minimum value of SSE,  $\frac{\partial(SSE)}{\partial a} = 0$

$$\frac{\partial(SSE)}{\partial b} = 0$$

# LEAST SQUARE METHOD TO ESTIMATE $\alpha$ AND $\beta$

Thus we set

$$na + b \sum_{i=1}^n x_i = \sum_{i=1}^n y_i$$

$$a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i$$

These two equations can be solved to determine the values of  $a$  and  $b$ , and it can be calculated that

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$a = \bar{y} - b\bar{x}$$

# $R^2$ : MEASURE OF QUALITY OF FIT

- A quantity  $R^2$ , is called **coefficient of determination** is used to measure the proportion of variability of the fitted model.
- We have  $SSE = \sum_{i=1}^n (y_i - \hat{y})^2$
- It signifies the **variability due to error**.
- Now, let us define the **total corrected sum of squares**, defined as

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

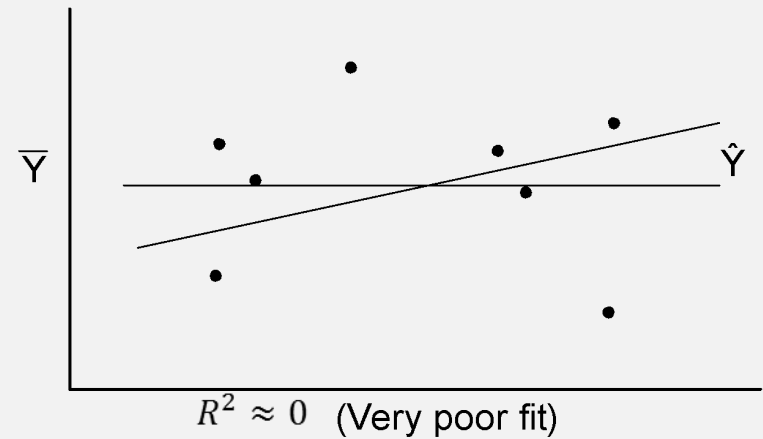
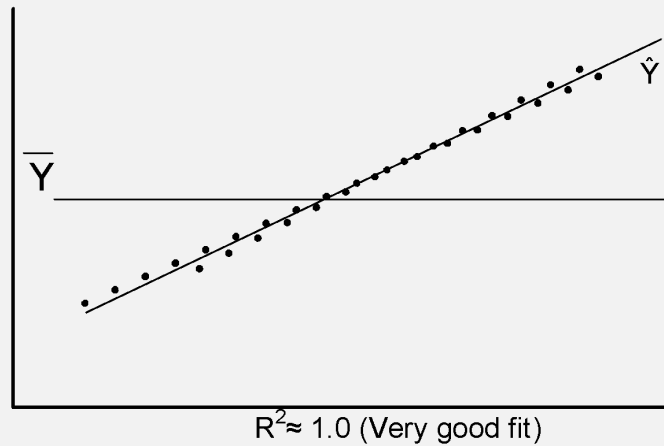
- SST represents the variation in the response values. The  $R^2$  is

$$R^2 = 1 - \frac{SSE}{SST}$$

## Note:

- If fit is perfect, all residuals are zero and thus  $R^2 = 1.0$  (very good fit)
- If SSE is only slightly smaller than SST, then  $R^2 \approx 0$  (very poor fit)

# $R^2$ : MEASURE OF QUALITY OF FIT



# MULTIPLE LINEAR REGRESSION

- When more than one variable are independent variable, then the regression can be estimated as a **multiple regression model**
- When this model is linear in coefficients, it is called **multiple linear regression model**
- If  $k$ -independent variables  $x_1, x_2, x_3, \dots, x_k$  are associated, the multiple linear regression model is given by

$$y_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_k x_k + \epsilon$$

- And the estimated response is obtained as

$$\hat{y}_i = b_0 + b_1 x_1 + b_2 x_2 + b_k x_k$$

# MULTIPLE LINEAR REGRESSION

## Estimating the coefficients

Let the data points given to us is

$$(x_{1i}, x_{2i}, x_{3i}, \dots, x_{ki}, y_i) \quad i = 1, 2, \dots, n, \quad n > k$$

where  $y_i$  is the observed response to the values  $x_{1i}, x_{2i}, x_{3i}, \dots, x_{ki}$  of  $k$  independent variables  $x_1, x_2, x_3, \dots, x_k$ .

Thus,

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_k x_{ki} + \epsilon_i$$

and 
$$\hat{y}_i = b_0 + b_1 x_{1i} + b_2 x_{2i} + b_k x_{ki} + e_i$$

where  $\epsilon_i$  and  $e_i$  are the random error and residual error, respectively associated with true response  $y_i$  and fitted response  $\hat{y}_i$ .

Using the concept of **Least Square Method** to estimate  $b_0, b_1, b_2, \dots, b_k$ , we minimize the expression

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

# MULTIPLE LINEAR REGRESSION

- Differentiating SSE in turn with respect to  $b_0, b_1, b_2, \dots, b_k$  and equating to zero, we generate the set of  $(k+1)$  normal **estimation equations for multiple linear regression**.

$$nb_0 + b_1 \sum_{i=1}^n x_{1i} + b_2 \sum_{i=1}^n x_{2i} + \dots + b_k \sum_{i=1}^n x_{ki} = \sum_{i=1}^n y_i$$

$$b_0 \sum_{i=1}^n x_{1i} + b_1 \sum_{i=1}^n x_{1i}^2 + b_2 \sum_{i=1}^n x_{1i} \cdot x_{2i} + \dots + b_k \sum_{i=1}^n x_{1i} \cdot x_{ki} = \sum_{i=1}^n x_{1i} \cdot y_i$$

$$\dots \quad \dots \quad \dots \quad \dots \quad \dots \quad \dots$$

$$\dots \quad \dots \quad \dots \quad \dots \quad \dots \quad \dots$$

$$b_0 \sum_{i=1}^n x_{ki} + b_1 \sum_{i=1}^n x_{ki} \cdot x_{1i} + b_2 \sum_{i=1}^n x_{ki} \cdot x_{2i} + \dots + b_k \sum_{i=1}^n x_{ki}^2 = \sum_{i=1}^n x_{ki} \cdot y_i$$

- The system of linear equations can be solved for  $b_0, b_1, \dots, b_k$  by any appropriate method for solving system of linear equations.
- Hence, the multiple linear regression model can be built.

# NON LINEAR REGRESSION MODEL

- When the regression equation is in terms of  $r$ -degree,  $r > 1$ , then it is called nonlinear regression model. When more than one independent variables are there, then it is called Multiple Non linear Regression model. Also, alternatively termed as polynomial regression model. In general, it takes the form

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_r x^r + \epsilon$$

- The estimated response is obtained as

$$\hat{y} = b_0 + b_1 x + b_2 x^2 + \dots + b_r x^r$$



## SOLVING FOR POLYNOMIAL REGRESSION MODEL

Given that  $(x_i, y_i); i = 1, 2, \dots, n$  are  $n$  pairs of observations. Each observations would satisfy the equations:

$$y_i = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_r x^r + \epsilon_i$$

and 
$$\hat{y}_i = b_0 + b_1 x + b_2 x^2 + \dots + b_r x^r + e_i$$

where,  $r$  is the degree of polynomial

$\epsilon_i$  = is the  $i^{th}$  random error

$e_i$  = is the  $i^{th}$  residual error

**Note:** The number of observations,  $n$ , must be at least as large as  $r+1$ , the number of parameters to be estimated.

The polynomial model can be transformed into a general linear regression model setting  $x_1 = x, x_2 = x^2, \dots, x_n = x^r$ . Thus, the equation assumes the form:

$$y_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_r x^r + \epsilon_i$$

$$\hat{y}_i = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_r x_r + e_i$$

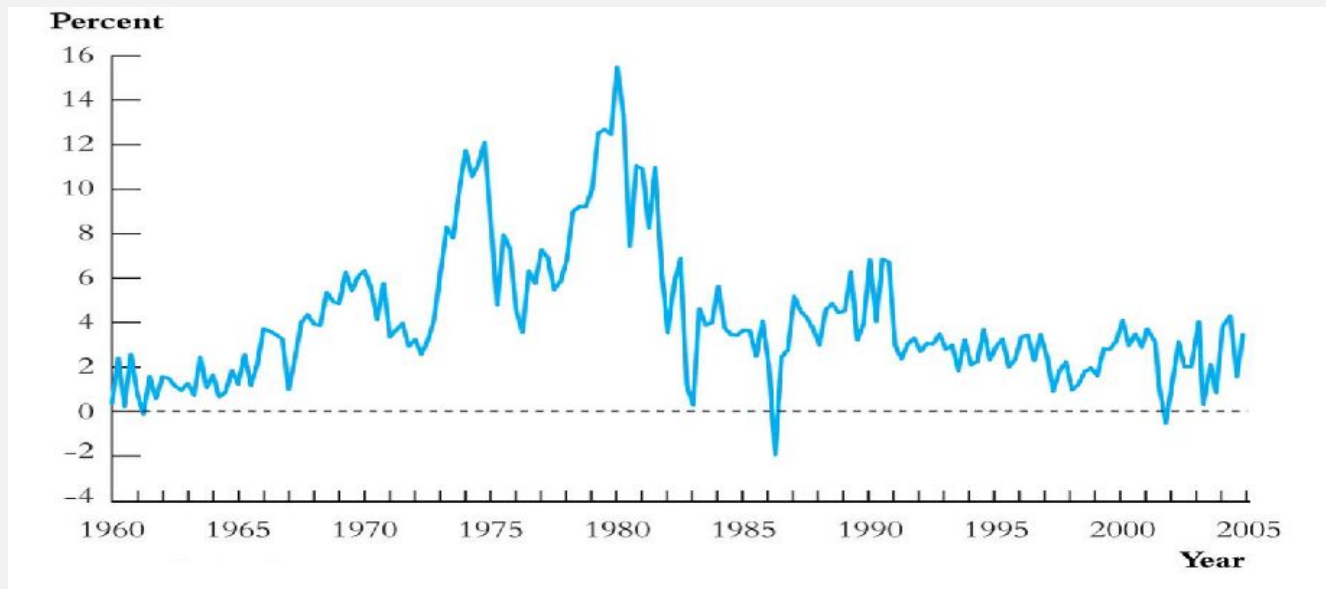
This model then can be solved using the procedure followed for multiple linear regression model.

# AUTO-REGRESSION ANALYSIS

# AUTO REGRESSION ANALYSIS

- Regression analysis for time-ordered data is known as **Auto-Regression Analysis**
- **Time series data** are data collected on the same observational unit at multiple time periods

Example: Indian rate of price inflation

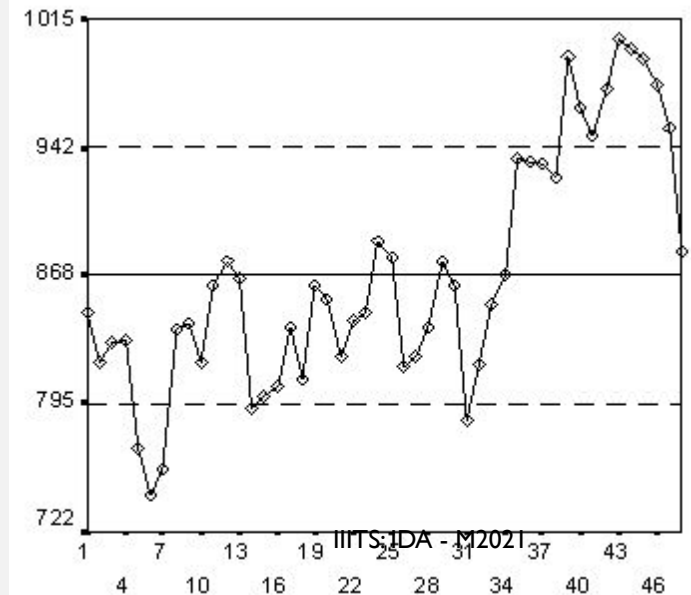
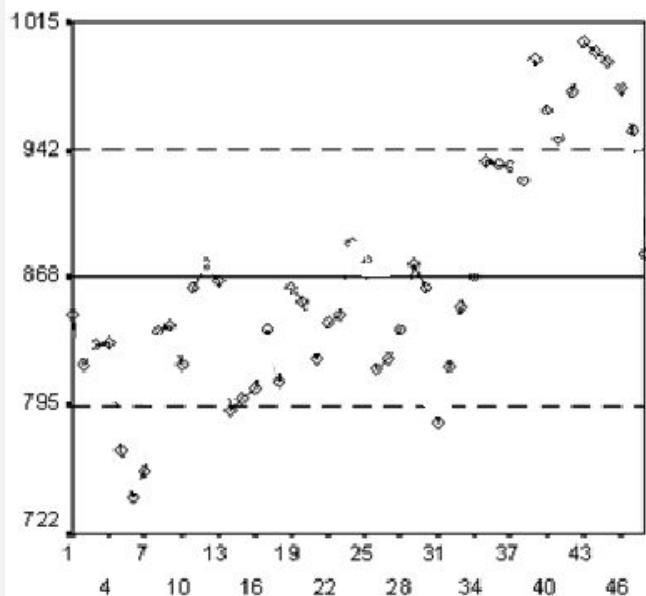
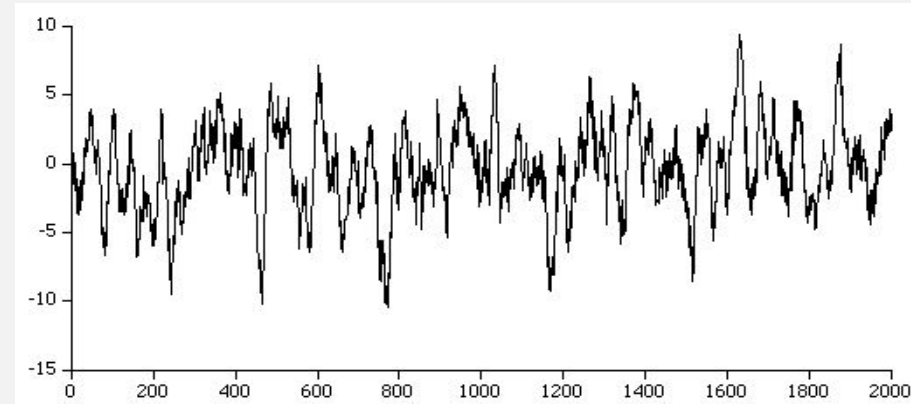
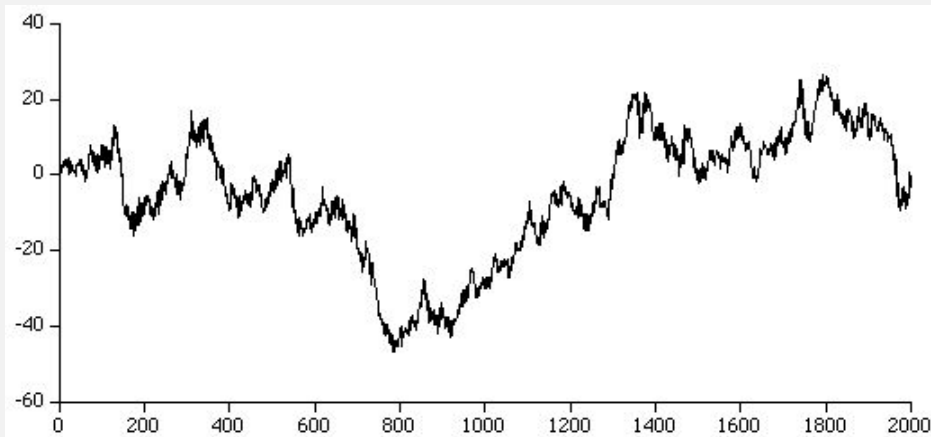


# AUTO REGRESSION ANALYSIS

- **Examples:** Which of the following is a time-series data?
  - Aggregate consumption and GDP for a country (for example, 20 years of quarterly observations = 80 observations)
  - Yen/\$, pound/\$ and Euro/\$ exchange rates (daily data for 1 year = 365 observations)
  - Cigarette consumption per capita in a state, by years
  - Rainfall data over a year
  - Sales of tea from a tea shop in a season

# AUTO REGRESSION ANALYSIS

- **Examples:** Which of the following graph is due to time-series data?



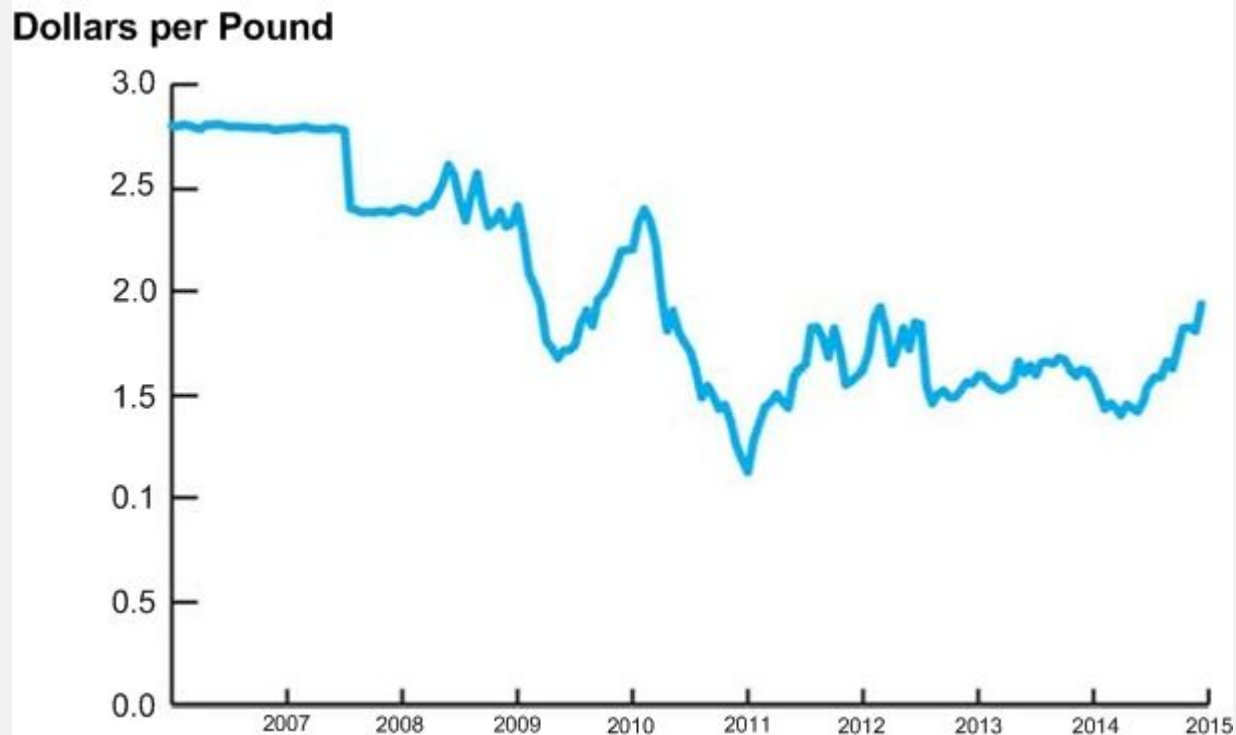
# USE OF TIME SERIES DATA

- To develop forecast model
  - What will the rate of inflation be next year?
- To estimate dynamic causal effects
  - If the rate of interest increases the interest rate now, what will be the effect on the rates of inflation and unemployment in 3 months? in 12 months?
  - What is the effect over time on electronics good consumption of a hike in the excise duty?
- Time dependent analysis
  - Rates of inflation and unemployment in the country can be observed only over time!

# MODELING WITH TIME SERIES DATA

- Correlation over time
  - Serial correlation, also called autocorrelation
  - Calculating standard error
- To estimate dynamic causal effects
  - Under which dynamic effects can be estimated?
  - How to estimate?
- Forecasting model
  - Forecasting model build on regression model

# AUTO-REGRESSION MODEL FOR FORECASTING



- Can we predict the trend at a time say 2022?



# SOME NOTATIONS AND CONCEPTS

- $Y_t$  = Value of  $Y$  in a period  $t$
- Data set  $[Y_1, Y_2, \dots, Y_{T-1}, Y_T]$ :  $T$  observations on the time series random variable  $Y$
- **Assumptions**
  - We consider only consecutive, evenly spaced observations
    - For example, monthly, 2000-2015, no missing months
  - A time series  $Y_t$  is **stationary** if its probability distribution does not change over time, that is, if the joint distribution of  $(Y_{i+1}, Y_{i+2}, \dots, Y_{i+T})$  does not depend on  $i$ .
    - Stationary property implies that history is relevant. In other words, Stationary requires the future to be like the past (in a probabilistic sense).
    - Auto Regression analysis assumes that  $Y_t$  is stationary.

# SOME NOTATIONS AND CONCEPTS

- There are four ways to have the time series data for AutoRegression analysis
  - **Lag:** The first lag of  $Y_t$  is  $Y_{t-1}$ , its  $j$ -th lag is  $Y_{t-j}$
  - **Difference:** The first difference of a series,  $Y_t$  is its change between period  $t$  and  $t-1$ , that is,  $y_t = Y_t - Y_{t-1}$
  - **Log difference:**  $y_t = \log(Y_t) - \log(Y_{t-1})$
  - **Percentage:**  $y_t = \frac{Y_{t-1}}{Y_t} \times 100$

# SOME NOTATIONS AND CONCEPTS

- **Autocorrelation**

- The correlation of a series with its own lagged values is called autocorrelation (also called serial correlation)

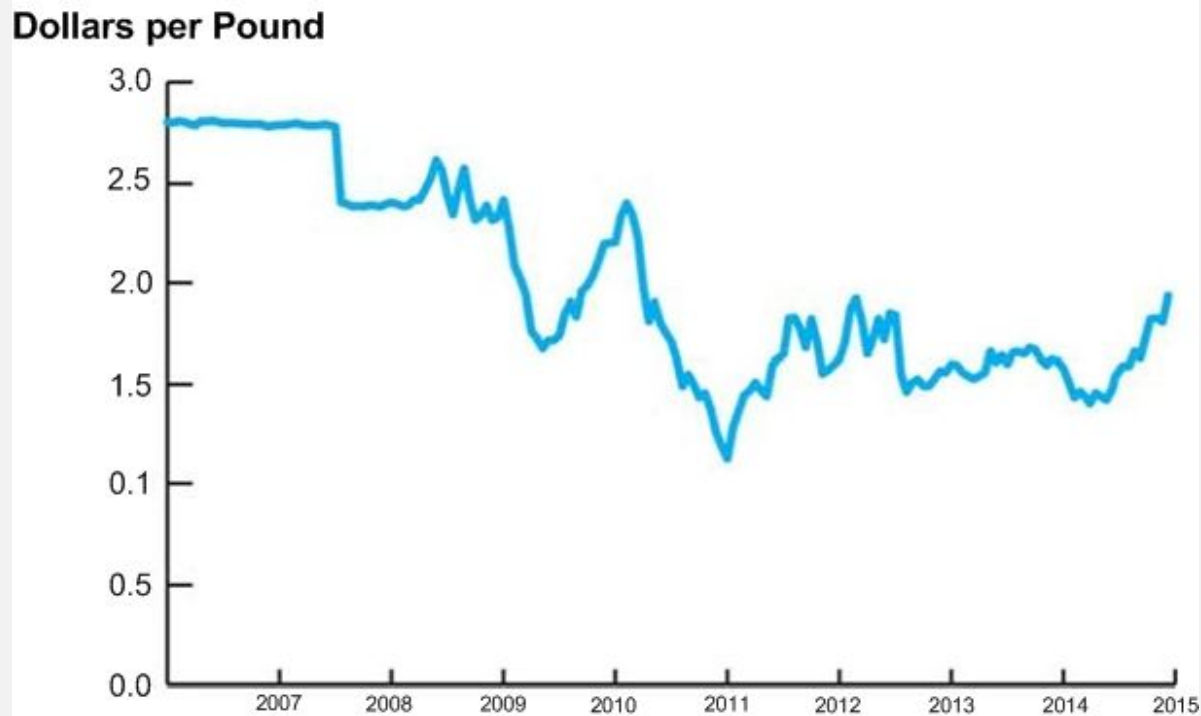
## Definition: ***j*-th Autocorrelation**

The *j*-th autocorrelation, denoted by  $\rho_j$  is defined as

$$\rho_j = \frac{COV(Y_t, Y_{t-j})}{\sqrt{\sigma_{Y_t} \sigma_{Y_{t-j}}}}$$

where,  $COV(Y_t, Y_{t-j})$  is the ***j*-th autocovariance**

# SOME NOTATIONS AND CONCEPTS



- For the given data, say  $\rho_1 = 0.84$ 
  - This implies that the Dollars per Pound is highly serially correlated
- Similarly, we can determine  $\rho_2, \rho_3, \dots$  etc., and hence different regression analyses

# AUTO-REGRESSION MODEL FOR FORECASTING

- A natural starting point for forecasting model is to use past values of  $Y$ , that is,  $Y_{t-1}$ ,  $Y_{t-2}$ , ... to predict  $Y_t$
- An autoregression is a regression model in which  $Y_t$  is regressed against its own lagged values.
- The number of lags used as regressors is called the **order of autoregression**
  - In first order autoregression (denoted as AR(1)),  $Y_t$  is regressed against  $Y_{t-1}$
  - In  $p$ -th order autoregression (denoted as AR( $p$ )),  $Y_t$  is regressed against,  $Y_{t-1}$ ,  $Y_{t-2}$ , ...,  $Y_{t-p}$

# P-TH ORDER AUTO-REGRESSION MODEL

## Definition: *p*-th AutoRegression Model

In general, the *p*-th order autoregression model is defined as

$$Y_t = \beta_0 + \sum_{i=1}^p \beta_i Y_{t-i} + \varepsilon_t$$

where,  $\beta_0, \beta_1, \dots, \beta_p$  is called autoregression coefficients and  $\varepsilon_t$  is the noise term or residue and in practice it is assumed to Gaussian white noise

- For example, AR(1) is  $Y_t = \beta_0 + \beta_1 Y_{t-1} + \varepsilon_t$
- The task in AR analysis is to derive the "best" values for  $\beta_i$   $i = 0, 1, \dots, p$  given a time series  $Y_t$ .

# COMPUTING AR COEFFICIENTS

- A number of techniques known for computing the AR coefficients
- The most common method is called **Least Squares Method** (LSM)
- The LSM is based upon the **Yule-Walker equations**

$$\begin{bmatrix}
 1 & r_1 & r_2 & r_3 & r_4 & \dots & r_{p-2} & r_{p-1} \\
 r_1 & 1 & r_1 & r_2 & r_3 & \dots & r_{p-3} & r_{p-2} \\
 r_2 & r_1 & 1 & r_1 & r_2 & \dots & r_{p-4} & r_{p-3} \\
 r_3 & r_2 & r_1 & 1 & r_2 & \dots & r_{p-5} & r_{p-4} \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\
 r_{p-1} & r_{p-2} & r_{p-3} & r_{p-4} & r_{p-5} & \dots & r_1 & 1
 \end{bmatrix}
 \begin{bmatrix}
 \beta_1 \\
 \beta_2 \\
 \beta_3 \\
 \vdots \\
 \vdots \\
 \vdots \\
 \vdots \\
 \beta_{p-1} \\
 \beta_p
 \end{bmatrix}
 =
 \begin{bmatrix}
 r_1 \\
 r_2 \\
 r_3 \\
 \vdots \\
 \vdots \\
 \vdots \\
 \vdots \\
 r_{p-1} \\
 r_p
 \end{bmatrix}$$

- Here,  $r_i$  ( $i = 1, 2, 3, \dots, p-1$ ) denotes the  $i$ -th auto correlation coefficient.
- $\beta_0$  can be chosen empirically, usually taken as zero.

# REFERENCE

- The detail material related to this lecture can be found in

The Elements of Statistical Learning, Data Mining, Inference, and Prediction (2<sup>nd</sup> Edn.), Trevor Hastie, Robert Tibshirani, Jerome Friedman, Springer, 2014.



Any question?