



**BHARATHIDASAN UNIVERSITY**

**Tiruchirappalli- 620024**

**Tamil Nadu, India.**

**Programme: M.Sc. Statistics**

**Course Title : Econometrics**

**Course Code: 23ST03DEC**

**Unit-III**

**Econometrics**

**Dr. T. Jai Sankar**

**Associate Professor and Head**

**Department of Statistics**

**Ms. E. Devi**

**Guest Faculty**

**Department of Statistics**

### UNIT III

Autocorrelation: Causes, consequences and testing for auto-correlated disturbances – Autoregressive series of order 1 (AR(1)) – Lagged variables and distributed lag methods – Errors in variable models and Instrumental variables. Economical Forecasting – long term and short term.

#### INTRODUCTION:

In econometrics, particularly in the study of time series data, the concept of autocorrelation occupies a central role in understanding the reliability of regression models. Autocorrelation, also known as serial correlation, refers to the correlation of a variable with its own past values or, more specifically in regression analysis, the correlation between the error terms (residuals) of different time periods. According to the classical linear regression model (CLRM), one key assumption is that the disturbance terms are independently distributed and uncorrelated across observations. When this assumption is violated, and the error term in one period is systematically related to that of another, autocorrelation is said to exist. The presence of autocorrelation is quite common in economic and financial data since such variables are rarely independent over time; for instance, today's inflation rate may depend on yesterday's inflation, and consumption patterns tend to follow a trend over several periods. While the ordinary least squares (OLS) estimators remain unbiased and consistent even in the presence of autocorrelation, they lose efficiency, meaning they are no longer the best linear unbiased estimators (BLUE).

Moreover, the standard errors derived under OLS become biased, which misleads hypothesis testing, leading to unreliable t-tests, F-tests, and confidence intervals. Autocorrelation can be caused by several factors, including omitted variables, incorrect functional form, inertia or persistence in economic activities, or data handling practices such as smoothing and aggregation. It can manifest as positive autocorrelation, where error terms follow the same direction, or negative autocorrelation, where errors alternate signs over time.

Detecting autocorrelation is therefore a crucial step in econometric analysis, commonly done through graphical inspection of residuals, the Durbin–Watson statistic, or the Breusch–Godfrey LM test. If detected, remedial measures such as model re-specification, generalized

least squares, Cochrane-Orcutt procedures, or the use of Newey–West robust standard errors may be applied to restore the efficiency of estimation and the validity of inference. In essence, understanding and addressing autocorrelation is vital in time series econometrics, as it ensures that the conclusions drawn from regression analysis remain credible and applicable in real-world economic decision-making.

### **AUTOCORRELATION:**

Autocorrelation is a mathematical representation of the degree of similarity between a given time series and a lagged version of itself over successive time intervals. It's conceptually similar to the correlation between two different time series, but autocorrelation uses the same time series twice: once in its original form and once lagged one or more time periods.

For example, if it's rainy today, the data suggests that it's more likely to rain tomorrow than if it's clear today. When it comes to investing, a stock might have a strong positive autocorrelation of returns, suggesting that if it's "up" today, it's more likely to be up tomorrow, too.

### **Causes of Autocorrelation:**

**Omission of important independent variable:** If the model does not include an important independent variable, then the error term captures its effects. This leads to dependencies between errors if the excluded variable is autocorrelated. Time series variables often exhibit autocorrelation due to their inherent nature. For example, income in the current period generally depends on the previous period's income. If we exclude the income variable from the consumption function, autocorrelation may affect the model. This is because the income variable we leave out is likely autocorrelated and is captured by the error terms. In other words, excluding the income variable can lead to autocorrelation impacting the model's accuracy.

**Nature of problem:** In some cases, error terms are autocorrelated due to the nature of economic phenomenon. Let us consider an example of agricultural production and droughts. Drought, here, is a random factor that might end up causing autocorrelation. This is because the drought will affect production not only in the current period but also in the next periods. Hence,

autocorrelation may exist because certain random factors have their effects over a longer period or across more than one period.

**Mis-specification of functional form:** A wrong functional form of the model may also cause autocorrelation. For example, if the true relationship between variables is cyclical, but the model uses a linear functional form. In such a case, the error terms might become correlated as the cyclical effects are not addressed by the explanatory variables.

**Non-stationarity:** A time series is stationary if its features (such as mean and variance) are constant over a given period of time. If the time series variables in a model are non-stationary, then the error term may also be non-stationary. As a result, varying characteristics (mean, variance, covariance etc.) over time might autocorrelate with the error term.

### **Consequences of Autocorrelation**

**Underestimation of Residual Variance:** the variance of the error term is underestimated if the errors are autocorrelated. If the autocorrelation is positive, then this problem can become even more serious.

**Variance of estimates:** in the presence of autocorrelation, the estimates of Ordinary Least Squares or OLS are still unbiased. However, the variance of the estimates will likely be larger than in the case of other methods. In the presence of autocorrelation, parameters exhibit a larger true variance, while the OLS formula underestimates the estimated variance. The variance of coefficients is smaller than their true variance. That is, we have an underestimation of the parameter or coefficient variance.

**Tests of significance:** usual tests of significance like standard errors, t-test and F-test will give misleading results and are no longer reliable because the variance of coefficients is underestimated. One cannot use these tests to establish the statistical significance of coefficients. The estimated standard errors will be smaller because of underestimated variance, thereby affecting t-values as well.

**Predictions:** in autocorrelation, the OLS model's predictions will be inefficient. This means that the predictions will have a larger variance as compared to predictions from other models like GLS.

### **Test for Autocorrelation**

The Durbin-Watson statistic is commonly used to test for autocorrelation. It can be applied to a data set by statistical software. The outcome of the Durbin-Watson test ranges from 0 to 4. An outcome closely around 2 means a very low level of autocorrelation. An outcome closer to 0 suggests a stronger positive autocorrelation, and an outcome closer to 4 suggests a stronger negative autocorrelation.

It is necessary to test for autocorrelation when analyzing a set of historical data. For example, in the equity market, the stock prices on one day can be highly correlated to the prices on another day. However, it provides little information for statistical data analysis and does not tell the actual performance of the stock.

Therefore, it is necessary to test for the autocorrelation of the historical prices to identify to what extent the price change is merely a pattern or caused by other factors. In finance, an ordinary way to eliminate the impact of autocorrelation is to use percentage changes in asset prices instead of historical prices themselves.

### **Testing For Autocorrelation – Durbin-Watson Test**

Durbin Watson test is a statistical test use to detect the presence of autocorrelation in the residuals of a regression analysis. The value of DW statistic always ranges between 0 and 4.

In stock market, positive autocorrelation (when  $DW < 2$ ) in stock prices suggests that the price movements have a persistent trend. Positive autocorrelation indicates that the variable increased or decreased on a previous day, there is a there is a tendency for it to follow the same direction on the current day. For example, if the stock fell yesterday, there is a higher likelihood it will fall today. Whereas the negative autocorrelation (when  $DW > 2$ ) indicates that if a variable increased or decreased on a previous day, there is a tendency for it to move in the opposite

direction on the current day. For example, if the stock fell yesterday, there is a greater likelihood it will rise today.

Assumptions for the Durbin-Watson Test:

- The errors are normally distributed, and the mean is 0.
- The errors are stationary.

Calculation of DW Statistics:

where;  $e_t$  is the residual of error from the Ordinary Least Squares (OLS) method.

The null hypothesis and alternate hypothesis for the Durbin-Watson Test are:

H0: No first-order autocorrelation in the residuals ( $\rho=0$ )

H1: Autocorrelation is present.

Formula of DW Statistics:

$$d = \frac{\sum_{t=2}^T (e_t - e_{t-1})^2}{\sum_{t=1}^T e_t^2}$$

Here,

- $e_t$  is the residual at time t.
- T is the number of observations.

Interpretation of DW Statistics:

- If the value of DW statistic is 2.0, it suggests that there is no autocorrelation detected in the sample.
- If the value is less than 2, it suggests that there is a positive autocorrelation.
- If the value is between 2 and 4, it suggests that there is a negative autocorrelation.

Decision Rule:

- If the Durbin-Watson test statistic is significantly different from 2, it suggests the presence of autocorrelation.
- The decision to reject the null hypothesis depends on the critical values provided in statistical tables for different significance levels.

### **Autoregressive:**

Autoregressive models are statistical models used for time series analysis, where current values are predicted based on a linear combination of past values. These models assume that past behavior influences future outcomes, making them useful for forecasting trends and patterns in data over time.

### **Types of Autoregressive Models:**

Autoregressive models vary based on the number of past values (lags) they use. The two most common types are:

1. AR(1) Model: This is a autoregressive model of order 1 which is the simplest form of an autoregressive model. In this model the current value of the time series depends only on its immediate past value along with a constant and some random noise. This model is particularly useful when the data shows strong autocorrelation at lag 1 i.e the current value depends only on the previous value. It is expressed as:

$$X_t = c + \phi_1 X_{t-1} + \varepsilon_T$$

2. AR(p) Model: It is the generalized form of the autoregressive model where the current value depends on the past P values. Choosing the correct order P is a crucial step and typically involves analyzing the Autocorrelation Function (ACF) and Partial Autocorrelation Function (PACF) plots

## AR(1) Process

The simplest autoregressive (AR) regression model is known as the AR(1) model, which stands for autoregressive of order 1. It expresses the current value of a variable as a linear combination of its lagged value.

In equation (14.1) we have presented an AR(1) model. For simplicity, let us assume that  $\mu = 0$ . Thus, equation (3.1) becomes

$$Y_t = \phi Y_{t-1} + \varepsilon_t \quad \dots (3.2)$$

If we consider time period (t-1), the AR(1) process given above, can be written as

$$Y_{t-1} = \phi Y_{t-2} + \varepsilon_{t-1} \quad \dots (3.3)$$

Let us substitute the value of  $Y_{t-1}$  from equation (3.3) in equation (3.2). we obtain

$$\begin{aligned} Y_t &= \phi (\phi Y_{t-2} + \varepsilon_{t-1}) + \varepsilon_t \\ &= \phi (\phi Y_{t-2} + \varepsilon_{t-1}) + \varepsilon_t \\ &= \phi^2 Y_{t-2} + \phi \varepsilon_{t-1} + \varepsilon_t \end{aligned}$$

If we substitute the value of  $Y_{t-1}$  ( on the basis of equation (3.2) in the above equation, then we obtain

$$\begin{aligned} Y_t &= \phi^2(\phi Y_{t-3} + \phi \varepsilon_{t-2} +) + \phi \varepsilon_{t-1} + \varepsilon_t \\ &= \phi^3 Y_{t-3} + \phi^3 \varepsilon_{t-2} +) + \phi \varepsilon_{t-1} + \varepsilon_t \quad \dots (3.4) \end{aligned}$$

On re-arrangement of terms, equation (3.4) can be written as

$$Y_t = \varepsilon_t + \phi \varepsilon_{t-1} + \phi^2 \varepsilon_{t-2} + \phi^3 Y_{t-3}$$

If we substitute the value of  $Y_{t-3}$  (from the AR(1) process) in the above equation we will obtain

$$Y_t = \varepsilon_t + \phi \varepsilon_{t-1} + \phi^2 \varepsilon_{t-2} + \phi^3 Y_{t-3} + \phi^4 Y_{t-4}$$

If we do successive substitutions of the value of  $Y$ , we obtain a series as follows:

$$Y_t = \sum_{s=0}^{t-1} \phi^s \varepsilon_{t-s} + \phi^t Y_0$$

If we assume that  $Y_0 = \varepsilon_0$ , then

$$Y_t = \sum_{s=0}^{t-1} \phi^s \varepsilon_{t-s}$$

From equation (3.6), we observe that the time series  $Y$  can be expressed as a sum of the white noise (random shocks). See that the impact of a shock on the time series ( $Y_t$ ) continues for all time to come. If  $\phi < 1$ , then the impact of a random shock dies down over time, i.e., its impact reduces gradually over time (If  $\phi = 0.7$ , then  $\phi^2 = 0.49$  and  $\phi^3 = 0.343$ ). Thus, distant observations are less important in AR models.

### **LAGGED VARIABLE:**

In time series forecasting, lagged variables play a critical role in capturing temporal dependencies and improving model accuracy. Essentially, lagged variables are past observations of your target variable, which are used as predictors for future values. By incorporating these historical data points, you can better understand patterns and trends that may influence future behavior.

To illustrate, consider a simple time series example: predicting future sales for a retail store. A lagged variable might be the sales figure from the previous month, which you would use to help forecast the current month's sales. In this context, the lagged variable serves as a key input that reflects habitual or seasonal trends in customer purchasing behavior.

Lagged variables are particularly useful in scenarios where historical patterns exert a strong influence on future outcomes. For instance, in financial markets, past stock prices (lagged variables) are often used to predict future movements. Similarly, in climatology, past temperature readings can help forecast future weather conditions.

The selection of the appropriate lag length is crucial. Too few lagged variables may lead to an under-specified model that misses critical information, while too many can result in an overfitted model that captures noise rather than signal. Analysts typically rely on domain knowledge, statistical tests, or model selection criteria to determine the optimal number of lags.

Incorporating lagged variables into your forecasting model can significantly enhance its predictive power. However, it is also important to preprocess your data appropriately, especially if your time series exhibits non-stationary characteristics. Common techniques include differencing or transformation to ensure that the statistical properties of your time series remain consistent over time.

In summary, lagged variables are a cornerstone of time series forecasting, providing essential insights from past data to predict future events. They allow models to learn from historical patterns, thus improving the reliability and accuracy of forecasts across various domains. Whether you are dealing with economic indicators, sales data, or environmental metrics, understanding and utilizing lagged variables can greatly enhance your analytical capabilities.

### **DISTRIBUTED-LAG MODEL:**

A distributed-lag model is a dynamic model in which the effect of a regressor  $x$  on  $y$  occurs over time rather than all at once. In the simple case of one explanatory variable and a linear relationship, we can write the model as

$$y_t = \alpha + \beta(L)x_t + U_t$$

where  $U_t$  is a stationary error term. This form is very similar to the infinite-moving-average representation of an ARMA process, except that the lag polynomial on the right-hand side is applied to the explanatory variable  $x$  rather than to a white-noise process  $\varepsilon$ . The individual coefficients  $\beta$  are called lag weights and the collectively comprise the lag distribution. They define the pattern of how  $x$  affects  $y$  over time.

We cannot, of course, estimate an infinite number of  $\beta$  coefficients in One practical method is to truncate the lag in to finite length  $q$ , which is appropriate if the lag distribution is effectively zero beyond  $q$  periods. Another approach is to use a functional form that allows the lag distribution in to decay gradually to zero. We saw above that a stationary autoregressive process can be expressed as an infinite moving average with declining lag weights, so a form with one or more lags of  $y$  on the right-hand side will allow infinite-length lag distributions while requiring estimation of only a small number of parameters.

One difficulty that is common to all distributed-lag models is choice of lag length, whether this be choosing the point  $q$  at which to truncate a finite lag distribution in or choosing how many lagged dependent variables to include. We defer this question until later in the chapter, after various distributed-lag models have been introduced.

## Errors in variable model:

In statistics, an errors-in-variables model or a measurement error model is a regression model that accounts for measurement errors in the independent variables. In contrast, standard regression models assume that those regressors have been measured exactly, or observed without error; as such, those models account only for errors in the dependent variables, or responses.

In the case when some regressors have been measured with errors, estimation based on the standard assumption leads to inconsistent estimates, meaning that the parameter estimates do not tend to the true values even in very large samples. For simple linear regression the effect is an underestimate of the coefficient, known as the attenuation bias. In non-linear models the

direction of the bias is likely to be more complicated.

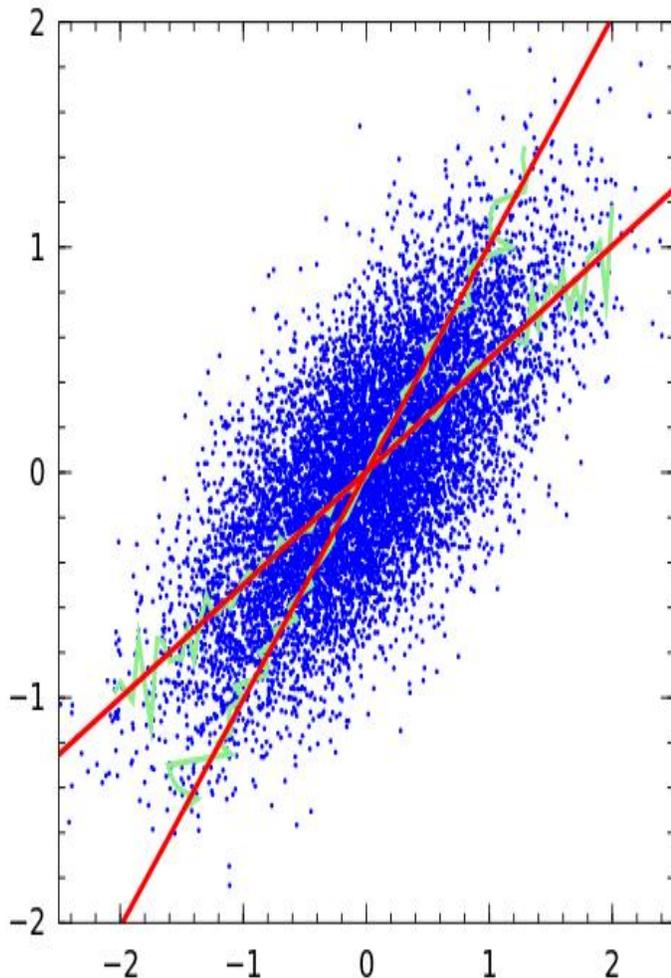


Illustration of regression dilution (or attenuation bias) by a range of regression estimates in errors-in-variables models. Two regression lines (red) bound the range of linear regression possibilities. The shallow slope is obtained when the independent variable (or predictor) is on the x-axis. The steeper slope is obtained when the independent variable is on the y-axis. By convention, with the independent variable on the x-axis, the shallower slope is obtained. Green reference lines are averages within arbitrary bins along each axis. Note that the steeper green and red regression estimates are more consistent with smaller errors in the y-axis variable.

## **Instrumental variables:**

In statistics, econometrics, epidemiology and related disciplines, the method of instrumental variables (IV) is used to estimate causal relationships when controlled experiments are not feasible or when a treatment is not successfully delivered to every unit in a randomized experiment. Intuitively, IVs are used when an explanatory (also known as independent or predictor) variable of interest is correlated with the error term (endogenous), in which case ordinary least squares and ANOVA give biased results. A valid instrument induces changes in the explanatory variable (is correlated with the endogenous variable) but has no independent effect on the dependent variable and is not correlated with the error term, allowing a researcher to uncover the causal effect of the explanatory variable on the dependent variable.

Instrumental variable methods allow for consistent estimation when the explanatory variables (covariates) are correlated with the error terms in a regression model. Such correlation may occur when:

1. changes in the dependent variable change the value of at least one of the covariates ("reverse" causation),
2. there are omitted variables that affect both the dependent and explanatory variables, or
3. the covariates are subject to measurement error.

Explanatory variables that suffer from one or more of these issues in the context of a regression are sometimes referred to as endogenous. In this situation, ordinary least squares produces biased and inconsistent estimates. However, if an instrument is available, consistent estimates may still be obtained. An instrument is a variable that does not itself belong in the explanatory equation but is correlated with the endogenous explanatory variables, conditionally on the value of other covariates.

In linear models, there are two main requirements for using IVs:

- The instrument must be correlated with the endogenous explanatory variables, conditionally on the other covariates. If this correlation is strong, then the instrument is said to have

a strong first stage. A weak correlation may provide misleading inferences about parameter estimates and standard errors.

- The instrument cannot be correlated with the error term in the explanatory equation, conditionally on the other covariates. In other words, the instrument cannot suffer from the same problem as the original predicting variable. If this condition is met, then the instrument is said to satisfy the exclusion restriction.

For example, suppose a researcher wishes to estimate the causal effect of smoking ( $X$ ) on general health ( $Y$ ). Correlation between smoking and health does not imply that smoking causes poor health because other variables, such as depression, may affect both health and smoking, or because health may affect smoking. It is not possible to conduct controlled experiments on smoking status in the general population. The researcher may attempt to estimate the causal effect of smoking on health from observational data by using the tax rate for tobacco products ( $Z$ ) as an instrument for smoking. The tax rate for tobacco products is a reasonable choice for an instrument because the researcher assumes that it can only be correlated with health through its effect on smoking. If the researcher then finds tobacco taxes and state of health to be correlated, this may be viewed as evidence that smoking causes changes in health.

## **ECONOMIC FORECASTING:**

**Definition:** Economic forecasting is the process of attempting to predict future conditions of the economy using a combination of indicators. Forecasting involves the building of statistical models with inputs of several key variables, typically in an attempt to come up with a future gross domestic product (GDP) growth rate. Primary economic indicators include inflation, interest rates, industrial production, consumer confidence, worker productivity, retail sales, and unemployment rates.

### **What is Economic Forecasting?**

Economic forecasting, the prediction of any of the elements of economic activity. Such forecasts may be made in great detail or may be very general. In any case, they describe the expected future behaviour of all or part of the economy and help form the basis of planning.

Formal economic forecasting is usually based on a specific theory as to how the economy works. Some theories are complicated, and their application requires an elaborate tracing of cause and effect. Others are relatively simple, ascribing most developments in the economy to one or two basic factors. Many economists, for example, believe that changes in the supply of money determine the rate of growth of general business activity. Others assign a central role to investment in new facilities—housing, industrial plants, highways, and so forth. In the United States, where consumers account for such a large share of economic activity, some economists believe that consumer decisions to invest or save provide the principal clues to the future course of the entire economy. Obviously the theory that a forecaster applies is of critical importance to the forecasting process; it dictates his line of investigation, the statistics he will regard as most important, and many of the techniques he will apply.

### **Limitations of Economic Forecasting:**

Economic forecasting is often described as a flawed science. Many suspect that economists who work for the White House, for instance, are encouraged to produce unrealistic projections in an attempt to justify legislation.

The challenges of economic forecasting are not limited to the government. Private-sector economists, academics, and even the Federal Reserve Board (FSB) have issued economic forecasts that were wildly off the mark.

In particular, economic forecasters have a history of neglecting to foresee crises. According to Prakash Loungani, assistant director and senior personnel and budget manager at the International Monetary Fund (IMF), economists failed to predict 148 of the past 150 recessions.

Loungani said this inability to spot imminent downturns is reflective of the pressures on forecasters to play it safe. Many, he added, prefer not to stray away from the consensus, mindful that bold projections could damage their reputations.

### **Forecasting techniques:**

Economic forecasters have a vast array of information to work with and a growing variety of techniques. A few economists, believing that just one or two key factors determine the

future course of the economy, limit their observations to these factors and develop forecasts based on them. A leading example of this is found in the school of thought that ascribes most importance to changes in the money supply. But most economists use a wider range of material.

### **Long-term forecasting:**

In recent years, increasing effort has been devoted to long-range forecasting for periods extending five, 10, or more years past the normal “short-term” forecast period of one or two years. Business has come to recognize the usefulness of such forecasts in developing plans for future expansion and financing.

Long-range forecasts usually are based on the assumption that activity toward the end of the period will reflect normal “full” employment. Given this assumption, the overall rate of growth depends on two principal factors: the number of people in the labour force and the rate at which productivity (output per worker) increases. The number of people of working age is known, barring some natural disaster (and excluding immigration), far into the future; they have already been born. Forecasters usually assume that productivity will continue to grow at the typical rates of recent decades. Expected technological developments, however, may alter the projected rate of change. The combination of changes in the labour force and productivity produces an estimate of the total growth rate for the economy.

A measure of total economic activity arrived at by such methods as these serves, in effect, as a control total for making long-range forecasts of the constituent elements of the economy. If estimates for spending by consumers, government, and business add up to more than the total of goods and services that can reasonably be expected, then the projection for one or more of these elements must be reduced. If the sum of the projected parts is less than the probable total, the analyst is likely to assume a shift in economic policy that will move the economy up to full employment by the end of the forecast period and adjust his various projections up to the appropriate total.

Long-range forecasts for individual parts of the economy depend on many of the same factors as do short-range forecasts, except that cyclical factors are usually ignored. Over the longer range, however, additional factors enter. Among the most obvious of these are growth in the population and shifts in its age composition. Changes in age composition have had a major

effect on both consumer and government spending patterns in many countries since World War II. The unusually large age cohorts born in the years following World War II had enormous influence on patterns of consumption and on labour-force composition. As young adults they tended to buy large amounts of durable goods and to add to the need for home construction; on the average, they saved less and borrowed more in relation to their incomes than most older people had. Their children constituted a secondary “baby boom,” who could expect to see their parents become the largest generation of retired persons ever known.

In addition to population pressures, a number of other trends and assumptions influence long-range forecasts. Assumptions about war and peace are obviously critical. Assumptions must be made about government spending programs; expensive new programs may bring higher taxes and less consumer spending, whereas slower growth in government spending may lead to tax reductions. Over longer periods of time, technological discoveries or changes in financial institutions can affect the overall economy. When the forecast is made for an industry or a firm, the expected introduction of new products is also important.

### **Short-Term Forecast:**

Short-term forecasting has several advantages, including its ability to provide real-time insights and facilitate quick decision-making. It is well-suited for tactical planning, inventory management, and responding to immediate changes in the market. However, short-term forecasts are susceptible to fluctuations caused by sudden external events, and they may lack accuracy when trying to predict long-term trends.

### **Understanding Short-Term Forecast Methodologies:**

Short-term forecasting aims to predict events or outcomes in the near future, typically ranging from a few days to a year. This approach is most suitable for scenarios where immediate action is required, and the influencing factors are expected to remain stable over the short period. Some common short-term forecasting methodologies include:

- 1. Time Series Analysis:** Time series analysis is a widely used method for short-term forecasting, especially when dealing with data points collected over regular intervals. This approach relies on historical data patterns to identify trends, seasonal variations, and cyclical

fluctuations. Techniques like moving averages and exponential smoothing are often employed to make short-term predictions based on past performance.

**2. Market Research and Surveys:** For businesses, short-term forecasts can heavily rely on market research and customer surveys. Collecting data on consumer behavior, preferences, and purchasing patterns can help anticipate demand fluctuations in the immediate future.

**3. Leading Indicators:** Short-term forecasts can also be derived from leading indicators, which are economic variables that tend to change before the overall economy does. Examples include stock market indices, consumer sentiment indexes, and building permits, which can provide insights into economic conditions for the next few months.

**4. Artificial Intelligence and Machine Learning:** With the advent of advanced technologies like AI and machine learning, short-term forecasting has benefited from more sophisticated prediction models. These technologies can analyze large datasets, identify patterns, and adjust forecasts in real-time, making them particularly valuable in fast-paced industries.