

Roll No: 67

NAME: HERAMB R. PAWAR

SUB: DAV

ASSIGNMENT

Q1 BRIEFLY EXPLAIN THE PHASES OF DATA ANALYTICS LIFECYCLE?

→ The Data analytic lifecycle is designed for complex data problems and Data Science projects. The cycle is iterative to portray a real project. It enables, work can return to earlier phase as new information is uncovered. The phases are as follows:

Phase 1: Discovery

This phase include gathering of:

a) Domain knowledge

b) Resources

c) Data

d) Planning & gathering
People having proper skills, who can work on the given problem and can help in solution building

e) Planning out Schedules and Deadlines

f) Problem statement formation

Phase 2: Data preparation

This phase mainly includes dealing with data

a) ETL: Extract, Transform, Load

It includes gathering data from various sources and transform it into a proper useable & acceptable format and loading it, into a database or warehouse

b) Data Conditioning:

Making Data Uniform through the process of normalization

c) Data Cleaning:

Data Cleaning includes removing of null values and irrelevant data from data available

d) Data Selection:

From the available data Selecting some of the attributes on the basis of which data analysis would be performed

Phase 3: Model Planning:

In this phase we decide the various methods that will be required for building of model such as clustering, association rule, regression, classification. It also includes selection of variables

Phase 4: Model building:

In this phase we build the model and execute it and adding up additional hardware and software resources if required

Phase 5: Communicating Results:

Communicating result whether the project was successful or not. Discussing the key findings and summarizing the narratives of stakeholders

Phase 6: Operationalize:

Preparing final reports, briefing, technical documents are submitted or delivered

Q5 What are the types of text analysis?

→ There are various types of text analysis techniques

a) Sentiment analysis:

This technique involves analyzing text to determine the emotional tone of the text. It can be used to identify positive, negative or neutral sentiments.

b) Text Categorization:

This involves categorizing text into predefined categories based on their content. For e.g. news article can be categorized into topics such as Politics, Sports, ~~Culture~~, Entertainment, etc.

c) Text clustering:

This technique involves grouping similar text together based on their content. Text clustering can be used to identify common themes or topics in large volume of text.

d) Named Entity Recognition:

This technique involves identifying and classifying entities mentioned in text such as people, ~~organizations~~, organization and locations.

e) Topic modelling:

This technique involves identifying topics or themes that are present in a collection of documents. Topic modelling can be used to understand the content of large volume of text and identify common patterns and themes.

Q7 what are 5 guidelines for better data visualization?

Ans)

1) Know your audience:

Before creating a data visualization it's important to understand who the audience is and what they are looking for. This can help to determine the type of visualization to use and level of detail to include.

2) Keep it simple:

A good data visualization should be simple and easy to understand. Avoid cluttering the visualization with too much information or using complicated charts and graphs that can confuse the audience.

3) Use appropriate visualization:

Different types of data require different type of visualization. For example bar charts are good for comparing values while line charts are better for showing trend over time. Choose the right type of visualization to effectively communicate information.

4) Highlight key insights:

A good data visualization should highlight the key insights or finding in the data. Use color annotations and other visual cues to draw attention to important point in visualization.

5) Be honest:

A data visualization should be honest and accurately represent the data. Avoid using tricks such as distorting the scale or using inappropriate axis labels that can mislead the audience.

Q8) Explain ten guidelines for better table design?

→

1) Keep it Simple:

Tables should be easy to understand and use. Avoid cluttering the table with unnecessary information or formatting that can make it difficult to read.

2) Choose the right layout:

Tables can be organized in different ways such as rows and columns or a more complex structure. Choose the layout that best suits the data being presented.

3) Use clear headings:

Headings should be clear and descriptive, providing summary of the data in the table.

4) Highlight key information:

Use formatting such as bold or color to highlight important information in table.

5) Avoid vertical scrolling:

Table should be fit ~~for single~~ screen, whenever possible to avoid vertical scrolling which makes it difficult to compare data.

FOR EDUCATIONAL USE

5) Use consistent Formatting:

Consistent formatting throughout table can make it easier to read and understand

7) Use white space effectively:

Use white space to separate different sections of table and make it easier to read

8) Keep data labels simple:

Data labels should be short and to the point avoiding long description or unnecessary details

9) Use visual cues:

Visual cues such as icons or arrows can help to draw attention to important information in the table.

10) Test the table:

Before finalizing the design, test the table with potential user to ensure that it is easy to understand

Q2 What is importance of regression model?

Ans Regression models are a type of statistical analysis used to model relationship between a dependent variable and one or more independent variables. The importance of regression models are as follows:

2) Predictive analysis:

Regression models are widely used in predictive analysis where they can be used to forecast values of the dependent variables based on the values of independent variables.

2) Relationship Identification:

Regression model can help to identify relationship between the dependent and the independent variable. This helps to understand the underlying factors that affects the dependent variable.

3) It helps in improved decision making and problem solving

3) Control for Confounding Variables:

Regression model allows to control confounding variables, which can help to isolate the effect of independent variable on dependent variable. This is important because it can help to eliminate bias and increase accuracy of the analysis.

4) Model Selection:

Regression models can help to identify the important variables and their relationships with the dependent variable. This can help to optimize the model and improve its predictive power.

5) Hypothesis testing:

Regression models can be used to test hypotheses about the relationship between the dependent variable and independent variable. This can help to identify the statistically significant relationship and provide insights to underlying mechanism that drive the relationships.

Q3 Explain logistic, simple and multiple linear regressions?

Ans → Regression is a statistical method used to model relationships between one or more independent variables and a dependent variable. The goal of regression is to find a mathematical equation that can predict the value of dependent variable based on value of independent variable. It is often used in situations where we want to understand how changes in one variable affect another. There are several types of regression they are as follows:

Simple linear regression:

In this type of regression there is only one independent variable and one dependent variable

Multiple linear Regression:

In this type of regression there are 2 or more independent variable and one dependent variable

Both simple and multiple employ a regression line known as best fit line. The linear connection is defined as $y = c + mx + e$ where 'c' denotes intercept, 'm' denotes slope of line and 'e' is the error term.

Logistic Regression:

When the dependent variable is discrete, the logistic regression technique is applicable. In other words this technique is used to compute the probability of mutually exclusive occurrences such as Pass/Fail, true/false, 0/1 and so forth. Thus the target variable can take only one of 2 values and a sigmoid curve represents its connection to independent variable and the probability has a value between 0 and 1.

- (Q4) a) What is time Series?
b) Explain ARIMA Model?
c) Explain Box Jenkins Methodology?

Ans a) Time Series Analysis is a statistical technique used to analyze data over time. In simpler terms it involves analyzing data that is collected at regular intervals over a period of time to identify patterns, trends and relationships.

Eg: time series analysis could be used to analyze monthly sales data of a company over the

Past few years to identify seasonal trends or forecast future sales. Time series analysis is a powerful tool for understanding how data changes over time and how it can be used to make informed decisions based on past trends and future predictions. It includes components like trend, season/cycle, irregularity & level.

Ans B) ARIMA Model:

- #) ARIMA stands for Autoregressive Integrated Moving Average. It is a time series model that combines the Autoregressive (AR), Integrated (I) and Moving Average (MA) model. It is a more powerful version of ARIMA model which ~~is~~ can be used to model time series data that exhibits trend and seasonality.
- #) The notation used to describe ARIMA model is ARIMA (P, I, Q) where " P " is order of autoregressive component, " I " is the degree of differencing and " Q " is the order of moving average component.
- #) The AR component models the influence of past value on current value, the "I" component models the trend and "MA" component models the influence of past error terms on current value.

The formula for ARIMA (P, I, Q) model is:
$$y_t - \hat{y}_t = c + \sum_{i=1}^P (a_i y_{t-i}) + \sum_{i=1}^q (b_i e_{t-i})$$

Ans(e)

The Box Jenkins methodology is widely used approach for time Series analysis and forecasting. It involves a series of steps that we follow to identify, model and forecast a time series.

Here are steps involved in Box Jenkins methodology:

1) Identify the time Series:

The first step is to identify the time ~~series~~ ^{series} data that needs to be analyzed. This involves collecting historical data on variable of interest such as sales or stock price over a period of time.

2) Visualize the data:

The next step is to visualize the data using plot such as time series plots, scatter plots or histogram. This helps to identify any patterns, trends or seasonality in data.

3) Check for stationarity:

Stationarity is key assumption in time series which means that the statistical properties of data remain constant over ~~time~~. The next step is to check whether time series is stationary or not. These can be done by checking for presence of trends, seasonality or periodicity in the data.

4) Make data stationary:

If the data is not stationary the next step is to make it stationary using techniques such as differencing, logarithmic transformation or detrending.

5) Identify order of differencing:

If differencing is required the next step is to identify appropriate order of differencing. This can be done by plotting autocorrelation function (ACF) and partial autocorrelation function (PACF) of the differenced data and selecting lag order that minimizes the AIC or BIC criterion.

6) Identify the ARMA model:

Once the data is stationary the next step is to identify appropriate autoregressive (AR) and moving average (MA) orders for ARMA model. This can be done by analyzing ACF and PACF plots of differenced data and selecting the orders that best fit the data.

7) Estimate model parameters:

Once the order are selected the next step is to estimate the model parameters using maximum likelihood estimation or another suitable method.

8) Check the model fit:

The next step is to check the fit of model by comparing model predictions to the actual data. This can be done using diagnostic plots such as residual plots, normal probability plots or Q-Q plots.

Q) Forecast Future Value.

Once the model is validated, the final step is to use it to forecast future value of time series.

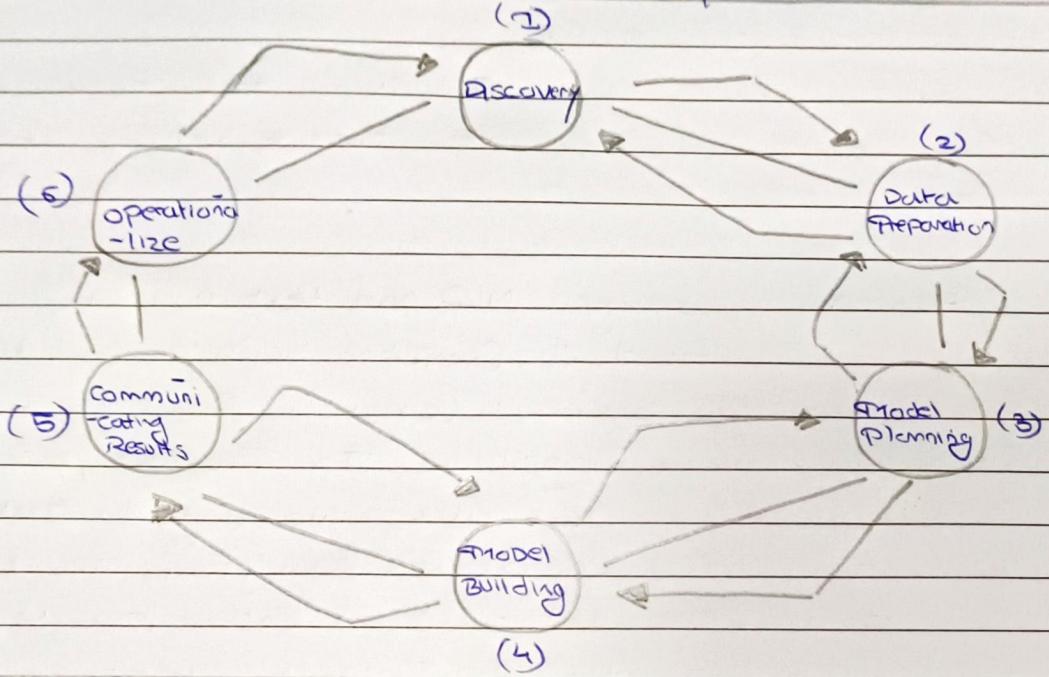
Q) Differentiate between NLP and text analysis?

Ans NLP is a subfield of Computer Science and AI concerned with interaction between computers and human language. It involves using algorithms to process, analyze and understand natural language text data. NLP aims to create computer programs that can understand and interpret human language in same way human do. This involve task such as speech recognition, language translation and sentiment analysis.

Text analysis on the other hand is a broader field that involves using statistical and ML technique to analyze and extract insights from unstructured text data. It includes various methods of processing text, such as pre-processing, feature extraction and text mining. Text analysis can be used for a variety of application including sentiment analysis, topic modelling and text classification.

(Q2)

PHASES OF DATA ANALYTICS LIFE CYCLE



(Q4) THE BOX JENKINS APPROACH:

