## Data Science Interview Questions and Answers

### Q1) What are the types of machine learning?

- Supervised learning
- Unsupervised learning
- Reinforcement Learning

### Q2) What is the Supervised learning in machine learning?

Supervised learning – When you know your target variable for the problem statement, it becomes Supervised learning. This can be applied to perform regression and classification.

Example: Linear Regression and Logistic Regression.

### Q3) What is the Unsupervised learning in machine learning?

Unsupervised learning – When you do not know your target variable for the problem statement, it becomes Unsupervised learning. This is widely used to perform Clustering.

Example: K-Means and Hierarchical clustering.

### Q4) What are the commonly used python packages?

- Numpy
- Pandas
- SCI-KIT Learn
- Matplot library

### Q5) What are the commonly used R packages?

- Caret
- Data.Table
- Reshape
- Reshape2
- E1071
- DMwR
- Dplyr
- Lubridate

### Q6) Name the commonly used algorithms.

- Linear regression
- Logistic regression
- Random Forest
- KNN

### Q7) What is precision?

The ration of predicted positive against the actual positive.

It is the most commonly used error metric is n classification mechanism.

The range is from 0 to 1, where 1 represents 100%.

### Q8) What is recall?

The ratio of the true positive rate against the actual positive rate.

The range is again from 0 to 1

### Q9) Which metric acts like accuracy in classification problem statement?

F1 Score –  2 * (Precision*Recall)/Precision + Recall

### Q10) What is a normal distribution?

When the data distribution is equally distributed as such the mean, median and mode are equal.

### Q11) What is overfitting?

Any prediction rate which has high inconsistency between the training error and the test error leads ta a high business problem, if the error rate in training set is low and the error rate ithe n test set is high, then we can conclude it as overfitting model.

### Q12) What is underfitting?

Any prediction rate which has provides low prediction in the training error and the test error leads to a high business problem, if the error rate in training set is high and the error rate inthe test set is also high, then we can conclude it as overfitting model.

### Q13) What is a univariate analysis?

An Analysis that can be applied to one attribute at a time is called as a univariate analysis.

Boxplot is one of the widely used univariate model.

Scatter plot and cook's distance are other methods used for bivariate and multivariate analysis.

### Q14) Name few methods for Missing Value Treatments.

Central Imputation – This method acts more like central tendencies. All the missing values will be filed with mean and median mode respective to numerical and categorical datatypes.

KNN – K Nearest Neighbour imputation.

Distance between two or multiple attributes are calculated using Euclidian's distance and the same will be used to treat the missing values. Mean and mode will agaibe n used as in CI.
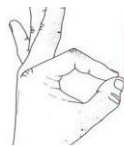
### Q15) What is the Pearson correlation?

Correlation between predicted and actual data can be examined and understood using this method.

The range is from -1 to +1.

-1 refers to negative 100% whereas +1 refers to positive 100%.

The formula is Sd(x)*m/Sd.(y)

### Q16) How and by what methods data visualizations can be effectively used?

In addition to giving insights in a very effective and efficient manner, data visualization can also be used in such a way that it is not only restricted to bar, line or some stereotypic graphs. Data can be represented in a much more visually pleasing manner.

One thing have to be taken care of is to convey the intended insight or finding correctly to the audience. Once the baseline is set. Innovative and creative part can help you come up with better looking and functional dashboards. There is a fine line between the simple insightful dashboard and awesome looking 0 fruitful insight dashboards.

### Q17) How to understand the problems faced during data analysis?

Most of the problem faced during hands on analysis or data science is because of poor understanding of the problem in hand and concentrating more on tools, end results and other aspects of the project.
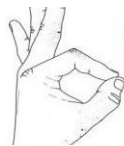
Breaking the problem down to a granular level and understanding takes a lot of time and practice to master. Coming back to square one in data science projects can be seen in lot of companies and even in your own project or kaggle problems.

### Q18) Advantages of Tableau Prep?

Tableau Prep will reduce a lot of time like how its parent software (Tableau) does when creating impressive visualizations. The tool has a lot of potentials in taking professionals from data cleaning, merging step to creating final usable data that can be linked to Tableau desktop for getting visualization and business insights. A lot of manual tasks will be reduced and the time can be used to make better findings and insights.

### Q19) What is the common perception about visualization?

People think visualization as just charts and summary information. But they are beyond that and drive business with a lot of underlying principles. Learning design principles can help anyone build effective and efficient visualizations and this Tableau prep tool can drastically increase our time on focusing more important part. The only issue with Tableau is, it is paid and companies need to pay for leveraging that awesome tool.

### Q20) What are the time series algorithms?

Time series algorithms like ARIMA, ARIMAX, SARIMA, Holts winters are very interesting to learn and use as well to solve a lot of complex problems for businesses. Data preparation for time series analysis plays a vital role. The stationarity, seasonality, cycles and noises need time and attention. Take as much time as you would like to make the data right. Then you can run any model on top of it.

### Q21) How to choose the right chart in case of creating a viz?

Using the right chart to represent data is one of the key aspects of data visualization and design principle. You will always have options to choose from when deciding on a chart. But fixing to the right chart comes only by experience, practice and deep understanding of end-user needs. That dictates everything in the dashboard.

### Q22) Where to seek help in case of discrepancies in Tableau?

When you face any issue regarding Tableau, try searching in the Tableau community forum. It is one of the best places to get your queries answered. You can always write your question and get the query answered with an hour or a day. You can always post on LinkedIn and follow people.
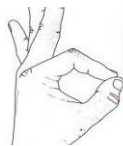
### Q23) Now companies are heavily investing their money and time to make the dashboards. Why?

To make stakeholders more aware about the business through data. Working on visualization projects helps you develop one of the key skills every data scientist should possess i.e. Thinking from the shoes of the end user.

If you're learning any visualization tool, download a dataset from kaggle. Building charts and graphs for the dashboard should be the last step. Research more about the domain and think about the KPIs you would like to see in the dashboard if you're going to be the end user. Then start building the dashboard piece by piece.

### Q24) How can I achieve accuracy in the first model that I built?

Building machine learning models involves a lot of interesting steps. 90% accuracy models don't come in the very first attempt. You have done a lot of better feature selection

techniques to get that point, which means it involves a lot of trial and error. The process will help you learn new concepts in statistics, math and probability.

### Q25) What is the basic responsibility of a Data Scientist?

As a data scientist, we have the responsibility to make complex things simple enough that anyone without context should understand, what we are trying to convey.

The moment, we start explaining even the simple things the mission of making the complex simple goes away. This happens a lot when we are doing data visualization.

Less is more. Rather than pushing too much information on to readers brain, we need to figure out how easily we can help them consume a dashboard or a chart.

The process is simple to say but difficult to implement. You must bring the complex business value out of a self-explanatory chart. It's a skill every data scientist should strive towards and good to have in their arsenal.

### Q26) How do I enhance a SAS analyst?

Step 1: Earn a College Degree. Businesses prefer SAS programmers who have completed a

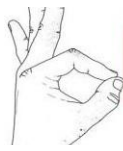statistics or computer science bachelor's degree program.

Step 2: Acquire SAS Certification.

Step 3: Consider Getting an Advanced Degree.

Step 4: Gain SAS Program Coding Work Background.

### Q27) What does SAS stand out to be the best over other data analytics tools?

Ease to understand: The provisions included in SAS are remarkably easy to learn. Further, it offers the most suitable option for those who already are aware of the SQL. On the other hand, R comes with a steep training cover which is supposed to be a low-level programming style.

Data Handling Capacities: it is at par the most leading tool which also includes the R& Python.

If it advances before handling the huge data, it is the best platform to engage Graphical Capacities: it comes with functional graphical capacities and has a limited knowledge field.

It is useful to customize the plots Better tool management: It benefits in a release the updates with regards to the controlled conditions.

This is the main reason why it is well tested. Whereas if you considered R&Python, it has open contribution also the risk of errors in the current development is also high.

### Q28) What is RUN-Group processing?

To practice RUN-group processing, you start the system and then submit many RUNgroups.

A RUN-group is a group of records that contain at least one product group including ends with a RUN statement. It can contain different SAS statements such as AXIS, BY, GOPTIONS,

LEGEND, Power, or WHERE.

### Q29) Definitions of is BY-Group processing?

Definitions for BY-Group Processing. is a method of preparing observations from one or numerous SAS data sets that are arranged or ordered by importance of individual or more

shared variables. All data sets that are being connected must include one or more BY variables.

### Q30) What is the right way to validate the SAS program?

The OPTIONS OBS=0 through the commencement of the code needs to be written but if yourself require to perform the same then their mind be any log which gets recognized by the

colors that get highlighted.

### Q31) Do you know any SAS functions and Call Routines?

Can be a mutable type, uniform, or any SAS expression, including different use. This product also a letter from contentions that SAS allows are called by special purposes. Multiple

arguments are separated with a comma.

### Q32) What is means by precision and Recall?

Recall:

It is known as a true real rate. The number of positives that your model has claimed related to the original defined number of positives available during this data.

Precision:

It is also known as a positive predicted value. This is more based on the prediction. That indicates a time like a number of accurate positives that the model needs when compared to the

number of positives it actually claims.

### Q33) What is deep learning?

Deep learning is a process where it is considered to be a subset of machine learning process.

### Q34) What is the F1 score?

The F1 score is defined as a measure of a model's performance.

### 35) How is F1 score is used?

The average of Precision and Recall of a model is nothing but F1 score measure. Based on the results, the F1 score is 1 then it is classified as best and 0 being the worst

### Q36) What is the difference between Machine learning Vs Data Mining?

Data mining is about working on unlimited data and then extract it to a level anywhere the unusual and unknown patterns are identified.

Machine learning is any method about a study whether it closely relates to design, development concerning the algorithms that provide an ability to certain computers to capacity to learn.

### Q37) What are confounding variables?

These are obvious variables in a scientific model that correlates directly or inversely with both the subject and the objective variable. The study fails to account for the confounding factor.

### Q38) How can you randomize the items of a list in place in Python?

Consider the example shown below:

from random import shuffle

x = ['Data', 'Class', 'Blue', 'Flag', 'Red', 'Slow']

shuffle(x)

print(x)

The output of the following code is as below.

['Red', 'Data', 'Blue', 'Slow', 'Class', 'Flag']

### Q39) How to get indices of N maximum values in a NumPy array?

We can get the indices of N maximum values in a NumPy array using the below code:

import numpy as np

arr = np.array([1, 3, 2, 4, 5])

print(arr.argsort()[-3:][::-1])

Output

[ 4 3 1 ]

### Q40) How make you 3D plots/visualizations using NumPy/SciPy?

Like 2D plotting, 3D graphics is beyond the scope of NumPy and SciPy, but just as in this 2D example, packages exist that integrate with NumPy. Matplotlib provides primary 3D plotting in

the mplot3d subpackage, whereas Mayavi produces a wide range of high-quality 3D visualization features, utilizing the powerful VTK engine.

### Q41) What are the types of biases that can occur during sampling?

Some simple models of selection bias are described below. Undercoverage occurs when some members of the population live badly represented inside the sample. … The survey relied on a service unit, drawn of telephone directories and car registration lists.

o Selection bias
o Under coverage bias
o Survivorship bias

### Q42) Which Python library is used for data visualization?

Plotly. The fifth tool is Plotly, also called as Plot.ly because of its main platform online. It is an interactive online visualization tool that is being used for data analytics, scientific graphs, and

other visualization. This contains some great API including one for Python

### Q43) Write code to sort a DataFrame in Python in descending order.

DataFrame.sort_values(by, axis=0, ascending=True, inplace=False, kind='quicksort',

na_position='last')[source]

Sort by the values along either axis

Parameters:

by: str or list of str

Name or list of names to sort by.

if an axis is 0 or 'index' then by may contain index levels and/or column labels

if the axis is 1 or 'columns' then by may contain column levels and/or index labels

Changed in version 0.23.0: Allow specifying index or column level names.

axis : {0 or 'index', 1 or 'columns'}, default 0

Axis to be sorted

ascending: bool or list of bool, default True

Sort ascending vs. descending. Specify list for multiple sort orders. If this is a list of bools, must
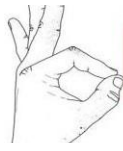
match the length of the by.

in place: bool, default False

if True, perform operation in-place

kind: {'quicksort', 'mergesort', 'heapsort'}, default 'quicksort'

Choice of sorting algorithm. See also array.np.sort for more information. mergesort is the only

stable algorithm. For DataFrames, this option is only applied when sorting on a single column or

label.

na_position : {'first', 'last'}, default 'last'

first puts NaNs at the beginning, last put NaNs at the end

Returns:

sorted_obj: DataFrame

## Q44) Why you should use NumPy arrays instead of nested Python lists?

Ans : let's say you have a list a of numbers, and you want to add 1 to every element of the list.

In regular python, you would do:

a = [6, 2, 1, 4, 3]

b = [e + 1 fore in a]

Whereas with numpy, you simply have to do:

import numpy as np
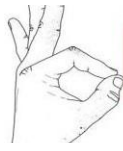
a = np.array([6, 2, 1, 4, 3])

b = a + 1

It also works for every numpy mathematics function: you can take the exponential of every

element of a list using np.exp for example.

## Q45) Why is an import statement required in Python?

Ans : To be able to use any functionality, the respective code logic needs to be accessible for the Python interpreter. With the help of the import statement, we can use specific scripts.

However, there are thousands of such scripts available and every script available cannot be used at once. Hence we import statement to use only the scripts that we want to use

- import pandas as pd
- import numpy as np

### Q46) What is alias in import statement? Why is it used?

Ans : Aliases are used in import statements for ease of usage. If the imported module has a large name, for example import multiprocessing . Everytime we want to access any scrript present in multiprocessing module, we need  to use the word multiprocessing.

However if an alias is used, import multiprocessing as mp, we can simply replace the words multiprocessing with mp

### Q47) Are the aliases used for a module fixed/static ?

Ans : No, the aliases are not pre-fixed. The alias can be named as per your convenience. However, the documentation of a respective module sometimes specifies the alias to be used for ease of understanding.

### Q48) How to access a specific script inside a module?

Ans : If the whole module needs to be imported, we simply can use from pandas import *

### Q49) What is a nonparametric test used for?

Ans : Non parametric tests do not assume that the data follows a specific distribution. They can be used whenever the data do not meet the assumptions of  parametric test.

### Q50) What are the pros and cons of Decision Trees algorithm?

Ans :

Pros – Easy to interpret. Will ignore irrelevant independant variables since information gain will be minimal. Can handle missing data. Fast modelling.

Cons – Many combinations are possible to create a tree. There are chances that it might not find the best tree possible.

**Q51) Name some Classification Algorithms.**

Ans : Linear Classifiers: Logistic Regression, Naive Bayes Classifier, Decision Trees, Random Forest, Neural Networks, K Nearest Neighbor.

**Q52) What are pros and cons of Naive Bayes algorithm?**

Ans :

Big sized data is handled easily

Multiclass perfomance is good and accurate

It is not process intensive

Cons: Assumea independence of predictor variables.

**Q53) What are the types of Skewness?**

Ans : A dataset that is skewed right or left are the two types.

**Q54) What is skewed data?**

Ans : A data distribution that is has skewed data towards the right or left.

**Q55) What is the skewness of this data? 27 ; 28 ; 30 ; 32 ; 34 ; 38 ; 41 ; 42 ; 43 ; 44 ; 46 ; 53 ; 56 ; 62**

Ans : The data set is skewed left

**Q56) What is an outlier?**

An outlier is a value that is very much away from the rest of the values in the data set.

**Q57) Mention the characteristics of symmetric data distribution?**

Ans : The mean is equal to the median and the tails of the distribution are balanced.

### Q58) What are the applications of data science?

Ans : Optical character recognition, recommendation engines, fitering algorithms, personal assistants, advertising, surveillance, autonomous driving, facial recognition and more.

### Q59) Define EDA?

Ans : EDA [exploratory data analysis] is an apporach to analysing data to summarise their main characteriscs, often with visual methods.

### Q60) What are the steps in exploratory data analysis?

Ans :

- Make summary of observations
- describe central tendencies or core part of dataset
- desribe shape of data
- identify potential associations
- develop insight into errors, misssing values and major deviations

### Q61) What are the types of data available in Enterprises?

Ans :

- Structured data
- unstructured data
- big data from social media, surveys, pictures, audio, video, drawings, maps.
- Machine generated data from instruments
- real time data feeds

### Q62) What are the various types of analysis on type of data?

Ans :

Univariate – 1 variable

bivariate – 2 variables

multivariate – more than 2 variables

**Q63) What is difference between primary data and secondary data?**

Ans :

Data collected by the interested/self is primary data. This data is collected afresh and first time.

Someone else has collected the data and being used by you is secondary data.

**Q64) What is the difference between qualitative & quantitative ?**

Ans : Quantitative method analyses the data based on numbers. Qualitative method analyses the data by attributes.

**Q65) What is histogram?**

Ans : Histogram is the accurate representation of numerical data based on their occurrences/frequencies.

**Q66) What are the common measures of central tendancies?**

Ans :

- Mean
- Median
- Mode

**Q67) What are quartiles?**

Ans : Quartiles are three points in the data, that divide the data into four groups. Each group consisting of a quarter of data.

**Q68) What are the commonly used error metrics in regression tasks?**

Ans :

MSE – Mean squared error – Average of square of errors

RMSE – Root mean square error – root of MSE

MAPE – Mean absolute percentage error

## Q69) What are the commonly used error metrics for classification tasks?

Ans :

- F1 score
- Accuracy
- Sensitivity
- Specificity
- Recall
- Precision

## Q70) What is it called when there are more than 1 explanatory variables in the regression task?

Ans : Multiple linear regression

## Q71) What are residuals in a regression task?

Ans : The difference between the predicted value and the actual value is called the residual.

## Q72) What are the main classifications in Machine learning?

Ans :

- Supervised learning
- Unsupervised learning
- Reinforcement learning

## Q73) What are the main types of supervised learning tasks?

Ans :

Classification task [categorical in nature]

Regression task [continuous in nature]

## Q74) Can Random forest be used for classification and regression?

Ans : Yes, it can be used

Give a simple representation for Linear Equation.

Ans : Y = mx + c ; where y is the dependant variable; c is the independant variable;m is slope

## Q75) What is R square value?

Ans : R squared values tells us how close the regression line is fit to the actual values.

## Q76) What are some common ways of imputation?

Ans : Mean imputation, median imputation, KNN imputation, Stochastic regression, substitution

## Q77) What is the difference between series and list

Ans :

list is size and data mutable

series is data mutable but not size mutable

## Q78) Which function is used to get descriptive statistics of a dataframe?

Ans : describe()

## Q79) What parameter is used to update the data without explicitly assigning data to a variable.
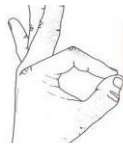
Ans : Inplace is used to assign result of function to itself. If inplace = True , there is no need to explicitly assign to a variable

## Q80) What is the difference between a dictionary and a set?

Ans :

Dictionary has key value pair

set does not have key value pairs

set has only unique elements

### Q81) How to create a series with letters as index?

Ans : Series({'a':1,'b':2}) will create a and b as index. 1 and 2 as their respective values.

### Q82) Which function can be used to filter a DataFrame?

Ans : The query function can be used to filter a dataframe.

### Q83) What is the function to create test train split?

Ans : From sklearn.metrics import test_train_split . This function is used to create test train split from the data.

### Q84) What is pickling?

Ans : Pickling is the process of saving a data structure into the physical drive or hard disk.

### Q85) What is unpickling?

Ans : Unpickling is used to read a pickled file from hard disk or physical storage drive.

### Q86) What are the most common web frameworks of Python?

Ans : Django and Flask.

### Q87) How to convert n number of series to a dataframe?

Ans : DataFrame(data = {'col1':series1,'col2':series2}).

### Q88) How to select a section of a dataframe?

Ans : Using iloc and loc functions the rows and columns can be selected.

### Q89) How are exceptions handled in Python?

Ans : Exceptions can be handled using the try except statements.

**Q90) Is multiprocessing possible in python?**

Ans : Yes it is possible using the multiprocessing module.

**Q91) Can the values be replaced in tuple?**

Ans : No values cannot be replaced in tuple as tuple is data immutable

**Q92) What are lambda function in Python and how it is different from def (defining functions) in Python?**

Ans : Lambda function in Python is used for evaluating an expression and then return a value. Where as def needs a function name, and the program logic is broken into smaller chunks. Lambda is an inline function consisting of only a single expression, It can take any number of arguments.

**Q93) Difference between supervised and unsupervised machine learning?**

Ans : Supervised learning is a method where it needs training specified data. When it gets to Unsupervised learning it doesn't need data labeling.

**Q94) How to differentiate from KNN and K-means clustering?**

Ans : KNN is standing for the K- Nearest Neighbours, it remains classified because a supervised algorithm.K-means is an unsupervised cluster algorithm.

**Q95) What is your opinion on our current data process &nbs p;?**

Ans : This type of questions signifies asked and the individuals must to carefully listen to their value case and at this same time, the return should be in a constructive also insightful manner. Based on your responses, the interviewer mind has a future to review and know whether you imply a vague reply to their team or not.

**Q96) Difference between Machine learning and Data Mining?**

Ans :

- Data mining is about going about unstructured data and when extracting this to a level anywhere that interesting also unknown patterns remain identified.

- Machine learning is any process or a concept whether it closely relates designing, development of the algorithms that give an experience within these machines on the capacity to learn.

### Q97) Explain about from capture of the correlation between continuous and categorical variable?

Ans : It is possible to that using ANCOVA technique. It exists for Analysis of Covariance. It is used to calculate this association between continuous including categorical variables.

### Q98) Difference between an Array and a Linked list?

Ans : An array is an established method of collection objects. A linked program is a group of objects that are prepared into sequential order.

### Q99) Difference between "long" and "wide" format data?

Ans : In the wide form, each subject's happened responses will remain in a separate row, and each answer is into a separate column. In the long format, each data is a one-time time by subject. You can understand data in wide form by that fact that columns usually design groups.

### Q100) What do you know by the term Normal Distribution?

Ans :

- Data is usually distributed under many ways including a bias on the port or over the benefit or it can all be jumbled up.
- However, there continue indications that data is distributed on a central position without bias to the left or right more gives natural order in some form of a bell-shaped curve.

### Q101) Differences between overfitting and underfitting

Ans :

- In statistics and machine learning, individual of that most basic tasks is to fit one model on a collection of training data, so doing to be ready to provide reliable predictions of general untrained data.
- Underfitting happens at a statistical design or machine learning algorithm cannot get this underlying trend of the data.

### Q102) Differentiate between univariate, bivariate and then multivariate analysis.

Ans :

- Univariate analyses are detailed statistical analysis methods which can be changed based upon the number of variables involved in a distributed period of time.
- The bivariate analysis tries to explain that difference between two variables at an individual time as in a scatterplot.
- Multivariate analysis contracts including the single study from and then a couple of variables to understand the effect from variables to some responses.

### Q103) Do you explain the word Botnet?

Ans : A botnet is a type of bot running on an IRC network created with a Trojan.

### Q104) What is data visualization?

Ans : Data visualization is a common word, which helps to understand the importance of data in a visual context.

### Q105) How to Clean Data is an Important Part of the Process?

Ans : Cleaning the data at the point of work is a great job. If we try to fix the sources of uncontrollable data like this plane, our time can take up to 80%.

### Q106) Which language is suitable for text analysis? R or Python?
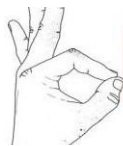
Ans : Because Python has a rich library called Python, researchers allow high quality data analysis tools and data structures, while R does not have this feature. So Python is more suited to text analysis.

### Q107) What Is a System?

Ans : A referrer system is still widely used in many fields such as film recommendations, music priorities, social references, research articles, search queries. Recommended systems operate by using recommendation based on filtering or personality based approach. This type of system based on a person's past behavior works to create a model for the future. It will predict future buying, movie viewing or reading the public book. Creating a filtering approach using the unique characteristic of the items when prescribing additional items.

### Q108) Compare SAS, R, and Python programming?

Ans : R and Python are the most important programming languages for machine learning algorithms.

SAS: One of the most widely used analytical tools used by the largest large companies on earth. It can come up with some of the best statistical functions, graphical user interface, but a price tag, so it will not be accepted by small companies immediately

R: The best part about R is an open source tool, so education and research community is generously used. It is a robust tool for statistical calculation, graphical representation and reporting. Due to its open source nature, the latest features are always updated and available to everyone.

Python: Python is a powerful open source programming language, it's very easy, works well with other tools and technologies. The best part about Python is that it has a lot of libraries and community-created blocks. It has statistical activity, model building and more functions.

### Q109) What are the different benefits of language?

Ans : R programming language is a software package used for graphical representation, statistical computing, data handling, and calculation.

The features of the R programming environment include the following:

- Detailed collection tools for data analysis
- Drives for arrays and array calculations
- Data analysis technique for graphical representation
- The most advanced yet simple and useful programming language
- It supports machine learning applications in detail
- It works as a link connecting between various software, tools and databases
- Create flexible and powerful high quality reviewing analysis
- Provides a strong set of ecosystems for various purposes
- It's useful when you need to solve a data-based problem

### Q110) What are the two main elements of the hottest architecture?

Ans :

- HDFS and YARN are two main components of the Hadoopo structure.
- HDFS-HOWTO distributed file system. This is the Hadoop top job distributed database. It is possible to save and retrieve the number of data at any time.
- YARN- still stands for another source of negotiation. It modifies resources and handles workloads.

### Q111) How do Data Scientists use statistics?

Ans : Statistics helps to see data scientists samples, data for late insights, and to convert large data to large intelligences. It helps customers get a good idea of what to expect. Data scientists can learn about consumer behavior, interest, involvement, retention, and last convertible statistics. It helps to create powerful data models to estimate some specifications and calculations. Everything can be changed into a powerful business idea by informing users exactly what they want.

### Q112) What is Logistic Recession?

Ans : It is a statistical technique or a model for analyzing the database and predicting binary effects. The effect must be zero or one or a binary effect of yes or no. Random forest is an important technique used for classification, resilience and other tasks in the database

### Q113) Why is data important in data analysis?

Ans : Since the data comes from many sources, it is important to ensure that data analysis is adequate. Data purification is very important. Data cleaning manipulates the process of detecting and repairing data data, ensuring that the data is complete and accurate, and if the components of the inappropriate data are omitted or modified according to requirement. This process will be compatible with data fights or batch processing.

Once the data has been purified, it confirms the rules of the data on the system. Data cleaning is an important part of data science because corruption is neglected due to human negligence, exchange or storage of other things. Data recovery data data is taken by a large portion of time and effort of the scientist, due to the speed and speed it receives from multiple data.

### Q114) Describe univariate, bivariate and multivariate analysis?

Ans : As the name suggests, these are single, double or multiple variables with analytical methods.

So a distinct analysis will have a variable, thereby causing no relationship and reasons. The key to the unique analysis is the briefing of the data to find the results of the process, and finding the forms within it.

Bivariate analysis deals with the relationship between the two segments of data. Connected data are related sources of this set, or models. While data has a relationship, there are various tools to analyze such data, including C-squire tests and t-tests. If data is measured, it can be analyzed using a graphical plot or a scattering graph. The strength of the connections between the two data sets will be tested in the Bivariate analysis.

### Q115) How is Machine Learning Used in Real World Scenes?

Ans : Here are some situations where machine learning can be found in real world applications:

- Online: Customer understanding, ad targeting and review
- Search Engine: Ranking pages depending on the search's personal choices
- Funding: Assessing Investment Opportunities and Risks, Finding Fraud Operations
- Medicare: Designing medicines depending on the patient's history and needs
- Robotics: Machine learning to handle situations outside of normal
- Social Media: Linking Understanding Relationships and Recommendations
- Extracting information: Creating questions to get answers from databases on the web

### Q116) What are the different features of the mechanical learning process?
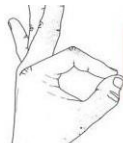
Ans : In this post I will discuss the elements involved to solve a problem using machine learning.

Domain knowledge

This is the first step to understanding various features from data and to learn more about the data we handle. The computer introduces more and more about the type of domain we're dealing with.

### Q117) What Is Normal Distribution?

Ans : It is a set of continuous variations in the form of a regular curve or in the form of a bell curve. This can be considered as a continuous probability distribution and useful in statistics. It is very useful for analyzing variables and their relationships while having a more general partition curve and normal sharing curve.

The normal distribution curve is symmetrical. Ordinary distribution models approach normal distribution to the extent of increase. Using the Central Limit Theory is very easy. This method helps to understand random data by creating the order, and helps to understand the results using the Bell shaped drawing.

### Q118) What is Linear Recreation?

Ans : This is the most commonly used method for predictive analysis. The linear replication method is used to describe the relationship between a dependent variable and one or the other variable. The main task of linear recursion is the method of applying a single line in a scattering plot.

Linear Recreation has the following three modes:

- Determining and analyzing data communication and direction
- Evaluating the model
- To ensure the use and validity of the model
- It is widely used in scenes that have a catch effect. For example, you should know the effect of a specific action to determine the various consequences.
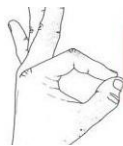
### Q119) What Is Interpolation and Extrapolation?

Ans : The interpolation and approval rules are important in any statistical analysis. Extraction is a valuation or evaluation of facts by determining it or taking an evaluation or to an unknown area or area. It is a technique that can penetrate something using the available data.

On the other hand, interpolation is the method of determining a certain value between a value of a certain value or a value of values. This is especially useful if you have data between the two sides of a particular region, but you do not have enough data points at the specified point. This is when you sort the interpolation to determine the required value.

### Q120) What is Power Analysis?

Ans : Energy analysis is an important part of the test design. It relates to the process of determining the sample size required to find a certain amount of effect with a certain degree of warranty. This allows the use of a particular probability in a sample.

### Q121) What is Q-Meaning? Can K Choose a K-Method?

Ans :K-material cluster is a fundamental supervised learning method. This is the method of classifying data using a specific set of clusters known as K clusters. It is used to group data to find data unity.

It defines K centers, each in a cluster. Clusters are defined as K groups before pre-defined K. K points are aligned to cluster centers. Objects are allocated to their closest cluster center. The objects in a cluster are closely interrelated to each other and the other clusters vary as much as possible. K is very good for large packages of data.

### Q122) How does data modeling change from database format?

Ans :

- Data modeling: This can be considered the first step for a database design. Data modeling creates a conceptual model based on the relationship between different data models. This process moves to the ideological model of the ideological model of the ideological model. It includes a systematic method for using data modeling techniques.
- Database design: This is the process of creating a database. Database design creates a publication of the detailed data model of the database. Database design which includes a detailed logical model of a database is strictly speaking, but it includes physical design options and storage parameters.

### Q123) The Difference Between Data Modeling and Database Design?

Ans : Data Modeling – Software Modeling Data modeling (or modeling) is the process of creating a data model for information systems by using data data modeling techniques.

### Q124) Can you briefly describe the scientific method?

Ans :

- Data is collected from sensors in the environment.
- Data is "cleaned up" or a data set (usually a data table) for processing.

### Q125) What is the recommended system?

Ans : A companion of a database filtering system to predict a customer's product or charts. Recommendations are widely used in movies, news, research articles, products, social tips, music, etc.

### Q126) Why data analysis is an important part of the analysis?

Ans : Because data sources are increasing, data increases due to the number of sources increases the amount of data generated in these sources: data data researchers or data scientists are a complex process that can convert data from multiple sources. It may take up to 80% to clean data that generates an important part of the analysis work.

### Q127) What is Linear Recreation?

Ans : The linear lag is the value of a variable Y, measured by the second variable X. X, which is the predictor variable and Y is referred to as the criterion variable.

### Q128) Explain what is the regulation and why it is useful. Regularization?

Ans : Organizing is an act of stimulating lowest in order to change the coincidence parameter. It is often done by plus a constant number of existing weight vectors. This standard is popular L1 (laso) or L2 (ridge). Sample predictions should then reduce the loss of functionality calculated in the formal training package.
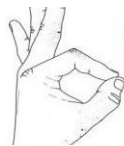
### Q129) What is TFT / ITF Vectation?

Ans: tf-idf The frequency-inverse file frequency is narrow, a numerical statistic intended to reflect how important a document or corpus is. It is often used as a weight factor in information retrieval and text mine. The Tf-idf value document increases the number of times the document appears in the document, but the word frequency in the corpus which helps to fix the fact that some words are normally more frequent.

### Q130) What is the Cluster Model?

Ans : Cluster model is a technique used when a wide area is hard to analyze widespread spaces, and a simple random sample is not used. Cluster model is a sample of a set or component of each modeling unit is a probability model.

E.g., a researcher wants to study the education program of Japanese high school students. He can split entire Japan into different clusters (towns). The researcher then selects several clusters based on his research with simple or systematic random sampling.

### Q131) What is the regulatory model?

Ans : The regulatory model is a statistical technique where elements are selected from a sorted sample frame. In the formal model, you can improve the lists circuitry, so when you finish the list, it comes back to top. The best example of a proper model is the probability of the equation.

### Q132) What are the agenciers and the agenwals?

Ans : In data analysis, we generally calculate the number of animators into a correlation or coordinate team. Eigenvectors are dynamically executed or stretched by turning a certain linear transition into directions.

Eigenvalue eigenvector is referred to as the factor of change in strength or compression.

### Q133) Can you quote some examples of false positives that are more false than negative ones?

Ans :

- Let's first see what the wrong positions and the wrong nexus are.
- It is wrong to say that you have incorrectly identified an event as a category a.k.a type I error.
- In case of non-occurrences, a.k.a Type II error, you are mistakenly mistaken.

Example 1: In the field of medicine, I think you should give chemotherapy to patients. A patient was sent to the hospital, studying the study, but she was not really cancer. This is a false positive event. This patient is more likely to start chemotherapy instead of cancer. In the absence of cancer cell, chemotherapy can cause specific damage to its normal healthy cells and can even cause serious illness.

Example 2: An e-commerce company may say that we have decided to award $ 1000 gift card to customers at least $ 10,000 worth of goods. Because they allow at least 20% of profits to be sold for over $ 10,000, they send free airmail to 100 customers at least in the buyer's position. Now we have purchased $ 1000 gift boxes for customers but have indicated $ 10,000 worth of purchase.

### Q134) Can you cite some of the worst negative examples of negative negative than negative ones?

Ans :

Example 1: An 'A' airport having high security threats is based on certain characteristics that identify whether or not specific passengers are threatened or not. Due to the shortage of employees, passengers predict the danger posed by their prediction model. What happens if the real threat is threatened by the customer's airport model?

Example 2: What if a judge or a judge decides to release a criminal?

Example 3: If you reject a good person based on your prediction model, if you meet him a few years later, do you realize that you are a wrong negative?

### Q135) Can you quote some examples of both false positives and misinformation?

Ans :The primary source of lending to banking industries, but your repayment rate is not good, you will not make any profits, but you will cause huge losses.

Banks do not want to lose good customers while at the same time do not want to get bad customers. In this situation, both the wrong position and the wrong negative are very important to measure.

### Q136) Can you explain the difference between a verification set and test set?
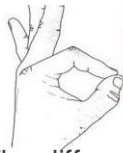
Ans :

- An assessor can be used as part of the training system to use the package parameter selection and skip the model specifically.
- On the other hand, a test set is used to test or evaluate the performance of a trained machine learning model.
- In simple terms, the differences may be brief; The package parameters of the training match

### Q137) What is a Pilab?

Ans :

- The pileup is a set, which integrates NumPy, SciPy and Matplotlib into single namespaces.

- The difference between tuples and lists in Python is the state.

- A list can be used to store multiple locations while Tuples is used in a dictionary to store notes in places. The lists are variable, while the doubles can be changed, ie they can not be edited.

- Specify some libraries in Python used for data analysis and scientific computing.

- NumPy, SciPy, Seaborn, Pandas, Matplotlib, SciKit Data Analysis and Scientific Calculations There are some libraries in Python used.

- Write a code to sort by column (n-1) in NumPy

- This can be achieved using the argsort () function. You can take an array X to sort the X (x-2) code (n-1).

### Q138) If you provide employees' first and last names, what type of data in Python stores them?

Ans :You can use a list of the first name and last name that an element contains, or the dictionary uses.

### Q139) Explain the use of decorators?

Ans : A decorator is a function that takes another function and extends the second functionality without explicitly changing it. They can be used to change the classes and functions of the code. With the help of decorators, a code code can be executed before or after the execution of the original code.

The output of the code below will be:

def foo (i = []):

i.append (1)

Come back

>>>foo ()

>>>foo ()

The output for the above code-

[1] [1, 1]

The argument for function foo is evaluated only once when function is defined. However, since this is a list, the entire list is replaced by the use of 1 in each step.

### Q140) Which tool should you use to find the bugs?

Ans : Tools for Finding Errors in Python isPhyllent and Bicenter. Pilliant is used to verify that a module satisfies all index standards. A standard analysis tool that helps find the bugs in the source code.

### Q141) What is the difference between the range () and xrange () in the Python?

Ans : The array () function returns a list. The Xrange () function provides an object that acts like a platform to generate numbers according to requirement.

### Q142) How to sort items of the list in Python?

Ans : You can use Shift (lst) function to rearrange the list of items in Python.

### Q143) What is PEP8?

Ans : PEP8 is a set of index guides in Python, which can be used by programmers to write code that is easy to use for other users.

### Q144) What is the monkey grafting in Python?

Ans : Programmers in Python can modify or extend the other code in the motion by using the chrome mechanics technique. It comes in handy during testing but the code is hard and there is not a good practice to use it in a production environment.
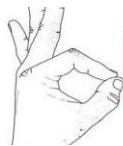
### Q145) What does it mean to understand the list?

Ans : Listening Understanding is a process of creating a list when doing some functions in data, so that it can be accessed using a distributor.

### Q146) You should find that data is stored in HDFS format and how the data is structured. Which command should you use to identify the names of HDFS keys?

Ans : In this case, the following command can be used

hf.keys ()

Note: The HDFS file is loaded as H5py as HF.

### Q147) What is the output of the code below?

Ans :

Word = 'aeioubcdfg'

Print word [: 3] + word [3:]

The above code output will be: 'aeioubcdfg'.

When the two pieces of the pieces collide and the "+" operator fits the string, it breaks the string into pieces.

### Q148) How do you see if a panda data information is empty or not?

Ans : The character df.empty is used to verify that the data in the panda data is empty.

### Q149) Which Python Library is used by Machine Leader?

Ans : SciKit-Learn is a crazy library.

You are given a list of numbers. Understand a single list of Python to create a new list that contains only the numbers that contain the numbers from the list of elements. For example, if the value is even [4] in the list, it must be added to the new release list because it is a symbol, even if it is in the list [5] even if it is not in the index.
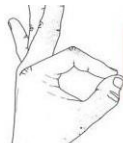
[x = 2 == 0 if x is x [1: 2]
The above code will take all the numbers in the code and reject the odd numbers.

### Q150) What are the basic assumptions for linear backlash?

Ans : Errors of error, statistical independent errors, interviewing and additional skills.

### Q151) Can the formula be written to calculate the R-square?

Ans :

R-squire can be calculated using the form below –

1 – (total squares / squares total)

### Q152) What is the benefit of dimension reduction before applying an SVM?

Ans : Support vector machine learning algorithm works best on low space. If the number of features is larger compared to the number of surveys, it will benefit from diminishing the dimension before the SVM is applied.

### Q153) What is Data Science ?

Ans: It's a science and methodology of acquiring data, pre-processing data, analyzing data , visualizing data and drawing meaningful conclusions from the data to drive the business need

### Q154) Who is a Data Scientist ?

Ans: A data scientist is a Person Trained In Mathemtaics, Statistics And Computer Science, who is adept in acquiring data from various sources, has the skills to clean and preprocess the data, analyze and visualize the data, draw inferences make predictions and present the results in the form of a convincing story to the client.

### Q155) What are the skills required in Data Science ?

Ans:

- Mathematics- College Arithmetic, Linear Algebra, Calculus
- Statistics- Data Types, Summary Statistics, Correlation, Regression, Central Limit Theorem, T-test, ANOVA
- Programming- ETL tools like Informatica, Querying in SQL, Data Analysis in R & Python , data visualization and creating dashboards using Tableau

### Q156) What is Machine Learning ?

Ans: Machine Learning is that part of data science which deals with making predictions. The prediction could be in the form of a number or a class or a group. It is the icing in the cake of data science. It come under the category known as predictive analytics.

### Q157) What is the difference between Traditional Programming and Machine Learning ?

Ans: In traditional programming, data is fed to a block of code and we get the desired output, whereas in Machine learning it's the other way round , ie. In MACHINE LEARNING DATA WRITES CODE and the output is a program/model.

### Q158) What' s the difference between a Regression and a Classification problem?

Ans: The difference is in the label. If the label (or target ) is a numerical value eg, a stock price , salary etc it is a regression problem whereas if the label is binary or multi-class like fraud/not fraud, yes/no , etc then it is a classification problem.

### Q159) What are the types of Machine Learning Algorithms?

Ans: Machine Learning can be broadly classified as Supervised and Unsupervised. In Supervised Learning, the data comes in the form of features and label. The Algorithm is trained on this data and a trained model is developed which is then used on the unseen data to make predictions. Supervised learning is used typically  . In unsupervised learning the data has only features, no output labels. The algorithm learns by itself and groups the subjects accordingly. This kind is used typically in customer segmentation problems.

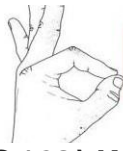### Q160) Name some Python Libraries used in Machine Learning .

Ans:

- Numpy
- pandas
- matplotlib
- seaborn
- scikitlearn

### Q161) Give examples of supervised and unsupervised ML algorithms.

Ans:

- Supervised- MLR, KNN, SVM, Logistic regression, Decision Tree, Random Forest
- Unsupervised- k-Means, Hierarchical, t-SNE

### Q162) What are the main components of a data science project ?

Ans:

- Understanding Business Requirement
- Data acquisition and preparation
- Data Analysis, Visualization & inference
- Project Management

### Q163) What percentage of time is usually required for each component in data science projects ?
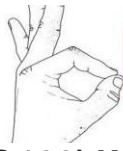
Ans:

- Understanding BR- 10%
- Data Acquisition and Preparation – 70%
- Data Analysis visualization and inference – 10%
- PMP- 10%

### Q164) What is Artifical Intelligence ?

Ans: It's the ability of a computer to learn by itself by being exposed to lots and lots of data. It uses repetitive iterations of predictions and error corrections to get better and better at predictions. ML is a subset of AI and DL is a subset of ML.

### Q165) What is Deep Learning (DL)?

Ans: DL is the ability of a computer to mimic the human brain. A typical deep learning architecture consists of an input layer, an output layer and hidden layer(s) of neurons. The basic difference between ML and DL is that in ML the programmer decides based on available data the features to be considered , whereas in DL, the algorithm itself detetcts the significant features by assigning weights to them and readjusting the weights using a principle known as back-prorogation. So in the final layer the features may not be even recognizable by the human. For example a Deep Learning classifier could very accurately ( almost 96% accuracy) predict whether a given brain scan has lesion or not. But doctors will not be informed as to why the model predicted it so.

### Q166) What is Backpropagation?

Ans: It is the method by which a neural network trains itself. Every time a data row is fed into the deep learning algorithm, weights are assigned to the synapses associated with each neuron. Based on these weights a prediction is given which is then compared to the ground truth and an error value is calculated using a Cost function. This value is then propagated backward through the entire neural network and the weights are re-adjusted and another prediction is made and again the error value is calculated. Through a number of iterations of this process, the weights are optimized to obtain minimum error and we say that the network is trained or optimized.

### Q167) What is Stochastic Gradient Descent?

Ans: The stochastic gradient descent is an optimizing method to locate the local minimum of the cost function. In this the weights updated are using the slopes of points in the Cost function. First weights are randomly chosen (close to 0 but not 0) and then the gradient(slope) of the point in the cost function is calculated. If it is negative then we descend right with steps proportional to the gradient. Now the slope of the new point will be positive. So we descend left with steps proportional to the new gradient. This process is iteratively done till the local minimum of the Cost function is reached. Thus we find the weights that minimizes cost function.

### Q168) Name some supervised and unsupervised deep learning algorithms.

Ans:

- Supervised-ANN, CNN, RNN, LSTM
- Unsupervised-SOMs, Autoencoders

### Q169) Name some Python libraries used in Deep Learning

Ans:

- Keras
- tensorflow
- pytorch
- cv2

### Q170) What is Data Science?

Ans:

- Data science or Data analytics is a process of analyzing large set of data points to get answers on questions related to that data.
- Combination of Predictive Analytics and Data Mining is called Data Science.

### Q171) How to Assign Code to the List?

Ans: Using this syntax continuation, we can assign symbolic value to any list.

Mylist = [None] * 10 (none of the 10's list)

### Q172) Give me two important tasks in the pants?

Ans:

- Series
- Data frame

### Q173) What is the difference between iloc and loc activity?

Ans:

- Take the pieces based on the lock labels (features).
- It uses the Index based position.

### Q174) What package is used to import data from the Oracle server?

Ans: We use CX_Oracle modules to link Python with Oracle server.

### Q175) Import of Flat File / CSV in Baidan

Ans:

- Read_csv
- Generatorrext

### Q176) How to read an Excel file without a file file in the Byndah?

Ans: Read the Excel file using the Xlsreader module and manipulate it.

### Q177) Data Science Carrier

Ans:

- Statisticians
- Business Analyst
- Mathematician
- Professor
- Risk Analyst
- Data Analyst
- Content Analyst
- Statistics Trainer
- Data Scientist
- Consultant
- Biostatistician
- Econometrician

### Q178) Data Science job areas

Ans:

- Census
- Ecology
- Medicine
- Election
- Crime
- Economics
- Education
- Film
- Sports
- Tourism

### Q179) Scope and Applications of Statistics

Ans:

- Statistics and actuarial science
- Statistics and Commerce
- Statistics and Economics
- Statistics and Medicine

- Statistics and Agriculture
- Statistics and Industry
- Statistics and Information Technology
- Statistics and Government
- Big Data

### Q180) What do the review process do?

Ans: You can change the data without changing the data.

### Q181) What do Dummies do?

Ans: It can alter the duplicate / cursor variables alternately.

### Q182) Difference between distinct, bivariate and multivariate analysis?

Ans: Analysis Data Simple analytics analysis of data analysis that contains only one variable. This is a variable because it does not cope with the causes or relationships. The main purpose of the unique analysis is to describe data and discover the forms inside it. Ex. Average, method, intermediate, range, variance, max, at least, quartz and standard deviation

Bivariate Analysis is used to find out if there is a relationship between two different variables.

Ex. Skater Plate, Co

Multidimensional analysis is analysis of three or more variables. There are a number of ways to analyze diversity according to your goals.

Ex. Cluster Analysis, Multiple Recreation Analysis

### Q183) What is the curtain?

Ans: In many setback analysis, one of the forecasts is in contrast to the other predictor / dependent, then this problem is known as collinearity.

### Q184) Why is a useful metric?

Ans: Can you measure a size from one person? Relationship (but some of the below examples can not be as often as we can see) can refer to the presence of a causal relationship. Many modeling techniques are used as a base size and base combination

### Q185) What is the removal of data backward in advance?

Ans: Disadvantage eliminates at least every significant aspect of each reaction that starts with all the features and improves the performance of the model. We do this again until there is no improvement in removing features.

### Q186) What is Unequal Data?

Ans: Balanced data sets for classification issues are special classes, and class distribution between classes is not uniform. In general, they are created by two types: the majority (negative) class and minority (positive) class. Because this type of set data mining is a new challenging problem, because standard classification protocols typically considers a consistent training package and most of it is thinking of a pro in class.

### Q187) The performance of the K modular system?

Ans:

The general procedure is as follows:

- Rotate the database randomly.
- Divide databases into k groups
- For each individual group:
- Take out as a group or take the test data package
- Take the remaining groups as training data set
- Apply a sample on training and evaluate the test package
- Retain valuation value and reject the model

Shorten the model's ability to model the model's rating scores

*Q188) What is standardization?*

Ans: Data rate evaluation is the process of restructuring one or more attributes, whereby they are the average value of the 0 and the standard nomination 1. Standard considers your data to be a gauge (bell curve) distribution. This will not be true, but your attribute distribution is a very effective technique