Self introduction

Hello,

I'm Satheesh, and I'm currently based in Bangalore, India. I've had a diverse career path, and I'd like to highlight the aspects that are most relevant to this job.

I currently work as a Python Developer at Tata Consultancy Services, a role I've held since September 2020. In this position, I've gained valuable experience in data analysis, machine learning, and Django. My work involves creating software that extracts essential information from documents. I'm responsible for ensuring the accuracy and reliability of our analytical models. I also specialize in data visualization, using tools like Matplotlib and Seaborn to present findings effectively. Additionally, I handle data preparation for machine learning, using Pandas and NumPy to clean and format data. One of my notable achievements in this role was improving data accuracy by 20% through the application of Natural Language Processing (NLP) and machine learning techniques on unstructured documents.

My role also extends to managing servers for deploying and testing our applications, including User Acceptance Testing (UAT) and System Integration Testing (SIT). I'm responsible for installing and configuring essential software packages and dependencies on servers to support application functionality. Moreover, I've created user-friendly interfaces within Django for real-time configuration of data extraction rules and progress monitoring.

In terms of education, I hold a B. Tech in Computer Science and Engineering from Srinivasa Ramanujan Institute of Technology, JNTUA.

In conclusion, I'm genuinely excited about the opportunity to work here. My skills in Python, machine learning, data analysis, and software development make me a valuable asset for this role. I'm enthusiastic about using my knowledge to contribute to the success of this company.

Thank you for considering my application.

# Objective and Additional Details

Objective:

"To leverage my expertise in developing innovative solutions and my proficiency in Django, Python, NLP, and machine learning for a dynamic role in software engineering. Committed to enhancing document information extraction processes, contributing to cutting-edge projects, and staying at the forefront of technological advancements."

Additional Details:

- Project Experience:
    - Led the development of a Document Information Extraction tool using Django, Python, NLP, and machine learning modules.
    - Achieved a 20% increase in data accuracy through the application of advanced NLP and machine learning techniques on unstructured documents.
    - Spearheaded feature engineering, model selection, and evaluation, ensuring the reliability of analytical models.
    - Utilized Matplotlib and Seaborn to create visually appealing charts, graphs, and dashboards for effective communication of findings.
- Data Processing and Management:
    - Conducted thorough data preprocessing with Pandas and NumPy, optimizing datasets for machine learning models.
    - Managed servers for deploying and testing the application, overseeing UAT and SIT.
    - Installed and configured essential software packages and dependencies, ensuring seamless functionality.
- User Interface and Experience:
    - Engineered user-friendly Django interfaces for real-time configuration of data extraction rules and progress monitoring.
    - Designed and implemented intuitive Django templates, ensuring an enhanced and responsive user experience within the application.
- Technological Proficiency:
    - Stayed current with NLP, machine learning, and web development advancements.
    - Developed RESTful APIs using Flask for enhanced interoperability and scalability.
    - Proficient in Anaconda and Jupyter Notebook interface.
- Collaboration and Problem-Solving:
    - Collaborated closely with system administrators to monitor server performance and troubleshoot issues promptly.

- Demonstrated a proactive approach to problem-solving and contributed innovative solutions for document information extraction challenges.

# What are your strengths and weaknesses as a data analyst?

I'm quick at data preprocessing, analyzing data, and creating data visualizations. I have a talent for spotting patterns in data. However, one area where I feel I could improve is in selecting the best features from a set of many features, even though I am effective in feature selection. Additionally, when a dataset has too many incomplete or inaccurate data points, it could impact the accuracy of my analyses. I'm aware of these challenges and committed to enhancing my skills in these specific areas.

[Top 65+ Data Analyst Interview Questions and Answers [2023]](#)

Data analysis vs data mining
Data analysis involves the process of cleaning, organizing, and using data to produce meaningful insights. Data mining is used to search for hidden patterns in the data.

Data validation is like giving your data a check-up to make sure it's accurate, reliable, and follows the rules. It's a process of ensuring that the data entered into a system is correct and meets certain criteria or standards. This helps prevent errors, inconsistencies, and ensures that the information you're working with is trustworthy. It's like having a data doctor making sure everything is in good health! Eg:

Sure, let's say you have a form on a website where users enter their age. Data validation in this context would involve checking that the age entered is a reasonable number (not negative, not a crazy high value), and that it's in the expected range for an age (let's say between 1 and 120 years).

So, if someone tries to enter an age of -5 or 150, the data validation process would catch that as an error and prompt the user to enter a valid age within the specified range. It's like a safety net for your data, ensuring that only sensible and accurate information gets through.

# How to know if a data model is performing well or not?

I'd check if the model is mostly right, using a kind of accuracy check. Then, there's this confusion map to see where it gets mixed up. For yes-or-no situations, there's a curve to show how well it separates things. I'd also make sure it doesn't take forever to do its job. And, most importantly, see if it's actually helping the business or just chilling there."

- A well-designed model should offer good predictability. This correlates to the ability to be easily able to predict future insights when needed.
- A rounded model adapts easily to any change made to the data or the pipeline if need be.
- The model should have the ability to cope in case there is an immediate requirement to large-scale the data.
- The model's working should be easy and it should be easily understood among clients to help them derive the required results.

# Data Cleaning in brief.

Data Cleaning is also called Data Wrangling. As the name suggests, it is a structured way of finding erroneous content in data and safely removing them to ensure that the data is of the utmost quality. Here are some of the ways in data cleaning:

- Removing a data block entirely
- Finding ways to fill black data in, without causing redundancies
- Replacing data with its mean or median values
- Making use of placeholders for empty spaces

# What are some of the problems that a working Data Analyst might encounter?

There can be many issues that a Data Analyst might face when working with data. Here are some of them:

- The accuracy of the model in development will be low if there are multiple entries of the same entity and errors concerning spellings and incorrect data.
- If the source the data being ingested from is not a verified source, then the data might require a lot of cleaning and preprocess before beginning the analysis.
- The same goes for when extracting data from multiple sources and merging them for use.
- The analysis will take a backstep if the data obtained is incomplete or inaccurate.

**Data profiling** is like getting to know your data on a personal level. It's the process of analyzing and summarizing the key characteristics of a dataset. Imagine it as a detective work for your data, figuring out what's in there, how clean it is, and what kind of shape it's in.

In simpler terms, data profiling involves examining the structure, quality, and content of your data. It helps you understand things like the types of values in your dataset, identifying any missing or duplicate entries, and checking if the data conforms to certain standards or expectations. It's like giving your data a thorough check-up before diving into analysis or using it for important decisions.

## What are the scenarios that could cause a model to be retrained?

Sure thing! When you want to impress your interviewer with your understanding of data model performance, you can say:

1. **Accuracy Check:** I'd make sure the model is getting things right most of the time. You know, hitting the bullseye!

2. **Confusion Matrix Talk:** I'd look at this matrix thingy that shows where the model is getting confused. It's like a map to see where it's doing well and where it's a bit lost.

3. **ROC Curve and AUC:** For binary choice situations, I'd check this curve and area thing to see how well the model is separating the good stuff from the not-so-good stuff.

4. **Time and Resources:** I'd also peek into how fast the model is working. We don't want it taking forever to tell us something.

5. **Real-Life Impact:** And, of course, I'd keep an eye on how the model is helping the business. Is it making things better or just adding confusion?

Throw in some smiles and confidence, and you'll have them nodding in agreement!

What are the steps involved when working on a data analysis project?
Embarking on a data analysis project is like taking a journey. Here are the steps to guide you along the way:

1. **Define Objectives:** Clearly outline the goals of your analysis. What are you trying to discover or achieve?

2. **Collect Data:** Gather the relevant data needed for your analysis. It's like packing your essentials for the journey.

3. **Clean and Preprocess:** Clean up your data. Remove errors, handle missing values, and transform it to a format that's ready for analysis. It's like making sure your gear is in top shape.

4. **Explore the Data:** Take a deep dive into your data. Look for patterns, trends, and outliers. It's like exploring the terrain before setting a course.

5. **Feature Engineering:** Create new features or modify existing ones to enhance your analysis. It's like upgrading your tools for better performance.

6. **Choose Analysis Methods:** Decide on the techniques and methods you'll use for analysis. It's like selecting your navigation tools for the journey.

7. **Apply Statistical Methods:** Use statistical tests or models to draw insights from the data. It's like using a compass to find your direction.

8. **Visualize Results:** Create visual representations of your findings. Charts and graphs can make the insights more understandable. It's like creating a map to guide others on your journey.

9. **Interpret and Conclude:** Draw conclusions from your analysis. What did you learn? How does it align with your objectives? It's like reaching the destination and reflecting on the journey.

10. **Document and Communicate:** Document your process, findings, and methodologies. Communicate your results to stakeholders in a clear and understandable way. It's like sharing your travelog with others.

11. **Iterate if Necessary:** If needed, go back and refine your steps. Maybe collect more data, try different analyses, or explore additional questions. It's like adjusting your route based on unexpected discoveries.

Remember, every data analysis project is unique, and flexibility is key. Enjoy the journey and the discoveries along the way!

**Can you name some of the statistical methodologies used by data analysts?**
**Certainly! Here's a brief explanation for each:**

1. **Regression Analysis:** Examines relationships between variables, useful for predicting outcomes.

2. **Hypothesis Testing:** Assesses whether sample data supports or contradicts a hypothesis about a population.

3. **Descriptive Statistics:** Summarizes main features of a dataset, including measures like mean and standard deviation.

4. **Cluster Analysis:** Groups similar data points, revealing natural patterns within the data.

5. **ANOVA (Analysis of Variance):** Analyzes differences among group means in a sample.

6. **Chi-Square Test:** Determines if there's a significant association between categorical variables.

7. **Time Series Analysis:** Studies data collected over time to identify patterns and trends.

8. **Bayesian Methods:** Updates probabilities based on new evidence, especially useful in handling uncertainty.

These methodologies help data analysts draw meaningful insights and make informed decisions from diverse datasets.

## Which are the types of hypothesis testing used today?

There are many types of hypothesis testing. Some of them are as follows:

- Analysis of variance (ANOVA): Here, the analysis is conducted between the mean values of multiple groups.
- T-test: This form of testing is used when the standard deviation is not known, and the sample size is relatively small.
- Chi-square Test: This kind of hypothesis testing is used when there is a requirement to find the level of association between the categorical variables in a sample.

## What are some of the data validation methodologies used in data analysis?

Many types of data validation techniques are used today. Some of them are as follows:

- Field-level validation: Validation is done across each of the fields to ensure that there are no errors in the data entered by the user.
- Form-level validation: Here, validation is done when the user completes working with the form but before the information is saved.
- Data saving validation: This form of validation takes place when the file or the database record is being saved.

- Search criteria validation: This kind of validation is used to check whether valid results are returned when the user is looking for something.

What are some of the data validation methodologies used in data analysis?
Sure, let's explore a couple of data validation methodologies with definitions and examples:

1. **Field-level validation:**
   - **Definition:** Checking individual fields to ensure data accuracy and consistency.
   - **Example:** Verifying that a "Date of Birth" field contains a valid date format and falls within a reasonable range.

2. **Form-level validation:**
   - **Definition:** Validating data integrity and coherence across multiple fields in a form before saving.
   - **Example:** Confirming that in a survey form, if a respondent selects "Male" as their gender, the "Pregnancy Status" field is logically disabled.

3. **Data saving validation:**
   - **Definition:** Ensuring the accuracy and reliability of data before it is saved to a database or file.
   - **Example:** Before storing customer information, validating that the provided postal code matches the selected city and state.

4. **Search criteria validation:**
   - **Definition:** Checking the validity of user-entered criteria to ensure meaningful search results.
   - **Example:** Verifying that in an e-commerce platform, the entered price range for a product search has the minimum price less than or equal to the maximum price.

These methodologies collectively contribute to maintaining data quality, preventing errors, and ensuring that data is reliable and meaningful for analysis. They help in catching inconsistencies and inaccuracies early in the data analysis process.

1. **Recall (Sensitivity or True Positive Rate):**
   - **Definition:** Recall measures the ability of a classification model to capture all the relevant instances or positives in the dataset. It is the ratio of true positives to the sum of true positives and false negatives.
   - **Formula:** Recall = True Positives / (True Positives + False Negatives)
   - **Example:** In a medical diagnosis scenario, recall would tell you the proportion of actual positive cases (e.g., patients with a disease) that the model correctly identified.

2. **Precision (Positive Predictive Value):**
   - **Definition:** Precision measures the accuracy of the positive predictions made by a model. It is the ratio of true positives to the sum of true positives and false positives.
   - **Formula:** Precision = True Positives / (True Positives + False Positives)

- **Example:** In a spam email classification task, precision would indicate the proportion of emails predicted as spam that are actually spam, avoiding false positives (legitimate emails marked as spam).

In summary, recall focuses on finding all the relevant instances, while precision emphasizes the accuracy of positive predictions. Depending on the task and its consequences, you might prioritize one metric over the other. Balancing recall and precision is often crucial in achieving an optimal model for a given problem.

**t test vs z test**

The t-test and z-test are both statistical methods used for hypothesis testing and making inferences about population parameters based on sample data. Here's a quick comparison between the two:

**T-Test:**
1. **Sample Size:** Suitable for small sample sizes (typically less than 30).
2. **Population Standard Deviation:** Used when the population standard deviation is unknown.
3. **Distribution:** Utilizes the t-distribution, which has fatter tails than the normal distribution, making it more appropriate for small samples.
4. **Example:** One-sample t-test is often used to compare the mean of a sample to a known or hypothesized population mean. Paired and independent two-sample t-tests compare means between groups.

**Z-Test:**
1. **Sample Size:** Preferable for large sample sizes (typically greater than 30).
2. **Population Standard Deviation:** Used when the population standard deviation is known.
3. **Distribution:** Utilizes the standard normal distribution (z-distribution), which approximates the normal distribution for larger sample sizes.
4. **Example:** Z-tests are used when comparing sample means to known population values, especially when the sample size is large. They are also employed for proportion tests.

**Key Considerations:**
- Use the t-test when dealing with small samples or when the population standard deviation is unknown.
- Use the z-test when dealing with large samples and when the population standard deviation is known.
- In general, if the sample size is large (typically >30), the t-distribution converges to the normal distribution, and the t-test approaches the z-test.

Both tests are valuable tools in statistics, and the choice between them depends on the specific characteristics of the data and the hypothesis being tested.

## What are the ideal situations in which t-test or z-test can be used?
The choice between a t-test and a z-test depends on various factors related to the sample size, the known or unknown population standard deviation, and the nature of the comparison. Here are some ideal situations for each:

**T-Test:**
1. **Small Sample Size:** Use a t-test when dealing with a small sample size (typically less than 30).

2. **Unknown Population Standard Deviation:** If the population standard deviation is unknown, the t-test is more appropriate.

3. **Paired Samples:** When comparing means of two related groups (paired or dependent samples), such as before-and-after measurements or matched pairs.

4. **One-Sample T-Test:** Used when comparing the mean of a sample to a known or hypothesized population mean.

**Z-Test:**
1. **Large Sample Size:** When dealing with a large sample size (typically greater than 30), the z-test is often preferred.

2. **Known Population Standard Deviation:** If you know the population standard deviation and have a sufficiently large sample, the z-test is appropriate.

3. **Comparing Two Independent Samples:** When comparing the means of two independent groups (unpaired or independent samples), and the sample size is large, the z-test may be used.

4. **Proportion Z-Test:** Specifically for comparing sample proportions to population proportions.

In practical terms, if you have a small sample or lack information about the population standard deviation, lean towards a t-test. If your sample is large and you know the population standard deviation, a z-test might be more suitable. Always consider the specific characteristics of your data and the hypothesis being tested when making the choice.


How are outliers detected?
Detecting outliers involves identifying data points that deviate significantly from the majority of the dataset. Here are some common methods for outlier detection:

1. **Visual Inspection:**
   - **Scatter Plots:** Visual examination of data points on a scatter plot to identify those that lie far from the general pattern.
   - **Box Plots:** Visualization of the distribution of data using box plots, with points outside the "whiskers" considered potential outliers.

2. **Descriptive Statistics:**
   - **Z-Score:** Calculation of the z-score for each data point, indicating how many standard deviations it is from the mean. Points with high absolute z-scores may be outliers.

- **IQR (Interquartile Range):** Use of the IQR to define a range and flagging data points outside this range as potential outliers.

3. **Machine Learning Models:**
   - **Clustering:** Application of clustering algorithms (e.g., k-means) and considering data points in small clusters as potential outliers.
   - **Isolation Forest:** Use of isolation forest algorithms that isolate outliers by recursively partitioning the data.

4. **Statistical Methods:**
   - **Grubbs' Test:** A statistical test for detecting a single outlier in a univariate dataset.
   - **Modified Z-Score:** Modification of the traditional z-score for increased robustness against outliers.

5. **Distance-Based Methods:**
   - **Mahalanobis Distance:** Calculation of the Mahalanobis distance, which accounts for correlations between variables.
   - **Euclidean Distance:** Identification of points with unusually large distances from the centroid.

6. **Domain-Specific Knowledge:**
   - **Subject Matter Expertise:** In some cases, domain knowledge is essential for recognizing outliers that may have specific significance.

It's crucial to consider the characteristics of the data and the context of the analysis when choosing an outlier detection method. Often, a combination of approaches is used to obtain a more comprehensive understanding of potential outliers in the dataset.

**How can one handle suspicious or missing data in a dataset while performing analysis?**

Handling suspicious or missing data is a crucial step in data analysis to ensure the accuracy and reliability of the results. Here are some common strategies:

1. **Identify and Understand the Causes:**
   - Investigate the reasons for suspicious or missing data. It could be due to data entry errors, system issues, or intentional omissions.

2. **Data Cleaning:**
   - Correct errors in the data by replacing or imputing incorrect values. This might involve using the mean, median, or mode for numerical data, or the most frequent category for categorical data.

3. **Outlier Detection and Treatment:**
   - Identify outliers using methods like z-scores or visualization tools. Depending on the nature of the data, you may choose to remove outliers or transform the data to reduce their impact.

4. **Imputation Techniques:**
   - Use imputation methods to fill in missing values. Common techniques include mean imputation, forward or backward filling, regression imputation, or more sophisticated methods like k-nearest neighbors imputation.

5. **Data Transformation:**
   - Transform data to handle outliers or skewed distributions. Techniques such as log transformation or box-cox transformation can be applied to make the data more suitable for analysis.

6. **Advanced Imputation Techniques:**
   - Consider more advanced imputation methods like multiple imputation, which accounts for the uncertainty associated with missing data by creating multiple imputed datasets.

7. **Domain Knowledge:**
   - Leverage domain expertise to make informed decisions about suspicious or missing data. Sometimes, certain patterns or reasons for missing data are best understood by those familiar with the subject matter.

8. **Create a Separate Indicator Variable:**
   - Create an additional variable indicating whether a value is missing or suspicious. This allows you to track and analyze the impact of missing or suspicious data on your results.

9. **Statistical Tests:**
   - Use statistical tests to assess whether missing data patterns are random or systematic. This can inform the choice of imputation methods or other strategies.

10. **Communication:**
    - Clearly document the steps taken to handle suspicious or missing data in your analysis. Transparent communication ensures the reproducibility and validity of your results.

Remember that the approach may vary depending on the specific characteristics of the dataset and the analysis goals. The key is to handle suspicious or missing data systematically, maintaining data integrity and preserving the reliability of the analysis.

### what should i tell for Why should we hire you for data analyst?

When responding to the question "Why should we hire you for a data analyst position?" it's essential to highlight your unique strengths, skills, and experiences that make you a strong fit for the role. Here's a structured response you can use as a guide:

---

1. **Strong Analytical Skills:**
   - Highlight your proficiency in analytical tools and techniques, such as Python, Pandas, NumPy, and Scikit-learn.
   - Mention specific instances where your analytical skills have been instrumental in extracting valuable insights from complex datasets.

2. **Proven Track Record:**
   - Discuss any notable projects or achievements, emphasizing your role in improving data accuracy and driving successful outcomes.
   - Provide concrete examples of how you've applied data analysis to solve real-world problems.

3. **Technical Expertise:**
   - Showcase your technical proficiency in relevant programming languages (Python, C, Java), frameworks (Flask, Django Rest API), and databases (MySQL, Postgres).
   - Highlight your experience with version control tools like Git and API development using tools like Postman.

4. **Machine Learning Skills:**
   - If applicable, emphasize your machine learning capabilities, especially if you've successfully implemented ML models in a previous role.
   - Discuss your experience with model selection, evaluation, and feature engineering.

5. **Web Technologies and Design:**
   - Mention your skills in web technologies (HTML, CSS, JS) and your ability to create user-friendly interfaces and templates.
   - Discuss any experience you have in designing and implementing Django templates.

6. **Continuous Learning and Adaptability:**
   - Express your commitment to continuous learning and staying updated with emerging technologies.
   - Demonstrate your adaptability to new tools and methodologies to ensure you remain at the forefront of data analysis practices.

7. **Communication Skills:**
   - Highlight your ability to communicate complex technical concepts to non-technical stakeholders.
   - Discuss instances where you've effectively presented findings and collaborated with cross-functional teams.

8. **Passion for Innovation:**
   - Convey your passion for leveraging data to drive innovation and informed decision-making.
   - Share how your enthusiasm for data analysis contributes to a positive and dynamic work environment.

9. **Conclusion:**

- Summarize by reiterating that your unique blend of skills, experiences, and passion make you the ideal candidate for the data analyst position.
  - Express your excitement about the opportunity to contribute to the team's success.

Remember to tailor your response based on the specific requirements of the job and the company culture. Providing concrete examples and metrics where possible can make your response more compelling.

# What is inheritance and what are its types?

Inheritance is a way to create a new class by inheriting properties and methods from an existing class. It allows code reusability and saves time. Types of inheritance are single, multiple, multilevel, and hierarchical. Single inheritance involves inheriting properties and methods from a single parent class. Multiple inheritance involves inheriting properties and methods from multiple parent classes. Multilevel inheritance involves inheriting properties and methods from a parent class, which in turn inherits from another parent class. Hierarchical inheritance involves multiple child classes inheriting properties and methods from a single parent class.

# What is the difference between variable and object?

Variables are names given to memory locations while objects are instances of a class with attributes and methods. Variables are used to store values while objects are used to represent real-world entities. Variables can be reassigned to different values while objects have a fixed identity. Variables are created when they are assigned a value while objects are created using constructors. Variables can be of different data types while objects are instances of a specific class. Example: x = 5 is a variable while car = Car() is an object. Example: x can be reassigned to a different value like x = 'hello' while car's identity remains the same.

# Have you heard about pickling and non pickling?

Pickling is a way to serialize Python objects, while non-pickling refers to objects that cannot be pickled. Pickling is used to convert Python objects into a byte stream that can be stored or transmitted. Non-pickling objects are those that cannot be serialized using the pickle module. Examples of non-pickling objects include file objects, network sockets, and database connections. To make an object picklable, it must be able to be reduced to a string of bytes and then reconstructed from those bytes. The pickle module is part of the Python standard library and can be used to serialize and deserialize objects.

[Python Interview Questions - CodinGame for Work](#)

9+0

I'm Satheesh from Bangalore, India, currently working as a Python Developer at Tata Consultancy Services with 3. Years of experience. My expertise spans data analysis, machine learning, and Django development. My work involves developing and maintaining DIET software that extracts essential information from documents. I excel in extracting insights from data, ensuring model accuracy, and visualizing data effectively with Matplotlib and Seaborn. Apart from these, i also manage servers for deployment and testing our application.Awarded on the spot awad last year and got appreciation from almarya client.
In terms of education i hold

In conclusion, I'm genuinely excited about the opportunity to work here.