

Course Two

Get Started with Python



Instructions

Use this PACE strategy document to record decisions and reflections as you work through this end-of-course project. You can use this document as a guide to consider your responses and reflections at different stages of the data analytical process. Additionally, the PACE strategy documents can be used as a resource when working on future projects.

Course Project Recap

Regardless of which track you have chosen to complete, your goals for this project are:

- Complete the questions in the Course 2 PACE strategy document
- Answer the questions in the Jupyter notebook project file
- Complete coding prep work on project's Jupyter notebook
- Summarize the column Dtypes
- Communicate important findings in the form of an executive summary

Relevant Interview Questions

Completing the end-of-course project will help you respond these types of questions that are often asked during the interview process:

- Describe the steps you would take to clean and transform an unstructured data set.

To clean and transform an unstructured data set I would follow these steps:

- **Data Understanding:**
 - **Assess Data Sources:** Identify the sources of the unstructured data (e.g., text files, social media feeds, emails).
 - **Explore Data:** Get a sense of the data by examining its structure, content, and quality.
- **Data Cleaning:**

- **Remove Noise:** Eliminate irrelevant data such as stop words, punctuation, and special characters.
- **Standardize Data:** Convert data to a consistent format (e.g., lowercasing text, removing duplicates).
- **Handle Missing Values:** Address missing data by filling in, removing, or estimating values.
- **Correct Errors:** Identify and correct spelling mistakes, formatting issues, and other inaccuracies.
- **Data Transformation:**
 - **Tokenization:** Break down text data into smaller units, such as words or sentences.
 - **Normalization:** Convert data to a standard format, such as stemming or lemmatization for text data.
 - **Feature Extraction:** Derive meaningful features from the data (e.g., keywords, sentiment scores, named entities).
 - **Vectorization:** Convert text data into numerical format using techniques like TF-IDF, word embeddings, or one-hot encoding.
- **Data Integration:**
 - **Merge Data Sources:** Combine multiple data sources to create a unified dataset.
 - **Align Data:** Ensure that data from different sources is aligned and consistent.
- **Data Validation:**
 - **Validate Transformations:** Check the transformed data for accuracy and consistency.
 - **Perform Quality Checks:** Conduct quality checks to ensure data integrity and reliability.
- **Documentation:**
 - **Document Steps:** Keep detailed documentation of the cleaning and transformation process.
 - **Create Metadata:** Generate metadata to describe the cleaned and transformed data.



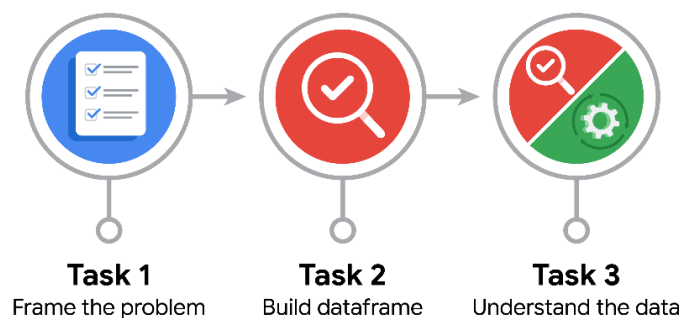
- What specific things might you look for as part of your cleaning process?
 - Locate the relevant information in the context of the analysis
 - Convert it into a consistent format:
 - ❖ identify missing data, duplicated entries, inconsistent data
 - ❖ identify noise, irrelevant data, errors and typos
 - Define the Data Types, normalize data, and check for integrity
 - Handle Categorical Data by encoding it into numerical values if required
 - Observe and understand the outlier values
- What are some of the outliers, anomalies, or unusual things you might look for in the data cleaning process that might impact analyses or ability to create insights?
 - ❖ Extreme Values:
 - Data points that are significantly higher or lower than the rest of the data. These can skew statistical analyses and models.
 - ❖ Inconsistent Data Entries:
 - Records that don't align with the expected format or values. For example, a numerical value where a categorical value is expected.
 - ❖ Unexpected Null Values:
 - Missing data in fields where you would normally expect complete information. This can indicate a problem with data collection or entry.
 - ❖ Temporal Anomalies:
 - Dates and times that don't make sense, such as future dates in historical data or timestamps that don't follow a logical sequence.
 - ❖ Duplicate Records:
 - Multiple entries for the same entity that can lead to overrepresentation in the data set.
 - ❖ Category Imbalance:
 - Categories within a variable that have disproportionately few or many records, which can affect the performance of certain models.
 - ❖ Irregular Patterns:
 - Unusual patterns or sequences that don't fit the expected trend. For example, a sudden spike in sales data without a clear reason.
 - ❖ Anomalous Relationships:



- Data points that don't follow the established relationships between variables. For example, a negative age value or a gender value that doesn't align with typical demographic data.
- ❖ Outlier Ratios:
 - Ratios or calculated values that are significantly different from the majority. For example, an extremely high or low debt-to-income ratio in financial data.
- ❖ Uncommon Values:
 - Rare values within a categorical variable that can indicate data entry errors or special cases that need to be handled separately.
- ❖ Geospatial Anomalies:
 - Location data that doesn't make sense, such as coordinates that fall outside the expected geographical area.
- ❖ Sensor or Instrument Errors:
 - In sensor data, look for readings that indicate malfunctioning or calibration issues.

Reference Guide

This project has three tasks; the visual below identifies how the stages of PACE are incorporated across those tasks.



Data Project Questions & Considerations



PACE: Plan Stage



- How can you best prepare to understand and organize the provided information?

By reviewing the column name and column description information that has been provided by our customer, I can get an idea on the type of data that I will be working with

- What follow-along and self-review codebooks will help you perform this work?

Panda fundamentals

DataFrames

Masking and Grouping

- What are some additional activities a resourceful learner would perform before starting to code?

Review the initial documentation, understand the goal of the project

Gather the function resources, conduct research on Python functions and classes that will be useful or necessary for the project

Use Gen AI tools to prompt for recommendations on how to tackle the problem, compare against my original idea, leverage from this information



PACE: **Analyze Stage**

- Will the available information be sufficient to achieve the goal based on your intuition and the analysis of the variables?

The number of records (408,294) with the details of the trip provided in the (18) fields will allow us to create the regression model.

- How would you build summary dataframe statistics and assess the min and max range of the data?

I would use the pandas function `describe()` which acts on dataframes and provides

Count: The number of non-null values.

Mean: The average of the values.

Std: The standard deviation of the values.

Min: The minimum value.

25%: The 25th percentile.

50% (Median): The 50th percentile or median.

75%: The 75th percentile.

Max: The maximum value.

- Do the averages of any of the data variables look unusual? Can you describe the interval data?

These are the average values for the variables I consider important, results don't look unusual to me:

passenger_count: 1.65

trip_distance: 2.9 miles

total_amount: 16.31 \$



PACE: Construct Stage

Note: The Construct stage does not apply to this workflow. The PACE framework can be adapted to fit the specific requirements of any project.



PACE: Execute Stage

- Given your current knowledge of the data, what would you initially recommend to your manager to investigate further prior to performing exploratory data analysis?

Get an explanation on trips that have a trip distance of 0

There are outlier values for total amount, get an explanation for those

Why are there negative values for the total amount of the trip

- What data initially presents as containing anomalies?

Trip distance and Total Amount

- What additional types of data could strengthen this dataset?

PULocationID: TLC Taxi Zone in which the taximeter was engaged

DOLocationID: TLC Taxi Zone in which the taximeter was disengaged



Executive Summary:

- **A summary of your tasks**
- **Information regarding the results of your data variable assessment**
- **Identify recommended next steps in order to build a predictive model**