

Data Prep & Exploration Checklist

This is a non-exhaustive checklist that can help you structure your analysis on Jupyter Notebook.

Step 1: Import & Observe

<input type="checkbox"/> Import necessary libraries	<pre>import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns</pre>
<input type="checkbox"/> Set notebook options	<pre># show all dataframe columns pd.set_option('display.max_columns', None) # set matplotlib global settings eg. figsize plt.rcParams['figure.figsize'] = (8.0, 6.0)</pre>
<input type="checkbox"/> Import your data (set delimiter and index column if applicable)	<pre># import data df_full = pd.read_csv('data.csv', sep=';', index_col=0) # check the first few rows of your dataframe df_full.head()</pre>
<input type="checkbox"/> Check the dimensions of your dataframe	<pre>df_full.shape</pre>
<input type="checkbox"/> Checkpoint #1: Make a copy of your dataset and remove any unnecessary columns (you might come back to this later as you start analysing the columns)	<pre>df = df_full.drop(columns=['B', 'C']).copy()</pre>
<input type="checkbox"/> Get a first overview of your dataset (look out for data types, missing values, column names and row range)	<pre>df.info()</pre>

Step 2: Format & Validate

<input type="checkbox"/> Check the output of <code>df.info()</code> . Does each column have the data type you would expect? If not, reformat column	<pre># some useful functions (click for more on string functions) df['col'].str.split() df['col'].str.replace() df['col'].astype()</pre>
<input type="checkbox"/> Format date columns	<pre># click to read more on how to specify date formats pd.to_datetime(df['date_col'], format='%Y%m%d')</pre>
<input type="checkbox"/> Check for suspicious values / outliers (Eg. what is the count? is the max value too far away from 75th percentile? Are there any negative values where there shouldn't be?)	<pre>df.describe()</pre>
<input type="checkbox"/> Create subsets for categorical / numerical variables to make your analysis easier	<pre># categorical variables df_cat = df.select_dtypes(include='object').copy() # numerical variables df_num = df.select_dtypes(include=['int','float']).copy()</pre>
<input type="checkbox"/> Check how many unique values are in each categorical column	<pre># count number of unique values per column df_cat.nunique()</pre>
<input type="checkbox"/> For variables with a small number of unique values, check value counts for consistency (eg. ideally you don't want a gender column to have values "M", "F", "Male", "Female")	<pre># check actual count df["categorical_col"].value_counts() # check proportion df["categorical_col"].value_counts(normalize=True) * 100</pre>

Step 3: Check for duplication

<input type="checkbox"/> Check for duplicated rows	<pre># check how many rows are duplicated df.duplicated().sum() # observe duplicated rows (understand why duplication occurs) df[df.duplicated(keep=False)]</pre>
<input type="checkbox"/> Check for duplicated columns (sometimes columns with different names might have the same values)	<pre>df.T.duplicated() # if data is too large you can just check the first few rows df.head(100).T.duplicated()</pre>
<input type="checkbox"/> Check that ID column does not contain duplicates (and if it does, understand why)	<pre>df['ID'].nunique() / len(df)</pre>
<input type="checkbox"/> Decide on duplication handling strategy (keep vs remove)	<pre># drop duplicates and keep the first occurrence df.drop_duplicates(keep='First')</pre>

Step 4: Check for missing values

<input type="checkbox"/> Check for missing data	<pre># to check number of missing values df.isna().sum() # to check proportion of missing values df.isna().mean() * 100</pre>
<input type="checkbox"/> Visualize missing data - Check if missing values coincide across columns - Sort by date and plot, do missing values occure before/after a certain date?	<pre>sns.heatmap(df.isnull(), yticklabels=False, cbar=False)</pre>
<input type="checkbox"/> Decide on missing data handling strategy (imputation vs deletion)	<pre># to drop missing values df.dropna() # to impute missing values eg. using mean df.filna(df[['co'l.mean()]])</pre>
<input type="checkbox"/> Checkpoint #2: Reset index, make a copy of your final clean dataframe and save it as a csv	<pre># create a copy to use for the rest of your analysis df_final = df.reset_index(drop=True).copy() # save a csv so that you can quickly import it next time df_final.to_csv('data_final.csv')</pre>

Step 5: Visualize data

<input type="checkbox"/> Univariate visualizations	<pre># categorical: which category has the most/least observations? sns.countplot(df, x="categorical_col") # numerical: what distribution do you observe? sns.histplot(df, x="numerical_col") sns.boxplot(df, x="numerical_col")</pre>
<input type="checkbox"/> Multivariate visualizations	<pre># create scatterplots for all numerical variable combinations sns.pairplot(df_num) # visualize correlations between numerical columns sns.heatmap(df_num.corr()) # check how distribution varies by different categories sns.violinplot(data=df, x="cat1", y="num", hue="cat2") sns.boxplot(data=df, x="cat1", y="num", hue="cat2")</pre>
<input type="checkbox"/> Save your visualizations	<pre># make sure to give your file a meaningful name plt.savefig('sales_by_quarter.png', dpi=200)</pre>

Step 6: Clean up notebook

<input type="checkbox"/> Remove any redundant / unused cells (or put them in an appendix)
<input type="checkbox"/> Make sure your notebook is well structured with clear sections defined by headings and subheadings
<input type="checkbox"/> Add a notebook introduction/executive summary at the beginning
<input type="checkbox"/> Make sure you have communicated your insights clearly (don't leave any plots or findings without explanations)
<input type="checkbox"/> Finally, reset your kernel and check that your notebook runs from top to bottom with no errors