# Python Data Analysis Code Cheatsheet

This cheatsheet is your post-program reference toolkit, this generic template is designed to help you quickly recall and apply the key code patterns you've learned throughout your data analysis journey.

## How to use it:

**1. During the program:** Use this as your personal playbook. When working on projects, solving problems, revisiting concepts later, come back to this sheet to copy, tweak, and run the code you need.

**2. When in doubt:** Think of this as your first stop before Googling or digging through notes.

Use this cheatsheet as your go-to companion whenever you need to work with data in Python. You've done the learning - now here's the code to get it done.

| Data Preparation Process | Workflow | Code |
|---|---|---|
| Getting Started | Import Python Modules | ```# Import the required Python Modules

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns``` |
| | loading Dataset | ```# Load the dataset
df = pd.read_csv("DATASETNAME.csv")

# Note: DF is the name of the data frame. If you are working with multiple data frames, you will need to use the appropriate data frame name based on how you defined them.``` |
| Exploring the data | Print Table | ```# Display the first 5 rows``` |

| frame | | ```df.head()```<br><br>```# Display the first 5 rows of specific COLUMNS. Note you can apply this concept of selecting specific columns for subsequent codes.```<br>```df[['COLUMN1', 'COLUMN2']].head()``` |
|---|---|---|
| | Row and column | ```# Check the shape of the dataset```<br>```print(f"Dataset contains {df.shape[0]} rows and {df.shape[1]} columns.")```<br><br>```# Note: [0] refers to rows, while [1] refers to columns.``` |
| | Identify missing values - using.info() | ```df.info()``` |
| Data Cleaning | Identify missing values - using isnull() | ```# Check for missing values```<br>```print("Missing Values Count:\n", df.isnull().sum())```<br><br>```# Display rows with missing values```<br>```df[df.isnull().any(axis=1)]``` |
| | Handling missing values - Removal method | ```#Remove rows with missing values in the COLUMN. Save this as a new variable called df2```<br>```df2 = df.dropna(subset=['COLUMN'])```<br><br>```#Verify if the process was completed by checking the number of remaining missing values```<br>```print("Missing values after cleaning:")```<br>```print(df2.isnull().sum())```<br><br>```#note that it is best practice to create new data frames when we make``` |

| | | |
|---|---|---|
| | | any changes so that we can go back to it |
| | Handling missing values - Imputation method | ```python<br># Fill missing values in the 'Price' column with the mean. Save this as a new variable called df2<br>df2['COLUMN'] = df['COLUMN'].fillna(df['COLUMN'].mean())<br><br># Note that you can change to median by changing .mean to .median<br><br># Verify the missing values are handled<br>print("Missing values after performing imputation:")<br>print(df2['Price'].isnull().sum())<br>``` |
| | Handling duplicates in data frames | ```python<br>#Check for duplicates<br>duplicates = df_cleaned2.duplicated().sum()<br>print(f"Number of duplicate rows: {duplicates}")<br><br>#Display duplicate rows if any exist<br>if duplicates > 0:<br>    print("Duplicate rows found:")<br>    display(df_cleaned2[df_cleaned2.duplicated(keep=False)])  # Show all duplicate rows (including first occurrence)<br><br># Remove duplicates if necessary<br>df_cleaned2 = df_cleaned2.drop_duplicates()<br><br># Verify removal<br>print(f"Number of duplicate rows after removal: {df_cleaned2.duplicated().sum()}")<br>``` |
| | Export the cleaned datasets for secondary storage | ```python<br># Export data frame into CSV<br>df2.to_csv('DATASET NAME.csv', index=False)<br>``` |
| Data Manipulation | Creating New Columns | ```python<br># Creating new columns in the table based on another column<br>df['NEW_COLUMN'] = df['OLD_COLUMN']<br>``` |

| | | |
|---|---|---|
| | | ```# Note: You can personalise the formulas accordingly. Ie: If we want to apply a 10% discount, you can multiply the OLD COLUMN by 0.9``` |
| | Encoding Categorical Data into Numerical Data | ```python<br># Define your custom mapping for categories<br>category_mapping = {<br>    'CATEGORY 1': 0,<br>    'CATEGORY 2': 1<br>}<br><br># Apply the mapping<br>df['NEW_COLUMN'] = df['OLD_COLUMN'].map(category_mapping)``` |
| | Extracting from datetime columns | ```python<br># Convert the 'DATE_COLUMN' into datetime format<br>df['DATE_COLUMN'] = pd.to_datetime(df['DATE_COLUMN'])<br><br># Create new columns that extracts the YEAR, MONTH, DAY, and<br>DAY_OF_WEEK from the DATE_COLUMN<br>df['YEAR'] = df['DATE_COLUMN'].dt.year<br>df['MONTH'] = df['DATE_COLUMN'].dt.month<br>df['DAY'] = df['DATE_COLUMN'].dt.day<br>df['DAY_OF_WEEK'] = df['DATE_COLUMN'].dt.day_name()<br><br># Display the transformed dataset<br>df[['DATE_COLUMN', 'YEAR', 'MONTH', 'DAY', 'DAY_OF_WEEK']].head()``` |
| Data Transformation | Converting data types in a new column | ```python<br># Create a new column and changes it into int data type.<br>df['COLUMN_INT'] = df['COLUMN'].astype(int)<br><br># Note: You can convert to other data types by changing 'int' with<br>others (ie: float, str)``` |
| | Convert date and time to string format | ```python<br>df['Date_of_Visit'] = pd.to_datetime(df['Date_of_Visit'],<br>errors='coerce', dayfirst=True).dt.strftime('%d-%m-%Y')``` |

| | | |
|---|---|---|
| | | ```df.head()``` |
| | Ensure consistent formatting of date & time | ```# Convert 'DATE_COLUMN' into datetime format and displays it as a DD-MM-YYYY format df['DATE_COLUMN'] = pd.to_datetime(df['DATE_COLUMN'], errors='coerce').dt.strftime('%d-%m-%Y')``` <br><br> ```# Note: You can personalise it based on the format of date that you want to display.``` |
| | Ensure consistent formatting of categorical data | ```# Standardise categorical data in a column with lower case df[COLUMN] = df['COLUMN'].str.lower().str.strip()``` <br><br> ```# Verify unique values after data transformation df['Weather'].unique()``` |
| | Information Range checks | ```# Replace null or negative prices with median price df['COLUMN'] = df['COLUMN'].apply(lambda x: x if pd.notnull(x) and x > 0 else df['COLUMN'].median())``` <br><br> ```# Verify if any prices are still missing or invalid df[df['COLUMN'] <= 0]``` <br><br> ```# Note: You can replace the highlighted parts with any range of your choice, and replace median with your preferred choice of handling data that is out of range.``` |
| Data Joining | GroupBy - on Numerical Data | ```# Standardise a categorical column # Replace <column_name> with the actual column header you want to clean df['<column_name>'] = df['<column_name>'].str.lower().str.strip()``` |

| | | ```
# Quick check
print(df['<column_name>'].unique())
``` |
|---|---|---|
| | Merge | ```
# Merge two DataFrames
# Replace <left_df>, <right_df>, and <join_key> with your actual names
merged_df = <left_df>.merge(<right_df>, on="<join_key>", how="left")

# Inspect the result
merged_df.head()
``` |
| Exploring Data Analysis (EDA) | Summarise Numerical variable- describe ( ) | ```
# STEP 1: Drop identifier or helper columns that you don't want in the summary
numerical_data = <df>.drop(columns=<id_columns>)

# STEP 2: Keep only numeric columns
numerical_columns = numerical_data.select_dtypes(include="number")

# STEP 3: Generate descriptive statistics
numerical_summary = numerical_columns.describe()

print(numerical_summary)
``` |
| | Explore categorical variable | ```
# STEP 1 ▶ Identify all categorical columns (dtype "object")
categorical_columns = <df>.select_dtypes(include="object").columns

print("\nFrequency distribution for categorical variables:")

# STEP 2 ▶ Loop through each categorical column
for col in categorical_columns:

    # STEP 3 ▶ Display absolute frequencies for that column
    print(f"\nColumn: {col}")
    print(<df>[col].value_counts())

    # OPTIONAL ▶ Uncomment to display relative frequencies
``` |

| | | |
|---|---|---|
| | | ```<br>(percentages)<br>    # print(<df>[col].value_counts(normalize=True))<br>``` |
| | Display the frequency of categories | ```python<br>import pandas as pd<br><br># STEP 1 ▶ Identify categorical columns (dtype "object")<br>categorical_columns = <DF>.select_dtypes(include="object").columns<br><br>print("\nFrequency distribution for categorical variables:")<br><br># STEP 2 ▶ Loop through each categorical column and print the counts<br>for col in categorical_columns:<br>    print(f"\nColumn: {col}")<br>    print(<DF>[col].value_counts())          # absolute frequencies<br><br>    # OPTIONAL ▶ Uncomment to show relative frequencies (percentages)<br>    # print(<DF>[col].value_counts(normalize=True))<br>``` |
| Data Visualisation | Visualising Trends : Line plot | ```python<br># STEP 1 ▶ Group data by <group_col> to calculate total visits (row count)<br>visits_by_category = (<br>    <DF><br>        .groupby('<group_col>')<br>        .size()                               # counts rows in each group<br>        .reset_index(name='<count_col>')      # renames the count column<br>)<br><br># OPTIONAL ▶ Sort categories in a custom order (e.g. calendar order)<br>custom_order = ['Jan', 'Feb', 'Mar', 'Apr', 'May']   # edit to your needs<br>visits_by_category['<group_col>'] = pd.Categorical(<br>    visits_by_category['<group_col>'],<br>    categories=custom_order,<br>    ordered=True<br>``` |

```
)
visits_by_category = visits_by_category.sort_values('<group_col>')

# STEP 2 ▶ Plot the trend with a line chart
plt.figure(figsize=(10, 5))
plt.plot(
    visits_by_category['<group_col>'],
    visits_by_category['<count_col>'],
    marker='o',
    linestyle='-'              # solid line
    # color parameter is optional; Matplotlib uses a default palette
)

plt.xlabel("<group_col>")
plt.ylabel("<count_col>")
plt.title(f"Total {<count_col>} by {<group_col>}")
plt.xticks(rotation=45)
plt.grid(True)
plt.tight_layout()
plt.show()
```

| | | |
|---|---|---|
| | Visualising comparison: Bar Chart | ```
# STEP 1 ▶ Group by <group_col> to calculate the average of <agg_col>
df_grouped = (
    <DF>
        .groupby('<group_col>')[ '<agg_col>' ]
        .mean()
        .reset_index()
)

# STEP 2 ▶ Sort the result in descending order of the average value
df_grouped_sorted = df_grouped.sort_values(by='<agg_col>',
ascending=False)

# STEP 3 ▶ Plot the bar chart
plt.figure(figsize=(8, 5))
``` |

| | | ```
plt.bar(
    df_grouped_sorted['<group_col>'],
    df_grouped_sorted['<agg_col>'],
)

# STEP 4 ▶ Customise the chart
plt.xlabel('<group_col>')
plt.ylabel(f'Average <agg_col>')
plt.title(f'Average <agg_col> by <group_col>')
plt.xticks(rotation=45)
plt.tight_layout()

# STEP 5 ▶ Display the plot
plt.show()
``` |
|---|---|---|
| | Visualising distribution: Histogram | ```
# Step 1: Create histogram
plt.figure(figsize=(8, 5))
sns.histplot(<DF>['<numeric_col>'], bins=20, kde=True, color='green')

# Step 2: Customize the plot
plt.xlabel("<numeric_col>")
plt.ylabel("Frequency")
plt.title("Distribution of <numeric_col>")

# Step 3: Show the plot
plt.tight_layout()
plt.show()
``` |
| | Visualising composition: Pie Chart | ```
# Step 1: Group by <group_col> and sum <value_col>
grouped_data = (
    <DF>
        .groupby('<group_col>')[ '<value_col>' ]
        .sum()
        .reset_index()
)
``` |

```python
# Step 2: Create a pie chart
plt.figure(figsize=(8, 8))
plt.pie(
    grouped_data['<value_col>'],
    labels=grouped_data['<group_col>'],
    autopct='%1.2f%%',
    startangle=140
)

# Step 3: Customize the chart
plt.title(f"Total {<value_col>} by {<group_col>}")
plt.tight_layout()

# Step 4: Display the chart
plt.show()
```