

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5
6

```

```

1 pip install pandas numpy scipy scikit-learn matplotlib seaborn
2

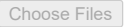
```

Requirement already satisfied: pandas in /usr/local/lib/python3.11/dist-packages (2.2.2)
 Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (1.26.4)
 Requirement already satisfied: scipy in /usr/local/lib/python3.11/dist-packages (1.14.1)
 Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (1.6.1)
 Requirement already satisfied: matplotlib in /usr/local/lib/python3.11/dist-packages (3.10.0)
 Requirement already satisfied: seaborn in /usr/local/lib/python3.11/dist-packages (0.13.2)
 Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas) (2.8.2)
 Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.1)
 Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.1)
 Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.4.2)
 Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.5.0)
 Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.3.1)
 Requirement already satisfied: cyclor>=0.10 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (0.12.1)
 Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (4.56.0)
 Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.4.8)
 Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (24.2)
 Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (11.1.0)
 Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (3.2.1)
 Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas) (1.17.0)

```

1 from google.colab import files
2 uploaded = files.upload()

```

 No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable

```

1 data = pd.read_csv('HistoricalData_1741898709939.csv')
2 data.head()

```

	Date	Close/Last	Volume	Open	High	Low
0	03/12/2025	\$216.98	62547470	\$220.14	\$221.75	\$214.91
1	03/11/2025	\$220.84	76137410	\$223.805	\$225.8399	\$217.45
2	03/10/2025	\$227.48	72071200	\$235.54	\$236.16	\$224.22
3	03/07/2025	\$239.07	46273570	\$235.105	\$241.37	\$234.76
4	03/06/2025	\$235.33	45170420	\$234.435	\$237.86	\$233.1581

```

1 missing_values = data.isnull().sum()
2 print("Missing values per column:")
3 print(missing_values)
4
5 duplicate_count = data.duplicated().sum()
6 print(f"Duplicate rows: {duplicate_count}")
7 # Optionally remove duplicates
8 data = data.drop_duplicates()
9
10 summary = data.describe()
11 print("Statistical summary:")
12 print(summary)
13
14

```

Missing values per column:

Date	0
Close/Last	0
Volume	0
Open	0
High	0
Low	0

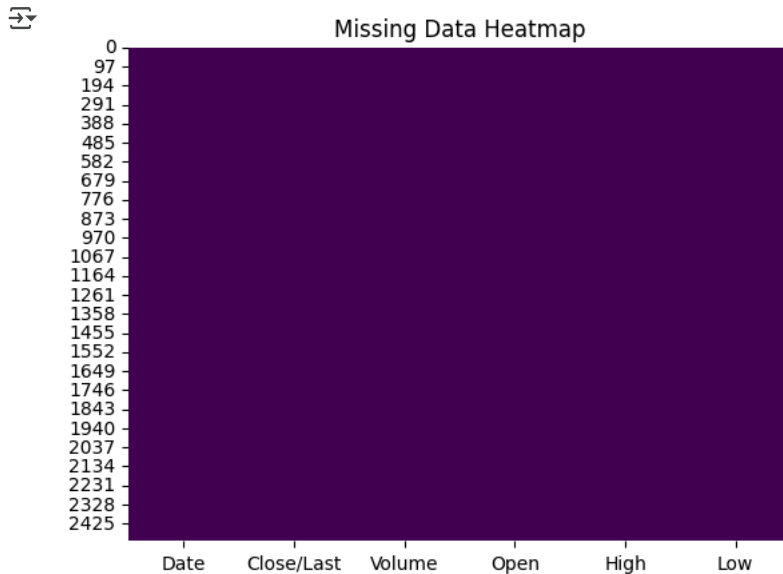
dtype: int64
 Duplicate rows: 0
 Statistical summary:

	Volume
count	2.515000e+03
mean	1.131585e+08
std	6.536536e+07
min	2.323471e+07
25%	6.927758e+07
50%	9.757607e+07
75%	1.376138e+08
max	6.475300e+08

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 sns.heatmap(data.isnull(), cbar=False, cmap='viridis')
5 plt.title("Missing Data Heatmap")
6 plt.show()
7

```



```

1 # Remove rows with missing values (for bootstrapping)
2 clean_data = data.dropna()
3
4 # Resample clean_data to create a synthetic dataset with the same number of rows
5 synthetic_data = clean_data.sample(n=len(clean_data), replace=True)
6
7 # Check that the synthetic dataset has similar statistical properties
8 print(synthetic_data.describe())
9

```

	Volume
count	2.515000e+03
mean	1.124608e+08
std	6.605131e+07
min	2.726298e+07
25%	6.905920e+07
50%	9.604190e+07
75%	1.348883e+08
max	6.475300e+08

```

1 report_content = f"""
2 Data Quality Report:
3
4 Missing Values per Column:
5 {missing_values.to_string()}
6
7 Duplicate Rows: {duplicate_count}
8
9 Statistical Summary:
10 {summary.to_string()}
11 """
12
13 with open('data_quality_report.txt', 'w') as file:
14     file.write(report_content)
15

```

```
16 print("Report generated as 'data_quality_report.txt'.")
17
```

Report generated as 'data_quality_report.txt'.

```
1 file_path = pd.read_csv('HistoricalData_1741898709939.csv')
2 class DataQualityTool:
3     def __init__(self, file_path ):
4         self.data = file_path
5
6     def check_missing(self):
7         return self.data.isnull().sum()
8
9     def check_duplicates(self):
10        return self.data.duplicated().sum()
11
12    def get_summary(self):
13        return self.data.describe()
14
15    def bootstrap_data(self):
16        clean_data = self.data.dropna()
17        return clean_data.sample(n=len(clean_data), replace=True)
18
19    def generate_report(self, report_filename='report.txt'):
20        missing = self.check_missing()
21        duplicates = self.check_duplicates()
22        summary = self.get_summary()
23        synthetic = self.bootstrap_data()
24
25        report = f"Missing Values:\n{missing}\n\nDuplicates: {duplicates}\n\nSummary:\n{summary}\n\nSynthetic Data Summary:\n{synthetic.
26
27        with open(report_filename, 'w') as f:
28            f.write(report)
29        print(f"Report generated as {report_filename}")
30
31 # Usage example:
32 dq_tool = DataQualityTool('your_dataset.csv')
33 dq_tool.generate_report()
34
```

AttributeError Traceback (most recent call last)
 <ipython-input-15-5923a7cf310b> in <cell line: 0>()
 31 # Usage example:
 32 dq_tool = DataQualityTool('your_dataset.csv')
 ----> 33 dq_tool.generate_report()

----- 1 frames -----
 <ipython-input-15-5923a7cf310b> in check_missing(self)
 5
 6 def check_missing(self):
 ----> 7 return self.data.isnull().sum()
 8
 9 def check_duplicates(self):


AttributeError: 'str' object has no attribute 'isnull'

```
1 import pandas as pd
2 import numpy as np
3 from scipy import stats
4
5 class DataQualityTool:
6     def __init__(self, file_path):
7         """
8         Constructor: Loads data from a CSV file.
9
10        Parameters:
11        - file_path (str): The path to the CSV file containing your dataset.
12        """
13        # Load the dataset from the given file path and store it in an instance variable.
14        self.data = pd.read_csv('HistoricalData_1741898709939.csv')
15
16    def check_missing(self):
17        """
18        Check for missing values in the dataset.
19
```

```

20     Returns:
21     - A pandas Series showing the count of missing values in each column.
22     """
23     return self.data.isnull().sum()
24
25     def check_duplicates(self):
26     """
27     Check for duplicate rows in the dataset.
28
29     Returns:
30     - An integer count of duplicate rows.
31     """
32     return self.data.duplicated().sum()
33
34     def get_summary(self):
35     """
36     Get basic statistical summary of the dataset.
37
38     Returns:
39     - A pandas DataFrame containing summary statistics.
40     """
41     return self.data.describe()
42
43     def bootstrap_data(self):
44     """
45     Generate a synthetic dataset by bootstrapping.
46     The process:
47     1. Drops rows with missing values.
48     2. Samples with replacement to generate a new dataset
49     with the same number of rows as the cleaned data.
50
51     Returns:
52     - A pandas DataFrame representing the bootstrapped dataset.
53     """
54     clean_data = self.data.dropna()
55     return clean_data.sample(n=len(clean_data), replace=True)
56
57     def generate_report(self, report_filename='report.txt'):
58     """
59     Generate a report that summarizes the data quality metrics.
60
61     Parameters:
62     - report_filename (str): The filename for saving the report.
63
64     The report includes:
65     - Missing values per column.
66     - Count of duplicate rows.
67     - Statistical summary of the dataset.
68     - Statistical summary of the bootstrapped (synthetic) dataset.
69     """
70     missing = self.check_missing()
71     duplicates = self.check_duplicates()
72     summary = self.get_summary()
73     synthetic = self.bootstrap_data()
74
75     # Create a report string using the computed metrics.
76     report = (
77         "Data Quality Report\n\n"
78         "Missing Values per Column:\n" + missing.to_string() + "\n\n"
79         f"Duplicate Rows: {duplicates}\n\n"
80         "Statistical Summary:\n" + summary.to_string() + "\n\n"
81         "Synthetic Data Summary:\n" + synthetic.describe().to_string()
82     )
83
84     # Write the report to a text file.
85     with open(report_filename, 'w') as f:
86         f.write(report)
87     print(f"Report generated as {report_filename}")
88
89 # Usage Example:
90 if __name__ == '__main__':
91     # Replace 'your_dataset.csv' with the actual path to your dataset file.
92     dq_tool = DataQualityTool('your_dataset.csv')
93     dq_tool.generate_report()
94

```

 Report generated as report.txt

```
1 with open('report.txt', 'r') as file:
2     content = file.read()
3     print(content)
4
5
```

Data Quality Report

Missing Values per Column:

Date	0
Close/Last	0
Volume	0
Open	0
High	0
Low	0

Duplicate Rows: 0

Statistical Summary:

	Volume
count	2.515000e+03
mean	1.131585e+08
std	6.536536e+07
min	2.323471e+07
25%	6.927758e+07
50%	9.757607e+07
75%	1.376138e+08
max	6.475300e+08

Synthetic Data Summary:

	Volume
count	2.515000e+03
mean	1.154672e+08
std	6.750152e+07
min	2.323471e+07
25%	7.062210e+07
50%	9.893191e+07
75%	1.414170e+08
max	6.475300e+08