# PANDAS PROFILING (PYTHON)

Python module for Exploratory Data Analysis (EDA))





**Links**

- ✓ https://pandas-profiling.github.io/pandas-profiling/docs/master/rtd/
- ✓ https://pypi.org/project/pandas-profiling/#modal-close
- ✓ https://www.kaggle.com/nulldata/intro-to-pandas-profiling-simple-fast-eda
- ✓ https://github.com/pandas-profiling/pandas-profiling

**License**      MIT
**Version**     2.6.0
**Last Update**  4/14/2020

**OS**   Linux, macOS, Windows

**System Requirements**   JRE 11, Python 3.8

**Description**   Pandas profiling is an open source Python module with which we can quickly do an exploratory data analysis with just a few lines of code.  It generates profile reports from a pandas DataFrame. The pandas df.describe() function is great but a little basic for serious exploratory data analysis. pandas profiling extends the pandas DataFrame with df.profile_report() for quick data analysis.

**Features**
- Type inference: detect the types of columns in a dataframe.
- Essentials: type, unique values, missing values
- Quantile statistics like minimum value, Q1, median, Q3, maximum, range, interquartile range
- Descriptive statistics like mean, mode, standard deviation, sum, median absolute deviation, coefficient of variation, kurtosis, skewness
- Most frequent values
- Histogram
- Correlations highlighting of highly correlated variables, Spearman, Pearson and Kendall matrices
- Missing values matrix, count, heatmap and dendrogram of missing values
- Text analysis learn about categories (Uppercase, Space), scripts (Latin, Cyrillic) and blocks (ASCII) of text data.
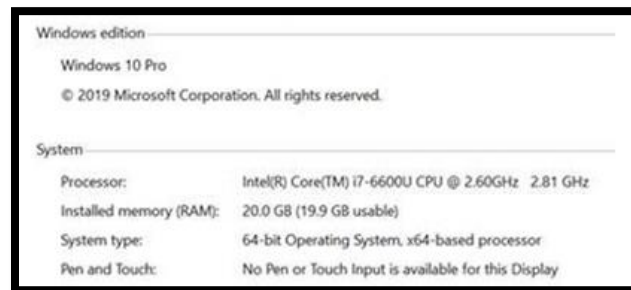
**Connectivity / Supported Data Sources & Formats**
- Text: - CSV, fixed-width test files, JSON, HTML, Clipboard, Excel
- Binary: OpenDocument, HDF5 Format, Feather Format, Parqeuet Format, ORC Format, Msgpak, Stata, SAS, SPSS, Python Pickle Format
- SQL, Google BigQuery

**Limitations**

With the increase in the size of the data the time to generate the report also increases a lot.  This problem can be solved by generating the report from only a part of the data or by using "minimum mode" introduced in version 2.4 for a simplified report. The tool requires people with technical knowledge.

**Performance**



Successfully able to load the 1.3M rows from CSV in < 5 seconds and generation of profiling report took ~1 min

## DATA PROFILING WITH PANDAS PROFILING (PYTHON)

### PYTHON:

1. Download Python:
    a. Download the latest version of Python 3 by running Python Installer from https://www.python.org/downloads/.
    b. Download Anaconda from "https://www.anaconda.com/products/individual".
2. Install pandas-profiling:
    a. the pip package manager by running "pip install pandas-profiling[notebook]"
    b. from Github: "pip install https://github.com/pandas-profiling/pandas-profiling/archive/master.zip"
    c. Using conda: conda install -c conda-forge pandas-profiling
3. The documentation for pandas_profiling can be found at https://pandas-profiling.github.io/pandas-profiling/docs/master/.
4. Run jupyter notebook or preferred Python IDE
5. Import required packages including pandas, pandas-profiling, ProfileReport as

    ```
    import pandas as pd

    import pyodbc

    import pandas_profiling

    import ProfileReport
    ```
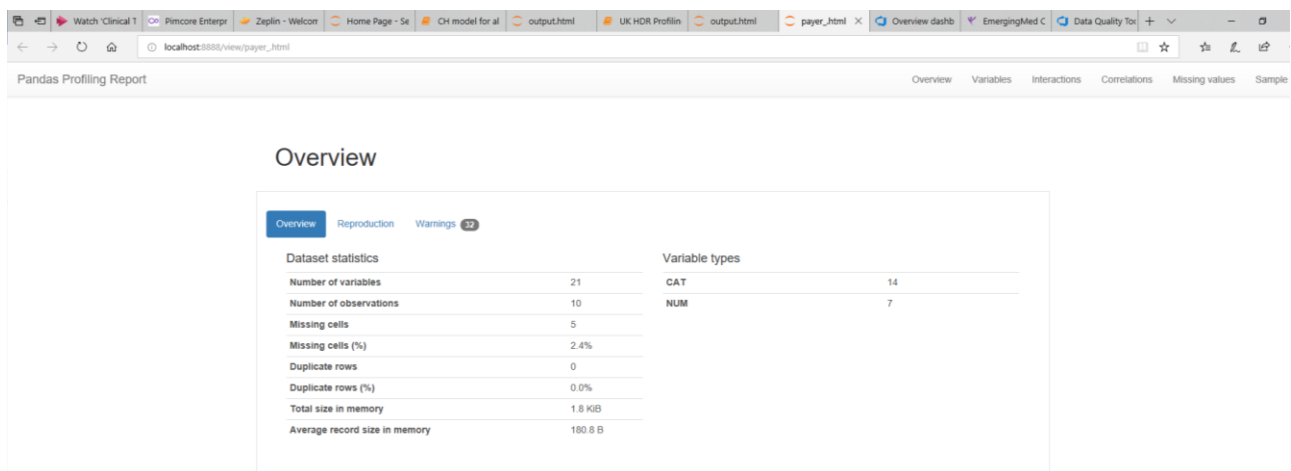
6. Connect to database by:
    ```
    pyodbc.connect('DRIVER={SQL Server}; SERVER=' + DB['servername'] + '; DATABASE=' + DB['database name'] + '; Trusted_Connection=yes')
    ```
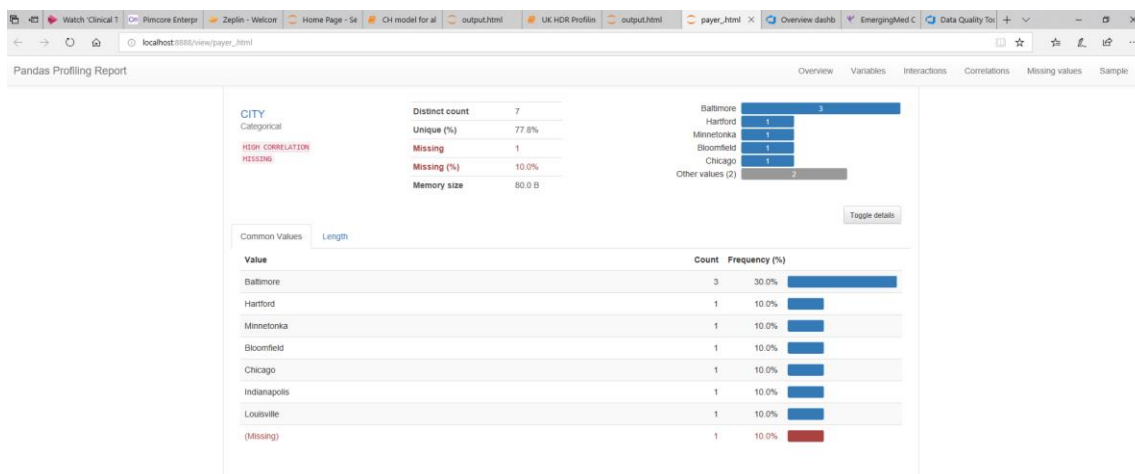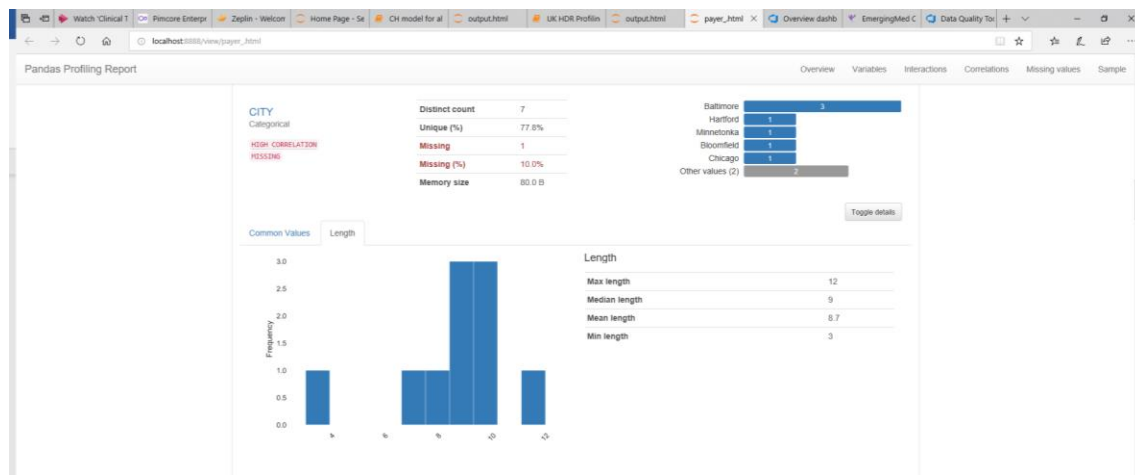
7. Import data from SQL database table:

    df = pd.read_sql_query(f"""select * from table_name_from_database """, conn)

8. To generate the report, run: pandas_profiling.ProfileReport(df)

9. Saving the report: report can be save as HTML or json by using to_file() function:

    a. save as HTML file: profile.to_file("your_report.html")
    b. save as JSON file: profile.to_file("your_report.json")

## PROFILING REPORT OVERVIEW: PROVIDES BASIC INFORMATION ABOUT DATA.
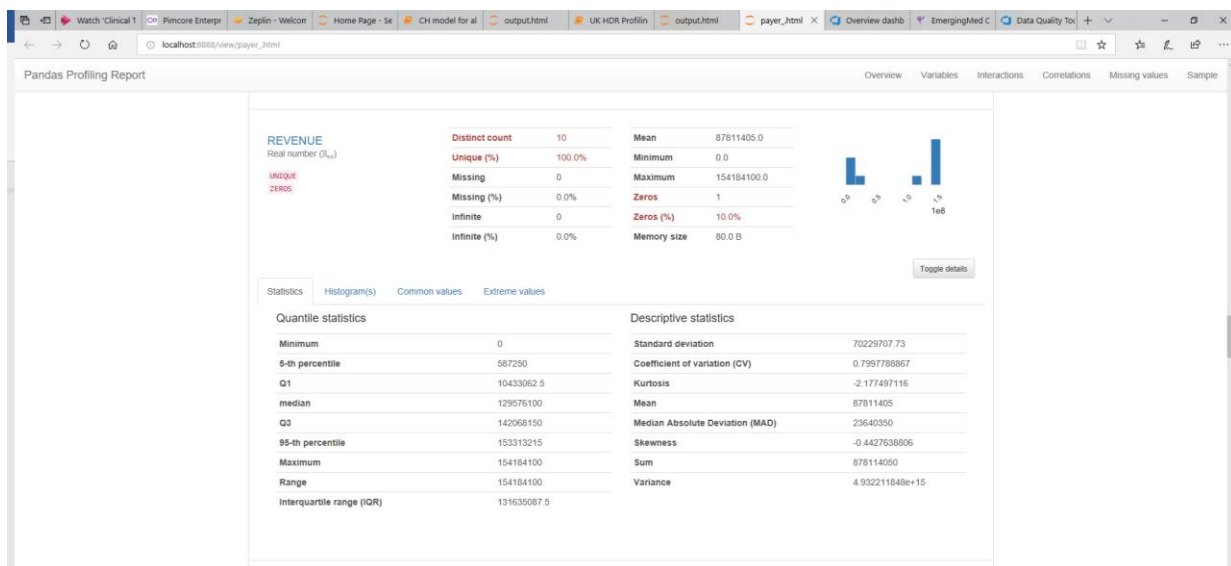


Profiling report categorical variables provides distinct values, unique (%), missing values, missing (%), histogram of count and frequency (%), etc.
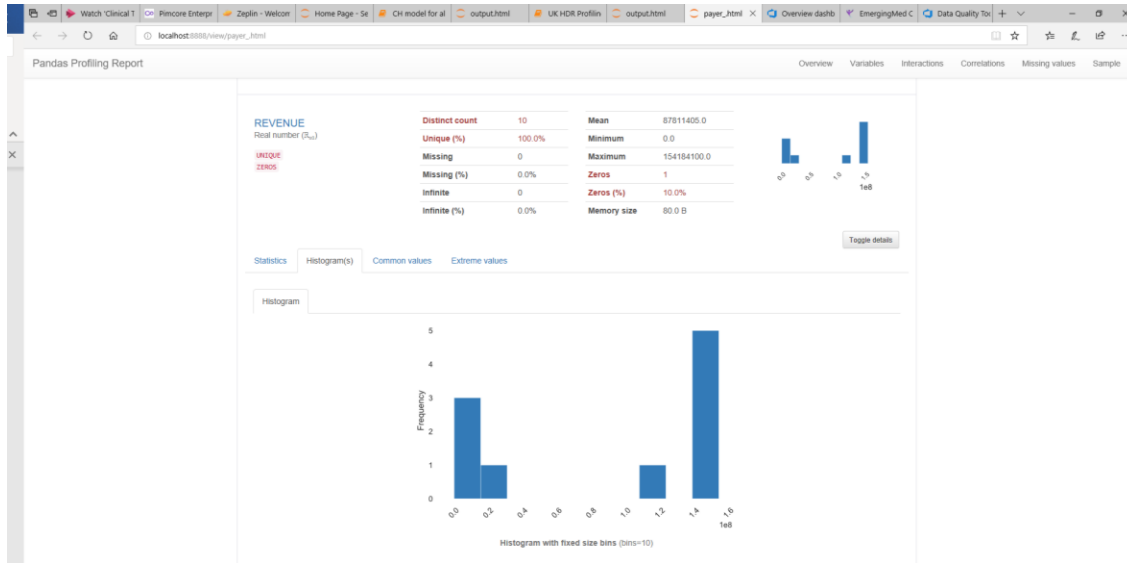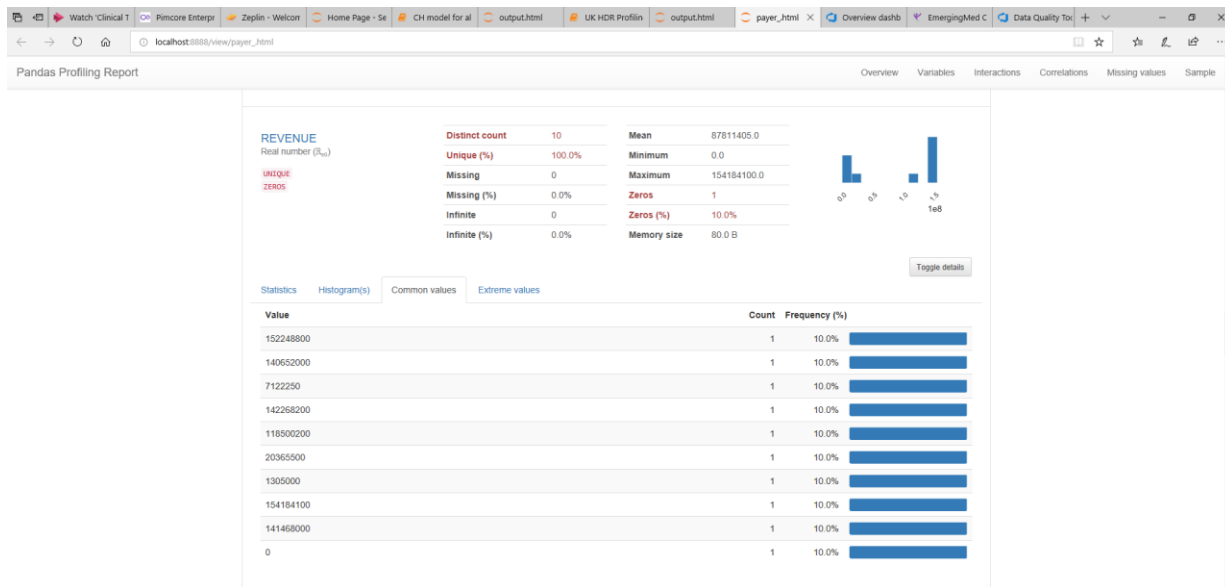
Profiling report numerical variables distinct values, unique (%), missing values, missing (%), histogram of count and frequency (%), mean, min, max, zeros, Quantile statistics (5th percentile, Q1, median, Q3, Max, Range, IQR), Quantile statistics (SD, CV, Kurtosis, Mean, MAD, Sum, Skewness, Variance) etc.
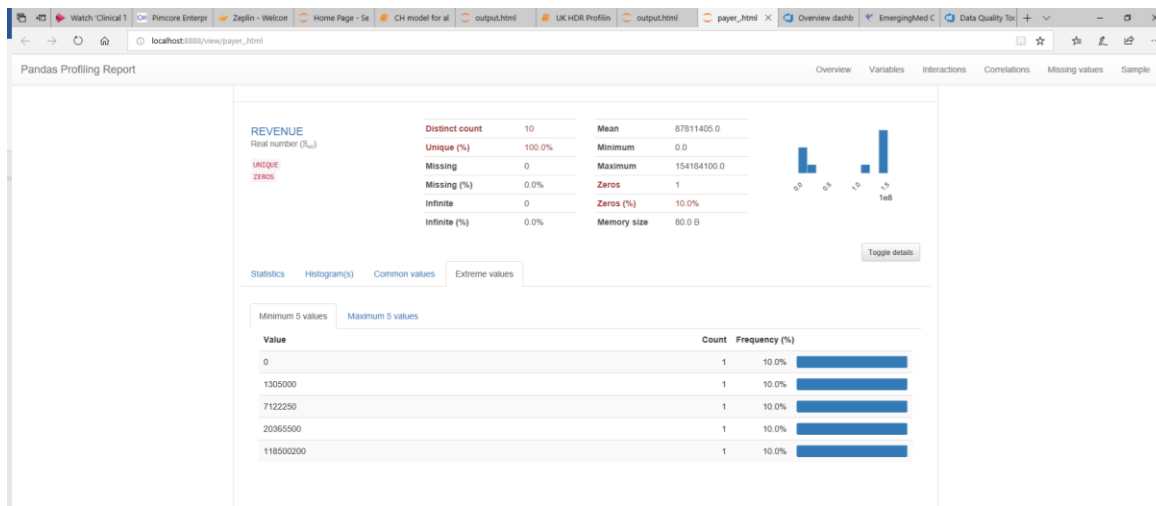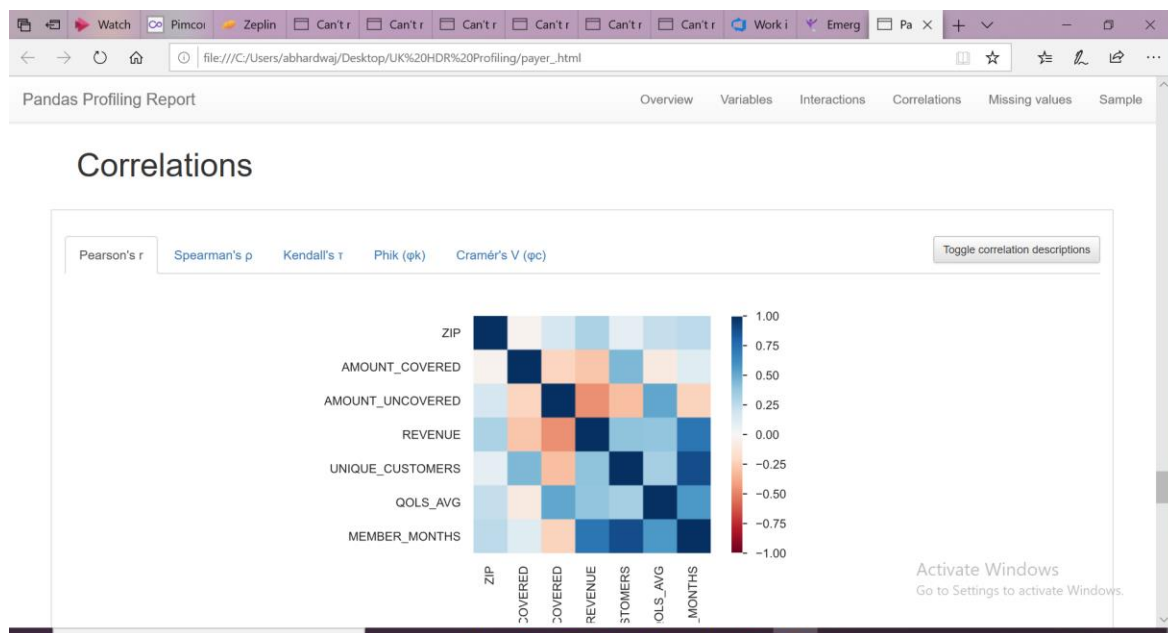
## FREQUENCY HISTOGRAM:



## COMMON VALUES COUNT AND FREQUENCY:

## EXTREME VALUES: COUNTS AND FREQUENCIES OF TOP 5 MIN AND MAX VALUES



## CORRELATIONS:

## INTERACTIONS BETWEEN VARIABLES: