

EDA Implementation using Automated EDA Tools

Submitted by: Nilutpal Das

```
In [13]: # Importing CSV dataset with a help of Pandas Library
import pandas as pd
df = pd.read_csv('D:\Data Science\Dataset\data5\Visadataset.csv')
```

```
In [14]: # head() used to provide top 5 data from dataset.
df.head()
```

```
Out[14]:
```

	case_id	continent	education_of_employee	has_job_experience	requires_job_training	no_of_employees	yr_of_estab	region_of_employment
0	EZVY01	Asia	High School	N	N	14513	2007	West
1	EZVY02	Asia	Master's	Y	N	2412	2002	Northeast
2	EZVY03	Asia	Bachelor's	N	Y	44444	2008	West
3	EZVY04	Asia	Bachelor's	N	N	98	1897	West
4	EZVY05	Africa	Master's	Y	N	1082	2005	South

1. Pandas Profiling Tools

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.

The **pandas_profiling** library generates a report having:

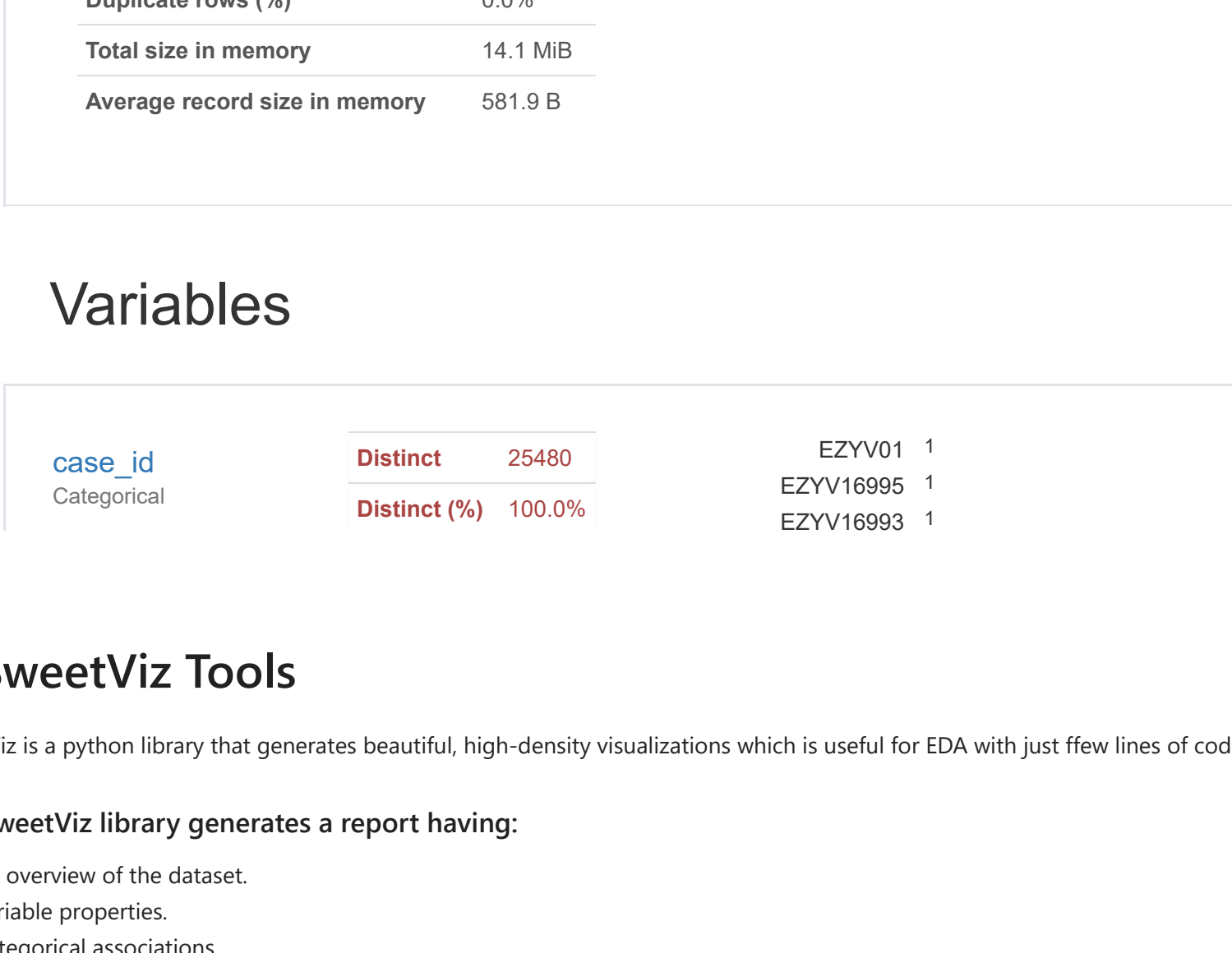
- An overview of the dataset.
- Variable properties.
- Interaction of variables.
- Correlation of variables.
- Sample data.
- Missing values.

```
In [15]: from pandas_profiling import ProfileReport
profile = ProfileReport(df, explorative=True, title='Report')
profile

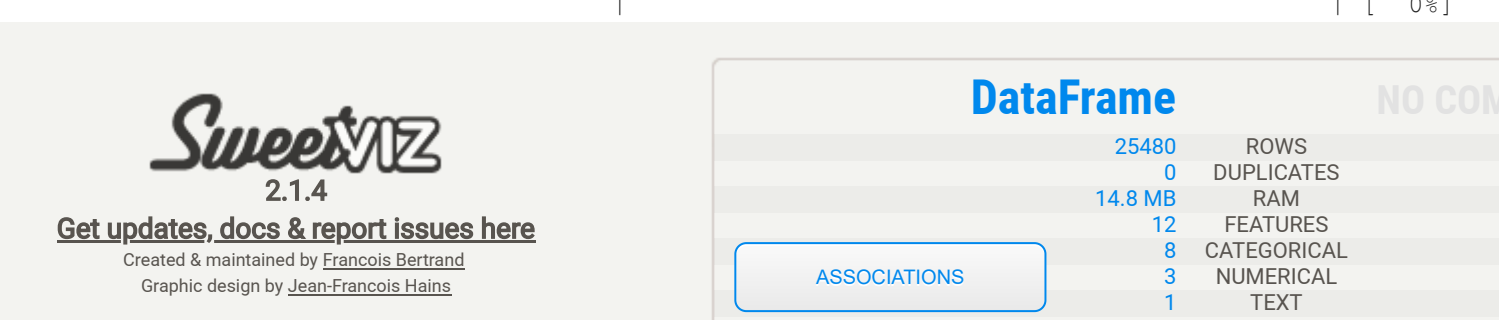
Summarize dataset:  0% | 0/5 [00:00<, ?it/s]
Generate report structure:  0% | 0/1 [00:00<, ?it/s]
Render HTML:  0% | 0/1 [00:00<, ?it/s]
```

Report Overview Variables Interactions Correlations Missing values Sample

Overview



Variables



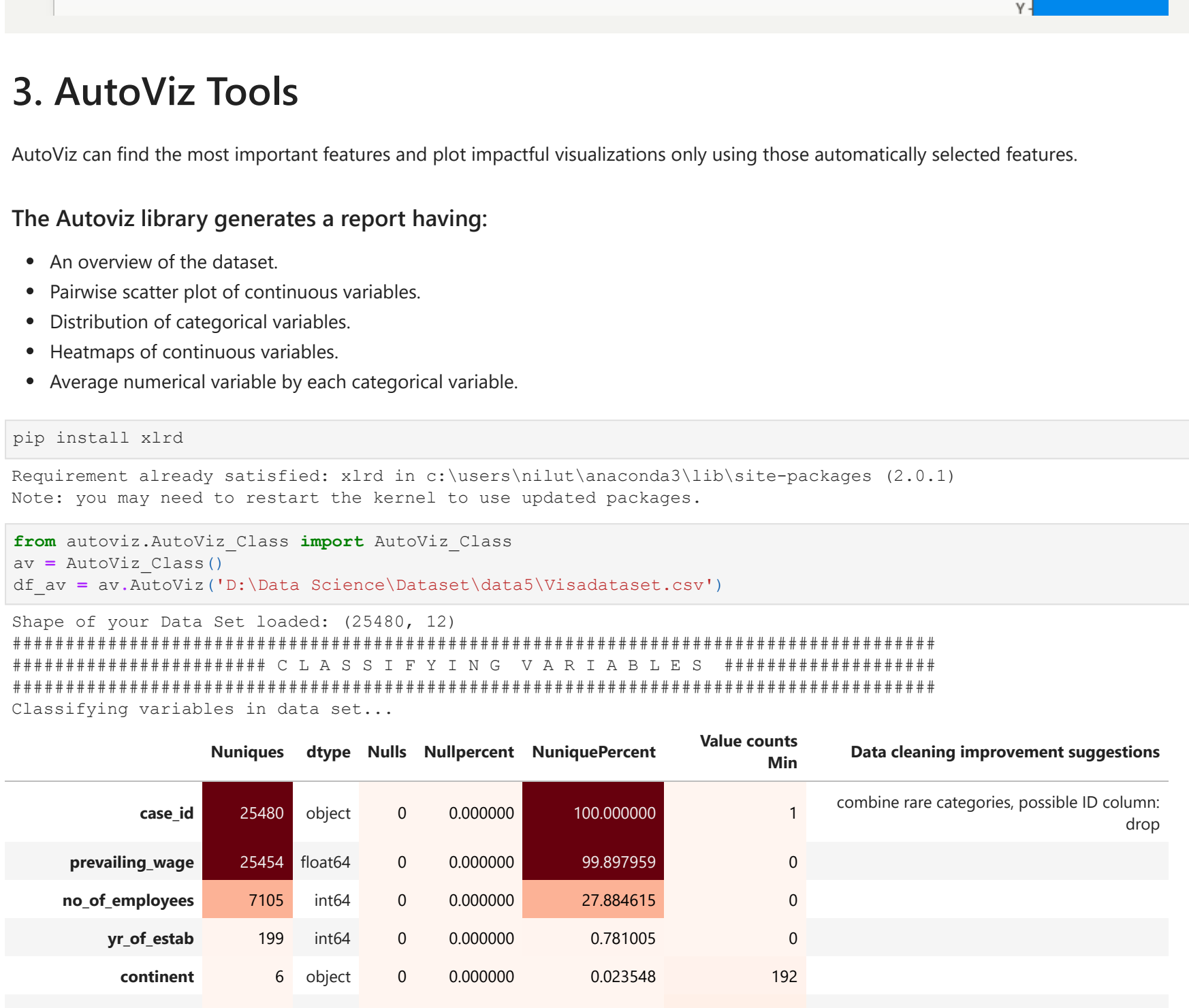
2. SweetViz Tools

SweetViz is a python library that generates beautiful, high-density visualizations which is useful for EDA with just few lines of code.

The **SweetViz** library generates a report having:

- An overview of the dataset.
- Variable properties.
- Categorical associations.
- Numerical associations.
- Most frequent, smallest, largest values for numerical features.

```
In [16]: import sweetviz as sv
report = sv.analyze(df)
report.show_notebook()
```



3. AutoViz Tools

AutoViz can find the most important features and plot impactful visualizations only using those automatically selected features.

The **AutoViz** library generates a report having:

- An overview of the dataset.
- Pairwise scatter plot of continuous variables.
- Distribution of categorical variables.
- Heatmaps of continuous variables.
- Average numerical variable by each categorical variable.

```
In [17]: pip install xlrd

Requirement already satisfied: xlrd in c:\users\nilut\anaconda3\lib\site-packages (2.0.1)
Note: you may need to restart the kernel to use updated packages.
```

```
In [18]: from autoviz.AutoViz_Class import AutoViz_Class
av = AutoViz_Class()
df_av = av.AutoViz('D:\Data Science\Dataset\data5\Visadataset.csv')
```

Shape of your Data Set loaded: (25480, 12)

C L A S S I F Y I N G V A R I A B L E S

Classifying variables in data set...

