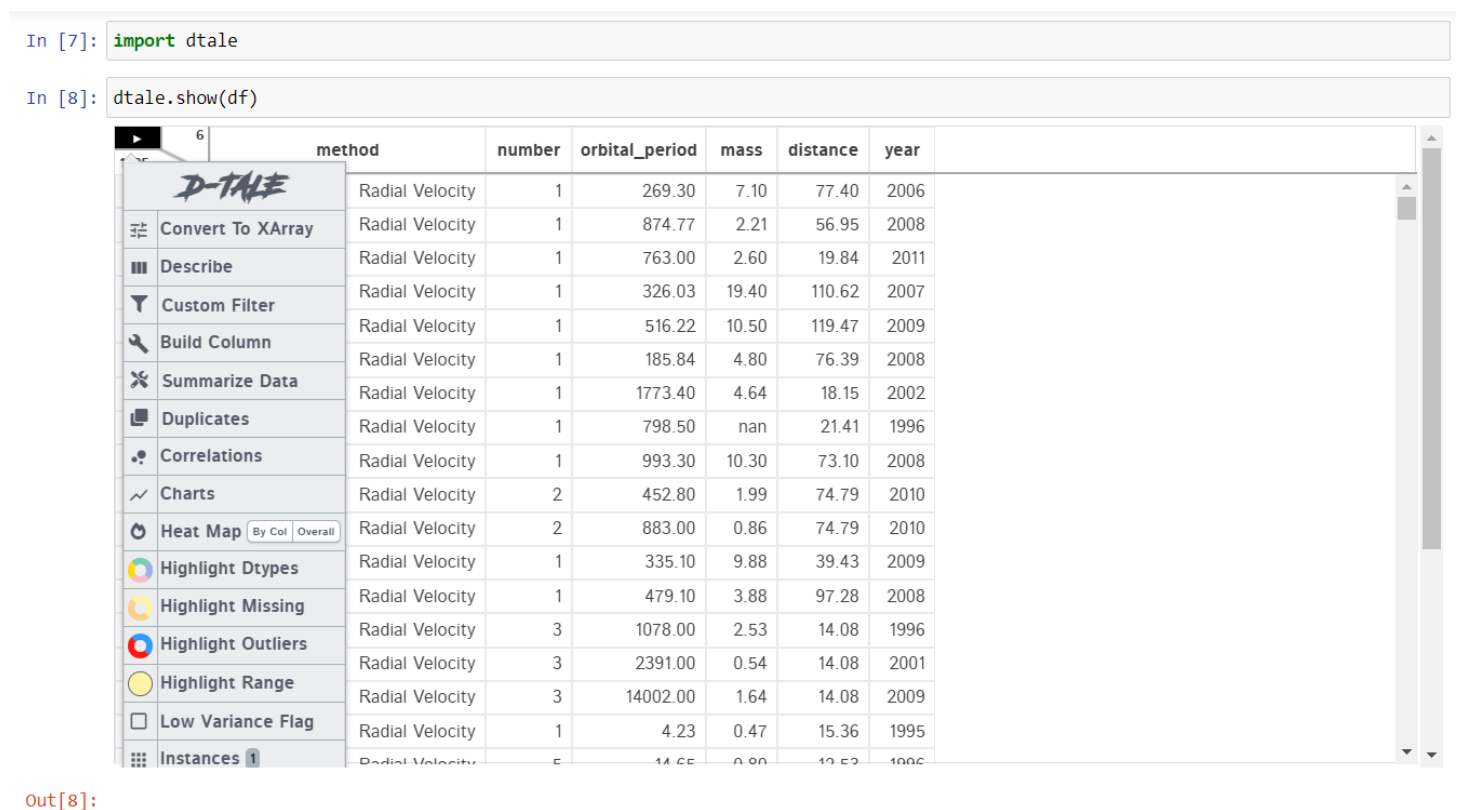


Project Name: EDA With Visualization Using DTale in 1 line - 3

Table of Contents

Demo
Overview
Motivation
Technical Aspect
Installation
Run/How to Use/Steps
Directory Tree/Structure of Project
To Do/Future Scope
Technologies Used/System Requirements/Tech Stack
Credits

Demo



Overview

This is an Exploratory Data Analysis with Visualization using dtale library.
EDA is essentially a type of storytelling for statisticians.
It allows us to uncover patterns and insights, often with visual methods, within data.
EDA is often the first step of the data modelling process.
This repository contains the code for EDA along with Visualization using python's various libraries.
It used numpy, pandas, dtale and seaborn libraries.

These libraries help to perform individually one particular functionality.

Numpy is used for working with arrays. It stands for Numerical Python.

Pandas objects rely heavily on Numpy objects.

Seaborn is data visualization library based on matplotlib.

D-Tale is the combination of a Flask back-end and a React front-end to bring you an easy way to view & analyze.

The purpose of creating this repository is to gain insights into EDA and Visualization.

These python libraries raised knowledge in discovering these libraries with practical use of it.

It leads to growth in my ML repository.

This above screenshot will help you to understand flow of output.

Motivation

My Aim for selecting is it brings data frame to life. And I have already worked with GUI in my previous projects so extending it. Since It is a GUI Based Exploratory Data Analysis tool that analyses and visualizes each and every aspect of the dataset. It is easy to use and blazingly fast. The reason behind building this is, to maximize I as analyst's insight into a data set and into the underlying structure of a data set, while providing all of the specific items that an analyst would want to extract. It is a way of visualizing, summarizing and interpreting the information that is hidden in rows and column format. EDA is understanding the data sets by summarizing their main characteristics often plotting them visually. This step is very important for me especially when I arrive at modelling the data in order to apply Machine learning. Another motive is, as a master's student I have learnt Data Mining Subject which has somewhere led me to also learn about EDA. Though, EDA and Data Mining has difference. EDA distinguishes itself from data mining, even though the two are closely related, as many EDA techniques have been adopted into data mining. Also, the goals of the two are very similar: EDA indeed makes sure that you explore the data in such a way that interesting features and relationships between features will become clearer. In EDA, you typically explore and compare many different variables with a variety of techniques to search and find systematic patterns. Data mining, on the other hand, is concerned with extracting patterns from the data. Those patterns provide insights into relationships between variables that can be used to improve business decisions. Also, in both cases, you have no a priori expectations or expectations that are not complete about the relations between the variables. However, in general, Data Mining can be said to be more application-oriented, while EDA is concerned with the basic nature of the underlying phenomena. In other words, Data Mining is relatively less concerned with identifying the specific relations between the involved variables. As a result, Data Mining accepts a "black box" approach to data exploration and doesn't only use techniques that are also used in EDA but also techniques such as Neural Networks to generate valid predictions but don't identify the specific nature of the relationships between the variables on which the predictions are based. Exploratory Data Analysis (EDA) is used on the one hand to answer questions, test business assumptions, generate hypotheses for further analysis. On the other hand, you can also use it to prepare the data for modelling. The thing that these two probably have in common is a good knowledge of your data to either get the answers that you need or to develop an intuition for interpreting the results of future modelling. There are a lot of ways to reach these goals: you can get a basic description of the data, visualize it, identify patterns in it, identify challenges of using the data, etc. Hence, I continue to gain knowledge while practicing the same and spread

intellectual wings in tech-heaven. Objective for doing EDA with dtale is, straightforward library. For client, if we can do work faster it will be better so it saves lot of time and win-win for all. I obtained grasp on this abstraction. It is a reward for me after practicing it manually. I will always recommend to practice manually first then once we get hang on it then can try with other libraries. I am trying many distinct techniques for training my visualization concept clarity because I also prefer viewing and retaining information rather than reading therefore this will also help while presenting in front of company's client and they will have finer perception of it.

Technical Aspect

Numpy contains a multi-dimensional array and matrix data structures. It works with the numerical data. Numpy is faster because is densely packed in memory due to its homogeneous type. It also frees the memory faster.

Pandas module mainly works with the tabular data. It contains Data Frame and Series. Pandas is 18 to 20 times slower than Numpy. Pandas is seriously a game changer when it comes to cleaning, transforming, manipulating and analyzing data.

Seaborn provides a high-level interface for drawing attractive and informative statistical graphics. It provides a variety of visualization patterns and visualize random distributions.

DTale is a Graphical Interface where we can select the data, we want to analyse and how to analyse using different graphs and plots.

Installation

Using intel core i5 9th generation with NVIDIA GFORCE GTX1650.

Windows 10 Environment Used.

Already Installed Anaconda Navigator for Python 3.x

The Code is written in Python 3.8.

If you don't have Python installed then please install Anaconda Navigator from its official site.

If you are using a lower version of Python you can upgrade using the pip package, ensuring you have the latest version of pip, *python -m pip install --upgrade pip and press Enter.*

Run/How to Use/Steps

Keep your internet connection on while running or accessing files and throughout too.

Follow this when you want to perform from scratch.

Open Anaconda Prompt, Perform the following steps:

```
cd <PATH>
```

```
pip install numpy
```

```
pip install pandas
```

```
pip install dtale
```

pip install seaborn

You can also create requirement.txt file as, pip freeze > requirements.txt
run files.

Follow this when you want to just perform on local machine.

Download ZIP File.

Right-Click on ZIP file in download section and select Extract file option, which will unzip file.

Move unzip folder to desired folder/location be it D drive or desktop etc.

Open Anaconda Prompt, write cd <PATH> and press Enter.

eg: cd C:\Users\Monica\Desktop\Projects\Python Projects 1\11)EDA_and_Visualization\dtale

In Anconda Prompt, pip install -r requirements.txt to install all packages.

Open in Jupyter Notebook, <filename>.ipynb

That is,

Open in Jupyter Notebook, 1)EDA_With_Visualization_using_dtale_in_1line_3.ipynb

Please be careful with spellings or numbers while typing filename and easier is just copy filename and then run it to avoid any silly errors.

Note: cd <PATH>

[Go to Folder where file is. Select the path from top and right-click and select copy option and paste it next to cd one space <path> and press enter, then you can access all files of that folder] [cd means change directory]

Directory Tree/Structure of Project

Folder: 11)EDA_and_Visualization > dtale

1)EDA_With_Visualization_using_dtale_in_1line_3.ipynb

To Do/Future Scope

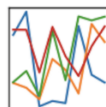
Can try another library.

Technologies Used/System Requirements/Tech Stack



NumPy

pandas
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$





	date	security_id	int_val	Col0	Col1	Col2	Col3	Col4	Col5	Col6	Col7	Col8	Col9	Col10	Col11
0	2018-09-13	102000	1370059006	1.65	1.68	-0.09	0.68	2.04	1.40	-0.74	-0.15	-0.40	-0.07	0.16	0.90
1	2018-09-13	102001	3277083703	0.87	-0.68	-2.27	-0.58	1.05	3.07	-0.31	0.40	-0.63	-0.74	-1.62	1.99
2	2018-09-13	102002	1620012202	0.27	1.76	-0.09	-0.79	0.57	-0.01	-0.28	0.14	-0.01	-0.01	-0.61	-0.61
3	2018-09-13	102003	1620000000	0.24	-0.04	-0.04	0.35	-0.04	0.85	-0.95	-0.01	-0.01	-0.14	-0.48	-0.48
4	2018-09-13	102004	05740007	-0.64	-1.36	-0.11	1.02	-0.64	-0.74	-1.84	-0.01	-0.01	-0.01	-0.01	-0.01
5	2018-09-13	102005	1700000000	-0.00	0.00	-0.73	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00
6	2018-09-13	102006	3041000000	0.00	0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00
7	2018-09-13	102007	1411000000	0.00	0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00
8	2018-09-13	102008	7207400000	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03
9	2018-09-13	102009	5400000000	0.78	0.36	0.86	-0.14	-0.57	-0.51	-0.99	-0.71	-0.48	-1.07	-0.41	-0.41
10	2018-09-13	102010	1570000000	-0.24	-0.04	-0.29	1.63	-0.39	-0.35	-0.79	-0.96	-1.02	-0.43	-1.40	-1.40
11	2018-09-13	102011	1967970000	-0.41	0.34	-1.21	0.17	-0.94	-0.35	-1.54	-0.81	-1.05	-0.57	-0.20	-0.20
12	2018-09-13	102012	3540400000	0.26	0.93	0.09	-0.11	0.77	-0.41	-0.25	-1.92	1.05	-1.74	-0.58	-0.58
13	2018-09-13	102013	5436620000	1.60	0.55	-1.83	-0.74	1.90	-1.42	0.52	-0.86	1.19	0.47	1.69	-0.17
14	2018-09-13	102014	2951000000	0.74	-1.22	-0.79	0.40	0.20	-0.81	-1.58	1.87	1.11	-0.48	0.21	0.12
15	2018-09-13	102015	5992000000	1.70	0.39	-0.51	-0.77	0.99	-0.69	0.34	-0.87	1.60	0.31	1.39	-0.09

Credits

<https://www.kdnuggets.com/2020/08/bring-pandas-dataframes-life-d-tale.html>

Krish Naik Channel