

Business Research and Data Analytics

Lecture 12: Exploratory Data Analysis in Python.

Igor Vyshnevskiy
Woosong University
May 29, 2024

Agenda

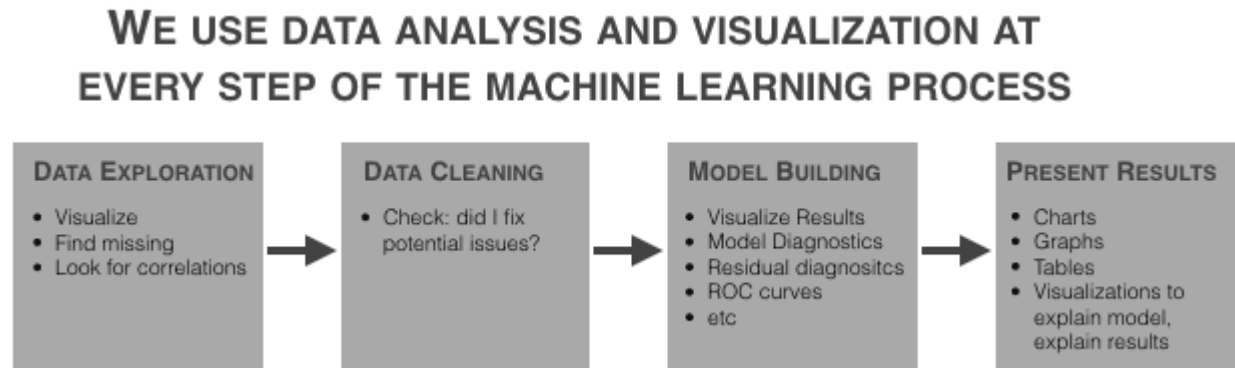
1. Intro to Exploratory Data Analysis (EDA)
2. Importance of EDA
3. Python packages for EDA
4. Common EDA Techniques
5. EDA Practicum
6. Final Exam instructions

1. Intro to Exploratory Data Analysis (EDA)

Intro

Exploratory Data Analysis (EDA) is the process of visually and quantitatively examining datasets to summarize their main characteristics and patterns, typically using statistical graphics, plots, and other data visualization techniques.

The **primary goal of EDA** is to understand the underlying structure, patterns, and relationships in the data to guide further analysis or model building.



Intro

- You can either explore data using graphs or through some *python* functions.
- There are two type of analysis. *Univariate and Bivariate*.
 - In the univariate, you analyzing a single attribute.
 - But in the bivariate, you analyzing an attribute with the target attribute.
- In the *non-graphical approach*, you will be using functions such as shape, summary, describe, isnull, info, datatypes and more.
- In the *graphical approach*, you will be using plots such as scatter, box, bar, density and correlation plots.
- **Remember:**
 - EDA is often an iterative process where initial findings might lead to further data cleaning, visualization, or modeling efforts, refining the analysis and insights.
 - Documenting EDA findings, insights, and the visualizations generated is crucial for sharing results and collaborating with other stakeholders.

Intro

Key Steps in EDA:

- Importing a dataset
- Understanding the big picture
- Preparation
- Understanding of variables
- Study of the relationships between variables
- Brainstorming

Intro

Key Steps of EDA in more details:

- *Data Loading*: Load the dataset into Python using libraries like Pandas, ensuring it is in a suitable format for analysis.
- *Data Cleaning*: Identify and handle missing values, outliers, and anomalies that could affect the analysis and conclusions.
- *Summary Statistics*: Compute basic statistics (mean, median, standard deviation, etc.) to summarize the central tendency and dispersion of the data.
- *Data Visualization*: Create various plots (histograms, scatter plots, box plots, etc.) to visualize the distribution, relationships, and trends in the data.
- *Exploratory Data Visualization*: Use advanced visualization techniques to delve deeper into relationships and patterns within the data.
- *Correlation Analysis*: Explore correlations between variables to identify potential predictive relationships.
- *Feature Engineering*: Derive new features or modify existing ones to improve the quality of input data for machine learning models.
- *Dimensionality Reduction*: Reduce the number of features while preserving essential information, aiding in model efficiency and interpretability.

2. Importance of EDA

Why EDA

- EDA helps in detecting errors, outliers, and anomalies.
- It provides a comprehensive understanding of the data's structure, improving modeling decisions.
- EDA guides feature selection and engineering, optimizing model performance.
- EDA aids in choosing appropriate machine learning algorithms based on the data's characteristics.
- EDA is *crucial* in various domains like finance (analyzing stock market trends), healthcare (clinical data analysis), marketing (customer segmentation), and more.

3. Python packages for EDA

What we use...

- *Pandas*: For data loading, cleaning, and initial data manipulation.
- *NumPy* and *SciPy*: For statistical analysis and mathematical computations.
- *Matplotlib* and *Seaborn*: For creating various visualizations.
- *Plotly*: For interactive and advanced visualizations.

4. Common EDA Techniques

What we use...

- *Histograms and Distributions*: To understand data distribution.
- *Scatter Plots*: To observe relationships between two variables.
- *Box Plots*: To detect outliers and distribution characteristics.
- *Heatmaps and Correlation Plots*: To visualize correlations between variables.
- *Pair Plots*: To visualize relationships across multiple variables.

5. EDA Example

What can be done...

- Please open the file “L13_work”.

5. EDA Practicum

What you will be doing...

- It's group work. You will be assigned to groups.
- Perform EDA for a given dataset in file “titanic.csv”, and finish by 2:40 pm.
- Present your findings briefly (3 min. per group).

You can get the file from LMS or running code

```
`import seaborn as sns`  
titanic = sns.load_dataset('titanic')`
```

Groups Allocation

Full name	ID number	Group
Gulnur Ismatullaeva	202112097	1
Mushfiqur Rahman Jeem	202212198	1
Islomjon Azizov	202212195	1
Tokhir Bakhtiyor Ugli Saydurasulov	202201176	1
Magar Ashik Rana	202312135	1
Sungmin Kim	202310004	1
Minseok Cho	202310019	2
Soyeong Park	202310009	2
Wonki Park	202110085	2
Md Mahadi Hassan	202312114	2
Mushtariy Bakhtiyor Kizi Ergashova	202312091	2
Shokhrukh Shaydullo Ugli Yodgorov	202212136	2
Khar Khar Hillary	202212224	3
Yoonhee Lee	202310014	3
Ekaterina Prokudina	202001112	3
Afnan Bin Zahid	202312146	3
Hermon Melese Mehret	202312361	3
Yonghwan Chun	202110097	3

Full name	ID number	Group
Masum Billah	202312093	4
Madina Khizbulaeva	202201135	4
Van Kip Lian Alfred	202212206	4
Seohyun Kim	202310003	4
Swikriti Ale	202312132	4
Sardorbek Bakhtiyor Ugli Makhmudov	202312086	4
Jin Young Kim	202310007	5
Othmane Salhi	202312392	5
Cho Kyeong-chan	202310017	5
Ahmmmed Syed Iftekhar	202312151	5
Betelhem Asrat Shiferaw	202312163	5
Zukhrabonu Otelloeva	202112157	6
Thant Kyaw Min	202212214	6
Abdulahadxon Jahongir Ugli Abbosov	202201181	6
Kim Jihyeok	202010173	6
Nana Kang	202310001	6

6. Final Exam instructions

Final Exam

- Group project report/presentation (Week 15): **30%**;
 - A group project report including the application of analytics principles and methods to a business issue will be required of the students.
 - A 20-minute PPT presentation and also an active participation in asking some questions to other presenters for a 15- minute discussion on the topic.
 - Students will be allocated to teams (up to 6 people in a group).
 - The assignment will be given today.
 - You have the remaining time and the time of our next class to work on this.
 - Please remember to upload your files in advance (script and ppt slides).
 - Only two groups will get the highest score.

Q & A

Thank you!