

DATA ANALYSIS AND PROGRAMMING WITH PYTHON

Lecture 3

Data preparation and visualization in Python

26 November 2023

1. Working with data in Python
2. Data preparation
 - a) Pandas
 - b) Pandas vs SQL
3. Data visualization
 - a) Matplotlib
 - b) Seaborn
 - c) Plotly
4. Practice in Jupyter Notebook

Pandas is a Python library used for working with data sets.

It has functions for

- analyzing,
- cleaning,
- exploring,
- and manipulating data.

The name "Pandas" has a reference to both "Panel Data", and "Python Data Analysis" and was created by Wes McKinney in 2008.



DATA PREPARATION

PANDAS DATAFRAME

	id	name	date	manner_of_death	armed	age	gender	race	city	state	signs_of_mental_illness	threat_level	flee	body_camera
0	3	Tim Elliot	2015-01-02	shot	gun	53.0	M	Asian	Shelton	WA	True	attack	Not fleeing	False
1	4	Lewis Lee Lembke	2015-01-02	shot	gun	47.0	M	White	Aloha	OR	False	attack	Not fleeing	False
2	5	John Paul Quintero	2015-01-03	shot and Tasered	unarmed	23.0	M	Hispanic	Wichita	KS	False	other	Not fleeing	False
3	8	Matthew Hoffman	2015-01-04	shot	toy weapon	32.0	M	White	San Francisco	CA	True	attack	Not fleeing	False
4	9	Michael Rodriguez	2015-01-04	shot	nail gun	39.0	M	Hispanic	Evans	CO	False	attack	Not fleeing	False
...
4890	5916	Rayshard Brooks	2020-06-12	shot	Taser	27.0	M	Black	Atlanta	GA	False	attack	Foot	True
4891	5925	Caine Van Pelt	2020-06-12	shot	gun	23.0	M	Black	Crown Point	IN	False	attack	Car	False
4892	5918	Hannah Fizer	2020-06-13	shot	unarmed	25.0	F	White	Sedalia	MO	False	other	Not fleeing	False
4893	5921	William Slyter	2020-06-13	shot	gun	22.0	M	White	Kansas City	MO	False	other	Other	False
4894	5924	Nicholas Hirsh	2020-06-15	shot	gun	31.0	M	White	Lawrence	KS	False	attack	Car	False

4895 rows × 15 columns

SQL	Python
<pre>SELECT * FROM titanic_test_data WHERE pclass = 1</pre>	<pre>titanic_df[titanic_df.pclass == 1]</pre>
<pre>SELECT * FROM titanic_test_data WHERE pclass = 1 OR pclass = 2</pre>	<pre>titanic_df[(titanic_df.pclass == 1) (titanic_df.pclass == 2)]</pre>
<pre>SELECT * FROM titanic_test_data WHERE pclass IN (1,2)</pre>	<pre>titanic_df[titanic_df.pclass.isin([1,2])]</pre>
<pre>SELECT name FROM titanic_test_data WHERE pclass = 1 AND gender = "male"</pre>	<pre>titanic_df[(titanic_df.pclass == 1) & (titanic_df.gender == "male")]["name"]</pre>
<pre>SELECT name, age FROM titanic_test_data WHERE pclass NOT IN (1,2)</pre>	<pre>titanic_df[~titanic_df.pclass.isin([1,2])] [["name", "age"]]</pre>

Pandas and **SQL** are both used for selecting filtering and grouping data.

While Pandas has more functionality, SQL performs better on large data sets.

As a part of Python - pandas is very convenient when you need not only extract data, but analyze it.

matplotlib Plot types User guide Tutorials Examples Reference Contribute Releases

Section Navigation

- Lines, bars and markers
- Images, contours and fields
- Subplots, axes and figures
- Statistics
- Pie and polar charts
- Text, labels and annotations
- Color
- Shapes and collections
- Style sheets
- Module - pyplot
- Module - axes_grid1
- Module - axisartist
- Showcase
- Animation
- Event handling
- Miscellaneous
- 3D plotting
- Scales
- Specialty plots
- Spines
- Ticks
- Units
- Embedding Matplotlib in graphical user interfaces
- Widgets
- Userdemo

Lines, bars and markers

Bar color demo

Bar Label Demo

Stacked bar chart

Grouped bar chart with labels

Horizontal bar chart

Broken Barh

CapStyle

Plotting categorical variables

Plotting the coherence of two signals

Cross spectral density (CSD)

Curve with error band

Errorbar limit selection

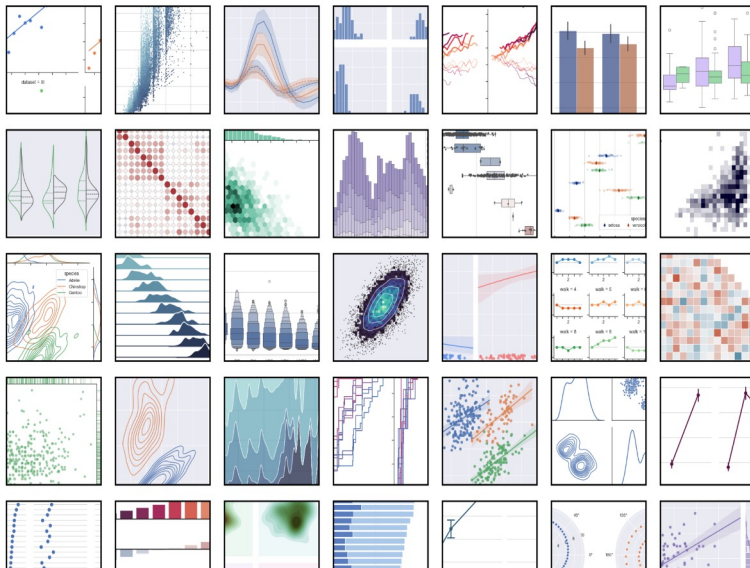
Matplotlib is the most common visualization tool. It perfectly balances simplicity and functionality and has lots of patterns.



Installing **Gallery** Tutorial API Releases Citing FAQ



Example gallery



Seaborn is an advanced statistical visualization tool. It has lot's of ready made visualizations in it's gallery that you can reuse for your research.

Maps

[More Maps »](#)



Mapbox Choropleth Maps



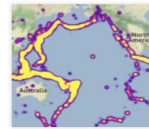
Lines on Mapbox



Filled Area on Maps



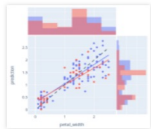
Bubble Maps



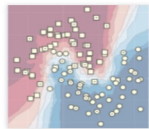
Mapbox Density Heatmap

Artificial Intelligence and Machine Learning

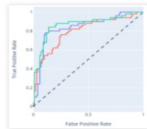
[More AI and ML »](#)



ML Regression



kNN Classification



ROC and PR Curves



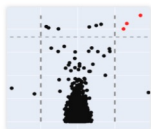
PCA Visualization



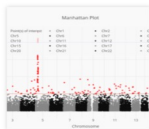
AI/ML Apps with Dash

Bioinformatics

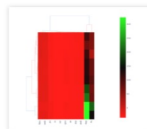
[More Bioinformatics »](#)



Volcano Plot



Manhattan Plot



Clustergram



Alignment Chart

Plotly is the most functional library and it's more dashboard based tool. It is very useful when you want to monitor the same data over time.

Let's practice in Jupyter Notebook!