

Data Visualization

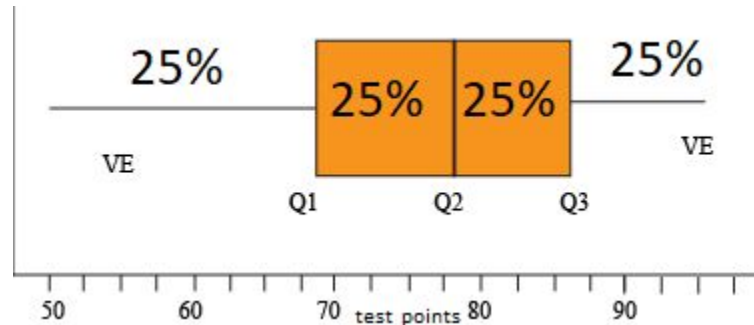
How to visualize data. Graphs, heatmap, pychart,...

How to call data.

Python Package named as **seaborn** (seaborn.pydata.org) statistical data visualization, based on **matplotlib**. You can install with `pip install seaborn` (documentation->Installing and getting started. It requires **numpy**, **scipy**, **matplotlib**, **pandas** package.

On gallery section of seaborn web we have a lot of different visualizations:

- ❑ Boxplots: source is sort the data and get the minimum, maximum, quartiles Q1(25%, lower median half), Q2(50% all data median) and Q3(75%, upper median half) and interquartile range ($IQR = Q3 - Q1$). Are a useful way to graph data divided into four quartiles, each with the same amount of values. The box diagram does not graph frequency or show individual statistics, but in them we can clearly see where half the data is. It is a good diagram to analyze the asymmetry in the data. Values below $Q1 - 1.5 \cdot IQR$ or higher than $Q3 + 1.5 \cdot IQR$ are considered outliers



- ❑ Heatmaps (anotated): color represents values, eg. pivot table. Color represents values of the data.
- ❑ Join Plot: Hexagonal bin plot with marginal distributions.
- ❑ Histograms: accurate representation of the distribution of numerical data, is an estimate of the probability distribution of a continuous variable. First step is to "bin" or "bucket" the range of values, that is, divide the entire range of values into a series of non overlapping intervals, that must be adjadcent, but not required of equal size; then count how many values fall into each interval. Histogram can be normalized to display relative frequencies; it shows the proportion of cases that fall into each of several categories, with the sum of the heights equaling 1. When bins are not equal width, vertical axis is not a frequency, is a frequency density.
- ❑ connell plots
- ❑

Assignment #3

It keeps simple: Download csv as dataframe, and make pivot table, after render this data with seaborn.

```
assignm3.py x
1 import pandas as pd
2 import seaborn as sns
3 # Set link to download csv
4 csv = 'https://raw.githubusercontent.com/resbaz/' \
5       'r-novice-gapminder-files/master/data/' \
6       'gapminder-FiveYearData.csv'
7 # Load Dataframe from csv that's in link
8 df = pd.read_csv(csv)
9 # Print Dataframe for test
10 print '- Printing Dataframe extracted from link -'
11 print df
12 print '- Making a pivot table, aggregation = average -'
13 dfpv = df.pivot_table(index='continent',
14                        columns='year',
15                        values='lifeExp',
16                        aggfunc='mean')
17 # Print Pivot Table for test
18 print '- Printing pivot table for test -'
19 print dfpv
20 # Generating Heat Map
21 print '- Generating heat map -'
22 sns.heatmap(dfpv).get_figure().savefig('assign3.png')
```

year	1952	1957	1962	...	1997	2002	2007
continent
Africa	39.135500	41.266346	43.319442	...	53.598269	53.325231	54.806038
Americas	53.279840	55.960280	58.398760	...	71.150480	72.422040	73.608120
Asia	46.314394	49.318544	51.563223	...	68.020515	69.233879	70.728485
Europe	64.408500	66.703067	68.539233	...	75.505167	76.700600	77.648600
Oceania	69.255000	70.295000	71.085000	...	78.190000	79.740000	80.719500

[1704 rows x 6 columns]
- Making a pivot table, aggregation = average -
- Printing pivot table for test -
[5 rows x 12 columns]
- Generating heat map -
>>>

A pivot table is an aggregation method to summarize info from several dimensions as rows and columns.

Rendering of heatmap: Expected Lifetime.

