



Workshop 1.2

Advanced Analytics and Applications [AAA]

I Multiple Choice

II Calculation

III Programming

Question 1:

Astronomers have been cataloguing distant objects in the sky using long-exposure CCD images. The objects need to be labeled as star, galaxy, nebula etc. The data is highly noisy, and the images are very faint. The cataloguing can take decades to complete. Which Method would you recommend the physicists use to automate the cataloguing process, and improve its effectiveness?

- Clustering
- Classification
- Regression

Question 1:

Astronomers have been cataloguing distant objects in the sky using long-exposure CCD images. The objects need to be labeled as star, galaxy, nebula etc. The data is highly noisy, and the images are very faint. The cataloguing can take decades to complete. Which Method would you recommend the physicists use to automate the cataloguing process, and improve its effectiveness?

- Clustering
- Classification
- Regression

Multiple Choice Questions

Question 2:

A customer of a supermarket owns a lot of cards, there is a lot of information about the supermarket like ID, age, gender, address, etc. assigned to each customer. Based on defined purchasing behavior and spending data, it can make a recommendation so the customer can make a decision to the supermarket which of his customers are most likely to react to a marketing campaign?

- Clustering
- Classification
- Regression

Through membership of a supermarket, customers like Customer A have a high spending Score was recommended to customers like customer B. Based on this recommendation using the supermarket which of his customers are most likely to react to a marketing campaign?

Question 3:

A large company named XYZ, employs, at any given point of time, around 4000 employees. However, every year, around 15% of its employees leave the company and need to be replaced with the talent pool available in the job market. The management believes that this level of attrition (employees leaving, either on their own or because they got fired) is bad for the company. Hence, the management has contracted an HR analytics firm to understand what factors they should focus on, in order to curb attrition. In other words, they want to know what changes they should make to their workplace, in order to get most of their employees to stay. Also, they want to know which of these variables is most important and needs to be addressed right away. Which method would you recommend using to model the probability of attrition where the results can be used by the management to understand what changes they should make to their workplace, in order to get most of their employees to stay?

- Clustering
- Classification
- Regression

Question 3:

A large company named XYZ, employs, at any given point of time, around 4000 employees. However, every year, around 15% of its employees leave the company and need to be replaced with the talent pool available in the job market. The management believes that this level of attrition (employees leaving, either on their own or because they got fired) is bad for the company. Hence, the management has contracted an HR analytics firm to understand what factors they should focus on, in order to curb attrition. In other words, they want to know what changes they should make to their workplace, in order to get most of their employees to stay. Also, they want to know which of these variables is most important and needs to be addressed right away. Which method would you recommend using to model the probability of attrition where the results can be used by the management to understand what changes they should make to their workplace, in order to get most of their employees to stay?

- Clustering
- Classification
- Regression

Multiple Choice

|| Calculation

||| Programming

Question 1: MSE

x	y	\hat{y}	$(y - \bar{y})^2$	$(\hat{y} - \bar{y})^2$	$(y - \hat{y})^2$
1	2	2,8			
2	4	3,4			
3	5	4			
4	4	4,6			
5	5	5,2			

Question 1: MSE

x	y	\hat{y}	$(y - \bar{y})^2$	$(\hat{y} - \bar{y})^2$	$(y - \hat{y})^2$
1	2	2,8	4	1,44	0,64
2	4	3,4	0	0,36	0,36
3	5	4	1	0	1
4	4	4,6	0	0,36	0,36
5	5	5,2	1	1,44	0,04

Calculation

Question 1: MSE

x	y	\hat{y}	$(y - \bar{y})^2$	$(\hat{y} - \bar{y})^2$	$(y - \hat{y})^2$
1	2	2,8	4	1,44	0,64
2	4	3,4	0	0,36	0,36
3	5	4	1	0	1
4	4	4,6	0	0,36	0,36
5	5	5,2	1	1,44	0,04

- a) Calculate MSE using the values from the above table (y = true value, \hat{y} = predicted value)

Question 1: MSE

x	y	\hat{y}	$(y - \bar{y})^2$	$(\hat{y} - \bar{y})^2$	$(y - \hat{y})^2$
1	2	2,8	4	1,44	0,64
2	4	3,4	0	0,36	0,36
3	5	4	1	0	1
4	4	4,6	0	0,36	0,36
5	5	5,2	1	1,44	0,04

- a) Calculate MSE using the values from the above table (y = true value, \hat{y} = predicted value)

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{2,4}{5} = 0,48$$

Question 1: MSE

x	y	\hat{y}	$(y - \bar{y})^2$	$(\hat{y} - \bar{y})^2$	$(y - \hat{y})^2$
1	2	2,8	4	1,44	0,64
2	4	3,4	0	0,36	0,36
3	5	4	1	0	1
4	4	4,6	0	0,36	0,36
5	5	5,2	1	1,44	0,04

b) When is MSE useful?

MSE can represent the difference between the actual observations and the observation values predicted by the model

Puts large emphasis on large deviations

Calculation

Question 2: R²

x	y	\hat{y}	$(y - \bar{y})^2$	$(\hat{y} - \bar{y})^2$	$(y - \hat{y})^2$
1	2	2,8	4	1,44	0,64
2	4	3,4	0	0,36	0,36
3	5	4	1	0	1
4	4	4,6	0	0,36	0,36
5	5	5,2	1	1,44	0,04

- a) Calculate R-Squared using the values from the above table (y = true value, \hat{y} = predicted value)

Question 2: R²

x	y	\hat{y}	$(y - \bar{y})^2$	$(\hat{y} - \bar{y})^2$	$(y - \hat{y})^2$
1	2	2,8	4	1,44	0,64
2	4	3,4	0	0,36	0,36
3	5	4	1	0	1
4	4	4,6	0	0,36	0,36
5	5	5,2	1	1,44	0,04

- a) Calculate R-Squared using the values from the above table (y = true value, \hat{y} = predicted value)

$$\begin{aligned}
 R^2 &= 1 - \frac{SSR}{SST} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{2,4}{6} \\
 &= \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{3,6}{6} = 0,6
 \end{aligned}$$

Question 2: R²

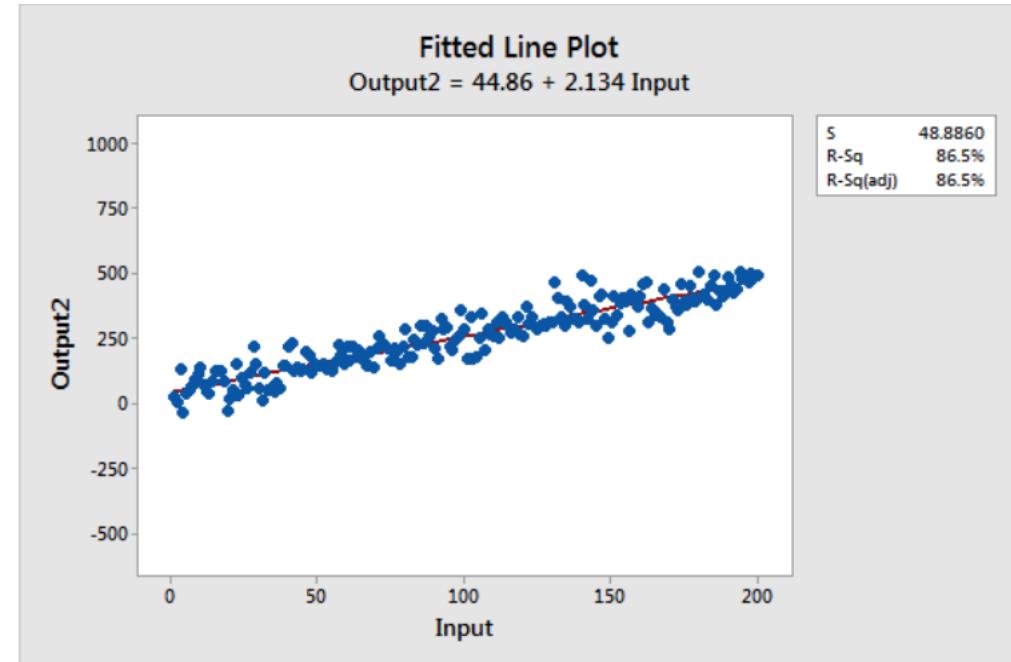
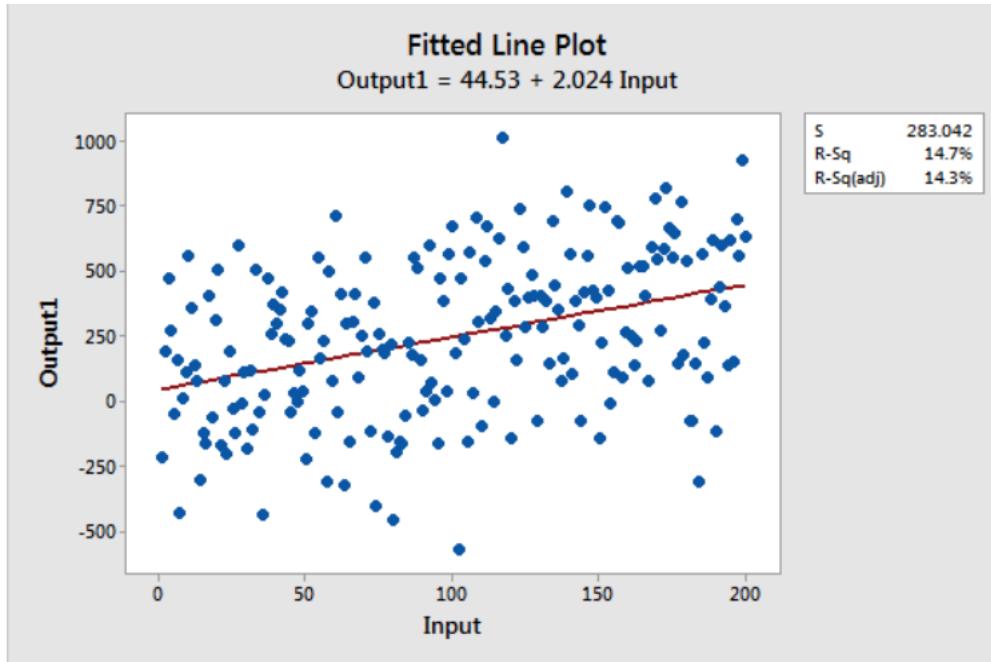
b) When is R-Squared useful? Compare it to MSE.

Question 2: R²

- b) When is R-Squared useful? Compare it to MSE.

R-squared represents the fraction of variance of response variable captured by the regression model rather than the MSE which captures the residual error.

Calculation Question 2: R²



I Multiple Choice

II Calculation

III Programming

Question 3: Working with Pandas

a) (Pandas) Create a python notebook that loads the *employee_survey_data.csv* into a Data Frame and find out the following stats:

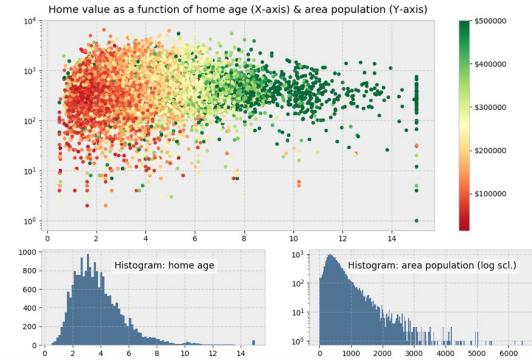
(1) Dimensionality: _____

(2) Sum of rows with missing values: _____

(3) Number of employees most satisfied with their jobs: _____

Tip: Find out max value. Filter employees by max value. Sum number of employees with max job satisfaction.

Question 3: Plotting with Matplotlib



(Plotting) In this task, we want to learn how ~~create~~ a plot using the python library matplotlib with multiple axes in one figure. To do so, create a ~~Jupyter~~ Notebook, then download and initialize the California housing data set.

Then, try to render a plot similar to the one depicted in Figure 1. In the first figure, we can want to render a scatter plot of “age” (Home Age) and “pop” (Population). In the second plot, we want to visualize a histogram of “age” (Home Age). Finally, we want to render a plot a histogram of the “pop” (Population). Follow the template from Figure 2 with 3 axis objects, and by using subplots (for more information, read the docs for matplotlib)

Matplotlib: Motivation for this workshop



The screenshot shows the official Matplotlib website at <https://matplotlib.org/3.2.1/>. The page features a large "matplotlib" logo with a circular icon containing a sunburst chart. Below the logo, the text "Version 3.2.1" is visible. A dark blue navigation bar contains links for "Installation", "Documentation", "Examples", "Tutorials", and "Contributing". Under "Documentation", there is a link to "home | contents ». The main content area has a light gray background. It starts with a heading "Overview" followed by a red "¶" symbol. Below this, there are two data entries: "Release: 3.2.1" and "Date: April 08, 2020". There is also a link to "Download PDF". A sidebar on the left lists several documentation sections: "User's Guide" (which is underlined in orange), "Installing", "Tutorials", "Interactive plots", "Whats New", "History", "GitHub Stats", "Previous What's New", "License", and "Credits".

Extensive documentation – The documentation of matplotlib is extensive and in detail BUT also hard for beginners.

Digging through the code is tedious and time-consuming (+70.000 lines).

Thus, we will cover important concept of matplotlib without commencing an "Indiana Jones" like adventure through the codebase.

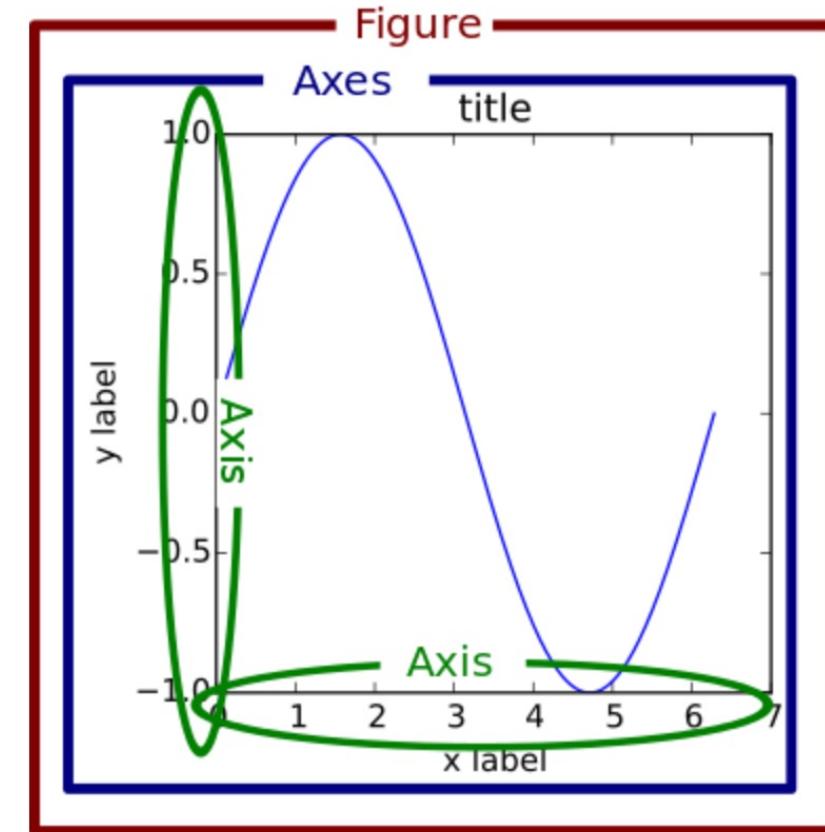
Matplotlib: History



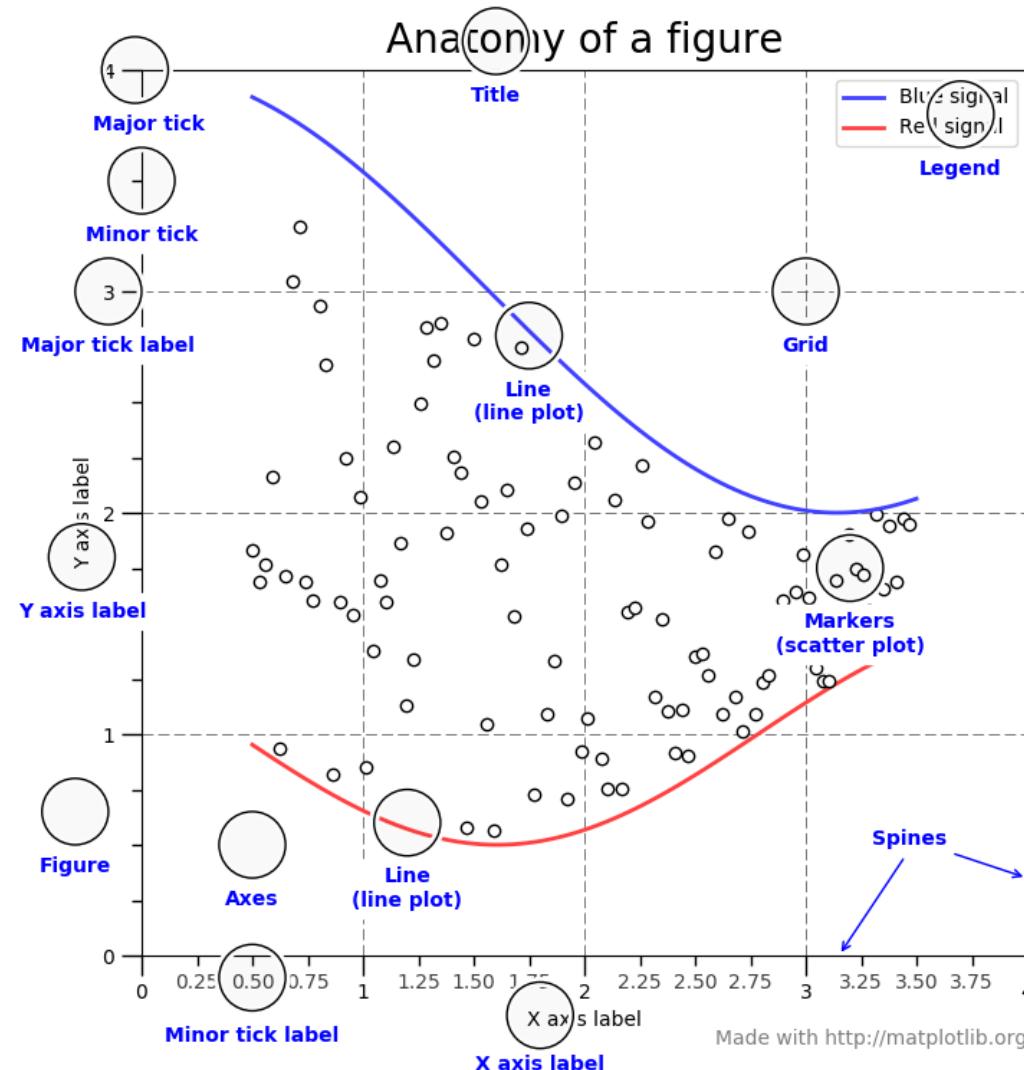
John D. Hunter (1968 – 2012) was an American neurobiologist and the original author of matplotlib.

Matplotlib: The object hierarchy

`plt.plot([1, 2, 3])` → Create tree-like object hierarchy



Matplotlib: Object hierarchy under the hood



Matplotlib: Stateful versus Stateless Approach

```
import matplotlib.pyplot as plt
```

```
plt.plot([1, 2, 3])  
plt.title(„Hello World“)  
plt.ylabel(„Emissions“)
```



Seems like a stateful "enviroment" but in reality it just mimics stateful behavior.

Matplotlib: Stateful versus Stateless Approach

Python

>>>

```
# matplotlib/pyplot.py
>>> def plot(*args, **kwargs):
...     """An abridged version of plt.plot()."""
...     ax = plt.gca()
...     return ax.plot(*args, **kwargs)

>>> def gca(**kwargs):
...     """Get the current Axes of the current Figure."""
...     return plt.gcf().gca(**kwargs)
```

Applies for almost any function in pyplot

Matplotlib: Stateful versus Stateless Approach

Bottom line: Matplotlib interface is object-oriented and thus is built upon the stateless approach

Contact



For general questions and enquiries on **research**, **teaching**, **job openings** and new **projects** refer to our website at www.is3.uni-koeln.de



For specific enquiries regarding this course contact us by sending an email to the **IS3 teaching** address at is3-teaching@wiso.uni-koeln.de