

Probability densities and data Visualization with PYTHON

Objective: Upon completing this exercise it is expected that you:

- Get an understanding of the many ways data can be visualized including histograms, boxplots, and scatter plots.

Piazza discussion forum: You can get help by asking questions on Piazza: piazza.com/dtu.dk/fall2019/october2019

Software installation: Extract the Python toolbox from the Dropbox folder . Start Spyder and add the toolbox directory (`<base-dir>/02450Toolbox_Python/Tools/`) to PYTHONPATH (Tools/PYTHONPATH manager in Spyder). Remember the purpose of the exercises is not to re-write the code from scratch but to work with the scripts provided in the directory `<base-dir>/02450Toolbox_Python/Scripts/` Representation of data in Python:

	Python var.	Type	Size	Description
	X	numpy.array	$N \times M$	Data matrix: The rows correspond to N data objects, each of which contains M attributes.
	attributeNames	list	$M \times 1$	Attribute names: Name (string) for each of the M attributes.
	N	integer	Scalar	Number of data objects.
	M	integer	Scalar	Number of attributes.
Classification	y	numpy.array	$N \times 1$	Class index: For each data object, y contains a class index, $y_n \in \{0, 1, \dots, C - 1\}$, where C is the total number of classes.
	classNames	list	$C \times 1$	Class names: Name (string) for each of the C classes.
	C	integer	Scalar	Number of classes.

4.1 Visualizing Fisher's Iris data

In this exercise we will reproduce most of the figures in "Introduction to Data Mining." section 3.3 in Matlab using the Iris flower dataset that was also introduced in Exercise 1.

- 4.1.1 The Iris data set is available in the Excel file `Data/iris.xls`. Load the data into Python and generate all the variables described in *Representation of data in Python* using the script `ex4_2_1.py`.

Script details:

-
- You can use the package `xlrd` which you have used before in exercise 2.2.

4.1.2 Inspect and run `ex4_2_2.py`. The script plots a histogram of each of the four attributes in the Iris data.

Script details:

- You can use the command `hist` to plot a histogram.
- Use indexing to extract each attribute. For example, `X[:,m-1]` extracts the `m`'th attribute.
- For multiple plots in one figure window you can use the command `subplot(n,m,i)`.

Show on the graph that the petal length is either between 1 and 2 cm. or between 3 and 7 cm. but that no flowers in the data set have a petal length between 2 and 3 cm. Do you think this could be useful to discriminate between the different types of flowers?

4.1.3 Inspect and run `ex4_2_3.py`. The script produces a boxplot of the four attributes in the Iris data as shown in Figure 3.11 in the book.

Script details:

- Take a look at the function `boxplot`.
- Type `help(boxplot)` to see how you can adjust the boxplot and add labels.

This boxplot shows the same information as the histogram in the previous exercise. Discuss the advantages and disadvantages of the two types of plots.

4.1.4 Inspect and run `ex4_2_4.py`. The scripts produces a boxplot for each attribute for each class as shown in Figure 3.12 in the book.

Script details:

- Use the functions `subplot()` and `boxplot()`.
- The variable `y` $\in \{0, 1, \dots, C - 1\}$ contains the class labels. To extract the data objects belonging to, say, class `c`, you can use `y` to index into `X` like this: `X[(y==c), :]`
- It is easier to compare the boxplots if they are all on the same axis. To do this, you can use the function `ylim()`.

Show on the graph that all the Iris-setosa in this data set have a petal length between 1 and 2 cm.

4.1.5 Inspect and run `ex4_2_5.py`. The scripts produces a matrix of scatter plots of each combination of two attributes against each other as shown in Figure 3.16 in the book.

Script details:

- To make a scatter plot, you can use the function `plot(x,y,s)` where `x` and `y` specify the coordinates and `s` is a string that specifies the line style and plot symbol, e.g., `s='.'` to make dots.
- To extract the data values for the `m`'th attribute in the `c`'th class, you can write `X[(y==c),m]`.
- You can use the command `hold all` to plot multiple plots on top of each other.

Say you want to discriminate between the three types of flowers using only the length and width of either sepal or petal. Show on the graph why it would be better to use petal length and width rather than sepal length and width.

4.1.6 Inspect and run `ex4_2_6.py`. The script produces a 3-dimensional scatter plot of three attributes as shown in Figure 3.17 in the book.

Script details:

- Read more about plotting in 3 dimensions:
matplotlib.sourceforge.net/mpl_toolkits/mplot3d/tutorial.html
- To plot in 3 dimensions you need to import `matplotlib.pyplot` as earlier, and additionally `Axes3D` from `mpl_toolkits.mplot3d`.

Try rotating the data. Can you find an angle where the three types of flower are separated in the plot?

4.1.7 Use the script `ex4_2_7.py` to plot the data matrix as an image as shown in Figure 3.23 in the book. The data matrix should be standardized to have zero mean and unit standard deviation.

Script details:

- You can use the function `imagesc()` to plot an image.
- By default, the image will be smoothed, what is not always desired when you look at the data. Use parameter `interpolation='None'` to display raw data.
- The function `zscore()` can be used to standardize the data matrix.

You are welcome to try out other plotting methods for the data. Matplotlib online repository is a good source of inspiration:

matplotlib.sourceforge.net/gallery.html

4.2 Visualizing Wine Data

We will in this part of the exercise consider two datasets related to red and white variants of the Portuguese "Vinho Verde" wine [1], the data has been downloaded from <http://archive.ics.uci.edu/ml/datasets/Wine+Quality>. Only physico-chemical and sensory attributes are available, i.e., there is no data about grape types, wine brand, wine selling price, etc. The data has the following attributes:

Attributes 1–11 are based on physicochemical tests and attribute 12 on human judging. Later in the course we will attempt to predict whether a wine is a red or

#	Attribute	Unit
1	Fixed acidity (tartaric)	g/dm ³
2	Volatile acidity (acetic)	g/dm ³
3	Citric acid	g/dm ³
4	Residual sugar	g/dm ³
5	Chlorides	g/dm ³
6	Free sulfur dioxide	mg/dm ³
7	Total sulfur dioxide	mg/dm ³
8	Density	g/cm ³
9	pH	pH
10	Sulphates	g/dm ³
11	Alcohol	% vol.
12	Quality score	0–10

white wine based on these attributes and we will also attempt to predict the wine quality. Unfortunately, the data set has many observations that can be considered outliers and in order to carry out analyses later it is important to remove the corrupt observations.

The aim of this exercise is to use visualization to identify outliers and remove these outliers from the data. It might be necessary to remove some outliers before other outlying observations become visible. Thus, the process of finding and removing outliers is often iterative. The wine data is stored in a Matlab file, `Data/wine.mat`.

4.2.1 Inspect and run the script `ex4_3_1.py`. The script loads the data into Python using the `scipy.io.loadmat()` function, as in previous exercises. This dataset contains many observations that can be considered outliers and the visualization tools you have worked with in the previous exercise is used to identify the outliers in the data set. How many outliers are identified by the script? How are the identified outliers removed from the data set?

Script details:

- You can use your solutions to the previous exercise as a starting point for making your visualizations.
- Say you want to find all data objects for which the alcohol percentage (attribute number 11) is not greater than 100%. You can mask them simply as `mask=(X[:,10]<=100)`.
- You can use the mask to eliminate the outlier observations (rows of data matrix). For instance you can write `X=X[mask,:]` where `mask` indicates the data objects that should be maintained. Remember also to remove them from the class index vector, `y=y[mask,1]` and to recompute `N`.

We will later in the course attempt to classify the type of wine (white or red) as well as predict the quality of wine based on the physicochemical tests. Visual inspection of the data can give an indication of the difficulty of these tasks.

- 4.2.2 Are there any of the measurements that seem to be well suited in order to discriminate between red and white wines? What plots are particular useful in order to investigate this?
- 4.2.3 Does any of the 12 attributes appear to correlate with each other? What plots are well suited to investigate this?
- 4.2.4 Can you identify any clear relationship between the various physicochemical measurements of the wines and the quality of the wines as rated by human judges?

4.3 Tasks for the report

After today, you can address the last questions for part one of the report:

- **Data visualization(s) based on suitable visualization techniques including a principal component analysis (PCA).**

Touch upon the following subjects, use visualizations when it appears sensible. *Keep in mind the ACCENT principles and Tufte's guidelines when you visualize the data.*

- Are there issues with outliers in the data,
 - do the attributes appear to be normal distributed,
 - are variables correlated,
 - does the primary machine learning modeling aim appear to be feasible based on your visualizations.
- **A discussion explaining what you have learned about the data.**

Summarize here the most important things you have learned about the data and give also your thoughts on whether your primary modeling task(s) appears to be feasible based on your visualization.

References

- [1] P. Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis. Modeling wine preferences by data mining from physicochemical properties. *In Decision Support Systems, Elsevier*, 47(4):547–553, 2009.