

CMM201_Week10_livewrangle

November 22, 2019

```
In [1]: # This cell is used to change parameter of the rise slideshow,  
# such as the window width/height and enabling a scroll bar
```

```
from notebook.services.config import ConfigManager  
cm = ConfigManager()  
cm.update('livereveal', {  
    'width': 1000,  
    'height': 600,  
    'scroll': True,  
})
```

```
Out[1]: {'width': 1000, 'height': 600, 'scroll': True}
```

```
In [ ]: # This command transforms the Jupyter notebook into a slideshow  
!jupyter nbconvert CMM201_Week10.ipynb --to slides --post serve  
# once a new browser opens, replace the "#" after the the_notebook.slides.html in the  
# ?print-pdf so that the url looks most likely like http://127.0.0.1:8000/the_notebook  
# finally, print to PDF file
```

1 Importing and Wrangling Data in Python

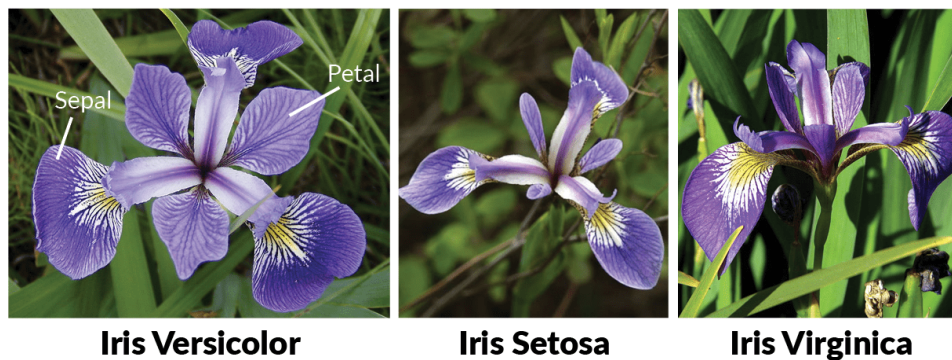
1.1 Aims of the Lecture

- Learn how to import numerical data to Python from different sources.
- Understand how to select certain parts of the imported data.

1.2 Example

1.2.1 Loading Data from a Module

- Python has a module called **scikit-learn** or *sklearn* which contains several datasets commonly used in data science and business analytics.
- For this exercise, we will use the **IRIS** database contained in this module.
- This dataset contains the sepal and petal lengths and widths from 150 samples of 3 different types of the iris flower.



iris.png

- Unlike last week, we will **NOT** work with the actual images, but rather with the numerical information extracted from samples.
- First, we need to install **sklearn**:

```
In [2]: !pip install sklearn
```

Collecting sklearn

Downloading <https://files.pythonhosted.org/packages/1e/7a/dbb3be0ce9bd5c8b7e3d87328e79063f8b>

Requirement already satisfied: scikit-learn in c:\anaconda\lib\site-packages (from sklearn) (0

Requirement already satisfied: scipy>=0.13.3 in c:\anaconda\lib\site-packages (from scikit-learn)

Requirement already satisfied: numpy>=1.8.2 in c:\anaconda\lib\site-packages (from scikit-learn)

Building wheels for collected packages: sklearn

Building wheel for sklearn (setup.py): started

Building wheel for sklearn (setup.py): finished with status 'done'

Stored in directory: C:\Users\CM8738\AppData\Local\pip\Cache\wheels\76\03\bb\589d421d27431bc

Successfully built sklearn

Installing collected packages: sklearn

Successfully installed sklearn-0.0

- Then, we can load the iris dataset:

```
In [3]: ## Load iris dataset
        from sklearn import datasets
        iris = datasets.load_iris()
        print(type(iris))
```

```
<class 'sklearn.utils.Bunch'>
```

- The dataset is contained on a **dictionary-like** structure referred to as **sklearn.utils.Bunch**.
- If you print it, you will see a lot of things contained:

```
In [4]: print(iris)
```

```
{'data': array([[5.1, 3.5, 1.4, 0.2],
[4.9, 3. , 1.4, 0.2],
[4.7, 3.2, 1.3, 0.2],
[4.6, 3.1, 1.5, 0.2],
[5. , 3.6, 1.4, 0.2],
[5.4, 3.9, 1.7, 0.4],
[4.6, 3.4, 1.4, 0.3],
[5. , 3.4, 1.5, 0.2],
[4.4, 2.9, 1.4, 0.2],
[4.9, 3.1, 1.5, 0.1],
[5.4, 3.7, 1.5, 0.2],
[4.8, 3.4, 1.6, 0.2],
[4.8, 3. , 1.4, 0.1],
[4.3, 3. , 1.1, 0.1],
[5.8, 4. , 1.2, 0.2],
[5.7, 4.4, 1.5, 0.4],
[5.4, 3.9, 1.3, 0.4],
[5.1, 3.5, 1.4, 0.3],
[5.7, 3.8, 1.7, 0.3],
[5.1, 3.8, 1.5, 0.3],
[5.4, 3.4, 1.7, 0.2],
[5.1, 3.7, 1.5, 0.4],
[4.6, 3.6, 1. , 0.2],
[5.1, 3.3, 1.7, 0.5],
[4.8, 3.4, 1.9, 0.2],
[5. , 3. , 1.6, 0.2],
[5. , 3.4, 1.6, 0.4],
[5.2, 3.5, 1.5, 0.2],
[5.2, 3.4, 1.4, 0.2],
[4.7, 3.2, 1.6, 0.2],
[4.8, 3.1, 1.6, 0.2],
[5.4, 3.4, 1.5, 0.4],
[5.2, 4.1, 1.5, 0.1],
[5.5, 4.2, 1.4, 0.2],
[4.9, 3.1, 1.5, 0.2],
[5. , 3.2, 1.2, 0.2],
[5.5, 3.5, 1.3, 0.2],
[4.9, 3.6, 1.4, 0.1],
[4.4, 3. , 1.3, 0.2],
[5.1, 3.4, 1.5, 0.2],
[5. , 3.5, 1.3, 0.3],
[4.5, 2.3, 1.3, 0.3],
[4.4, 3.2, 1.3, 0.2],
[5. , 3.5, 1.6, 0.6],
[5.1, 3.8, 1.9, 0.4],
[4.8, 3. , 1.4, 0.3],
[5.1, 3.8, 1.6, 0.2],
[4.6, 3.2, 1.4, 0.2],
```

[5.3, 3.7, 1.5, 0.2],
 [5. , 3.3, 1.4, 0.2],
 [7. , 3.2, 4.7, 1.4],
 [6.4, 3.2, 4.5, 1.5],
 [6.9, 3.1, 4.9, 1.5],
 [5.5, 2.3, 4. , 1.3],
 [6.5, 2.8, 4.6, 1.5],
 [5.7, 2.8, 4.5, 1.3],
 [6.3, 3.3, 4.7, 1.6],
 [4.9, 2.4, 3.3, 1.],
 [6.6, 2.9, 4.6, 1.3],
 [5.2, 2.7, 3.9, 1.4],
 [5. , 2. , 3.5, 1.],
 [5.9, 3. , 4.2, 1.5],
 [6. , 2.2, 4. , 1.],
 [6.1, 2.9, 4.7, 1.4],
 [5.6, 2.9, 3.6, 1.3],
 [6.7, 3.1, 4.4, 1.4],
 [5.6, 3. , 4.5, 1.5],
 [5.8, 2.7, 4.1, 1.],
 [6.2, 2.2, 4.5, 1.5],
 [5.6, 2.5, 3.9, 1.1],
 [5.9, 3.2, 4.8, 1.8],
 [6.1, 2.8, 4. , 1.3],
 [6.3, 2.5, 4.9, 1.5],
 [6.1, 2.8, 4.7, 1.2],
 [6.4, 2.9, 4.3, 1.3],
 [6.6, 3. , 4.4, 1.4],
 [6.8, 2.8, 4.8, 1.4],
 [6.7, 3. , 5. , 1.7],
 [6. , 2.9, 4.5, 1.5],
 [5.7, 2.6, 3.5, 1.],
 [5.5, 2.4, 3.8, 1.1],
 [5.5, 2.4, 3.7, 1.],
 [5.8, 2.7, 3.9, 1.2],
 [6. , 2.7, 5.1, 1.6],
 [5.4, 3. , 4.5, 1.5],
 [6. , 3.4, 4.5, 1.6],
 [6.7, 3.1, 4.7, 1.5],
 [6.3, 2.3, 4.4, 1.3],
 [5.6, 3. , 4.1, 1.3],
 [5.5, 2.5, 4. , 1.3],
 [5.5, 2.6, 4.4, 1.2],
 [6.1, 3. , 4.6, 1.4],
 [5.8, 2.6, 4. , 1.2],
 [5. , 2.3, 3.3, 1.],
 [5.6, 2.7, 4.2, 1.3],
 [5.7, 3. , 4.2, 1.2],

[5.7, 2.9, 4.2, 1.3],
 [6.2, 2.9, 4.3, 1.3],
 [5.1, 2.5, 3. , 1.1],
 [5.7, 2.8, 4.1, 1.3],
 [6.3, 3.3, 6. , 2.5],
 [5.8, 2.7, 5.1, 1.9],
 [7.1, 3. , 5.9, 2.1],
 [6.3, 2.9, 5.6, 1.8],
 [6.5, 3. , 5.8, 2.2],
 [7.6, 3. , 6.6, 2.1],
 [4.9, 2.5, 4.5, 1.7],
 [7.3, 2.9, 6.3, 1.8],
 [6.7, 2.5, 5.8, 1.8],
 [7.2, 3.6, 6.1, 2.5],
 [6.5, 3.2, 5.1, 2.],
 [6.4, 2.7, 5.3, 1.9],
 [6.8, 3. , 5.5, 2.1],
 [5.7, 2.5, 5. , 2.],
 [5.8, 2.8, 5.1, 2.4],
 [6.4, 3.2, 5.3, 2.3],
 [6.5, 3. , 5.5, 1.8],
 [7.7, 3.8, 6.7, 2.2],
 [7.7, 2.6, 6.9, 2.3],
 [6. , 2.2, 5. , 1.5],
 [6.9, 3.2, 5.7, 2.3],
 [5.6, 2.8, 4.9, 2.],
 [7.7, 2.8, 6.7, 2.],
 [6.3, 2.7, 4.9, 1.8],
 [6.7, 3.3, 5.7, 2.1],
 [7.2, 3.2, 6. , 1.8],
 [6.2, 2.8, 4.8, 1.8],
 [6.1, 3. , 4.9, 1.8],
 [6.4, 2.8, 5.6, 2.1],
 [7.2, 3. , 5.8, 1.6],
 [7.4, 2.8, 6.1, 1.9],
 [7.9, 3.8, 6.4, 2.],
 [6.4, 2.8, 5.6, 2.2],
 [6.3, 2.8, 5.1, 1.5],
 [6.1, 2.6, 5.6, 1.4],
 [7.7, 3. , 6.1, 2.3],
 [6.3, 3.4, 5.6, 2.4],
 [6.4, 3.1, 5.5, 1.8],
 [6. , 3. , 4.8, 1.8],
 [6.9, 3.1, 5.4, 2.1],
 [6.7, 3.1, 5.6, 2.4],
 [6.9, 3.1, 5.1, 2.3],
 [5.8, 2.7, 5.1, 1.9],
 [6.8, 3.2, 5.9, 2.3],

[illegible]

- Therefore, we need to extract each index of this dictionary into a different variables to understand and analyse them separately.
- First, we will import the data:

```
In [5]: data = iris['data']
        print(data, type(data), data.shape)
```

[5.1 3.5 1.4 0.2]
[4.9 3. 1.4 0.2]
[4.7 3.2 1.3 0.2]
[4.6 3.1 1.5 0.2]
[5. 3.6 1.4 0.2]
[5.4 3.9 1.7 0.4]
[4.6 3.4 1.4 0.3]
[5. 3.4 1.5 0.2]
[4.4 2.9 1.4 0.2]
[4.9 3.1 1.5 0.1]
[5.4 3.7 1.5 0.2]
[4.8 3.4 1.6 0.2]
[4.8 3. 1.4 0.1]
[4.3 3. 1.1 0.1]
[5.8 4. 1.2 0.2]
[5.7 4.4 1.5 0.4]
[5.4 3.9 1.3 0.4]
[5.1 3.5 1.4 0.3]
[5.7 3.8 1.7 0.3]
[5.1 3.8 1.5 0.3]
[5.4 3.4 1.7 0.2]
[5.1 3.7 1.5 0.4]
[4.6 3.6 1. 0.2]
[5.1 3.3 1.7 0.5]
[4.8 3.4 1.9 0.2]
[5. 3. 1.6 0.2]
[5. 3.4 1.6 0.4]

[5.2 3.5 1.5 0.2]
 [5.2 3.4 1.4 0.2]
 [4.7 3.2 1.6 0.2]
 [4.8 3.1 1.6 0.2]
 [5.4 3.4 1.5 0.4]
 [5.2 4.1 1.5 0.1]
 [5.5 4.2 1.4 0.2]
 [4.9 3.1 1.5 0.2]
 [5. 3.2 1.2 0.2]
 [5.5 3.5 1.3 0.2]
 [4.9 3.6 1.4 0.1]
 [4.4 3. 1.3 0.2]
 [5.1 3.4 1.5 0.2]
 [5. 3.5 1.3 0.3]
 [4.5 2.3 1.3 0.3]
 [4.4 3.2 1.3 0.2]
 [5. 3.5 1.6 0.6]
 [5.1 3.8 1.9 0.4]
 [4.8 3. 1.4 0.3]
 [5.1 3.8 1.6 0.2]
 [4.6 3.2 1.4 0.2]
 [5.3 3.7 1.5 0.2]
 [5. 3.3 1.4 0.2]
 [7. 3.2 4.7 1.4]
 [6.4 3.2 4.5 1.5]
 [6.9 3.1 4.9 1.5]
 [5.5 2.3 4. 1.3]
 [6.5 2.8 4.6 1.5]
 [5.7 2.8 4.5 1.3]
 [6.3 3.3 4.7 1.6]
 [4.9 2.4 3.3 1.]
 [6.6 2.9 4.6 1.3]
 [5.2 2.7 3.9 1.4]
 [5. 2. 3.5 1.]
 [5.9 3. 4.2 1.5]
 [6. 2.2 4. 1.]
 [6.1 2.9 4.7 1.4]
 [5.6 2.9 3.6 1.3]
 [6.7 3.1 4.4 1.4]
 [5.6 3. 4.5 1.5]
 [5.8 2.7 4.1 1.]
 [6.2 2.2 4.5 1.5]
 [5.6 2.5 3.9 1.1]
 [5.9 3.2 4.8 1.8]
 [6.1 2.8 4. 1.3]
 [6.3 2.5 4.9 1.5]
 [6.1 2.8 4.7 1.2]
 [6.4 2.9 4.3 1.3]

[6.6 3. 4.4 1.4]
 [6.8 2.8 4.8 1.4]
 [6.7 3. 5. 1.7]
 [6. 2.9 4.5 1.5]
 [5.7 2.6 3.5 1.]
 [5.5 2.4 3.8 1.1]
 [5.5 2.4 3.7 1.]
 [5.8 2.7 3.9 1.2]
 [6. 2.7 5.1 1.6]
 [5.4 3. 4.5 1.5]
 [6. 3.4 4.5 1.6]
 [6.7 3.1 4.7 1.5]
 [6.3 2.3 4.4 1.3]
 [5.6 3. 4.1 1.3]
 [5.5 2.5 4. 1.3]
 [5.5 2.6 4.4 1.2]
 [6.1 3. 4.6 1.4]
 [5.8 2.6 4. 1.2]
 [5. 2.3 3.3 1.]
 [5.6 2.7 4.2 1.3]
 [5.7 3. 4.2 1.2]
 [5.7 2.9 4.2 1.3]
 [6.2 2.9 4.3 1.3]
 [5.1 2.5 3. 1.1]
 [5.7 2.8 4.1 1.3]
 [6.3 3.3 6. 2.5]
 [5.8 2.7 5.1 1.9]
 [7.1 3. 5.9 2.1]
 [6.3 2.9 5.6 1.8]
 [6.5 3. 5.8 2.2]
 [7.6 3. 6.6 2.1]
 [4.9 2.5 4.5 1.7]
 [7.3 2.9 6.3 1.8]
 [6.7 2.5 5.8 1.8]
 [7.2 3.6 6.1 2.5]
 [6.5 3.2 5.1 2.]
 [6.4 2.7 5.3 1.9]
 [6.8 3. 5.5 2.1]
 [5.7 2.5 5. 2.]
 [5.8 2.8 5.1 2.4]
 [6.4 3.2 5.3 2.3]
 [6.5 3. 5.5 1.8]
 [7.7 3.8 6.7 2.2]
 [7.7 2.6 6.9 2.3]
 [6. 2.2 5. 1.5]
 [6.9 3.2 5.7 2.3]
 [5.6 2.8 4.9 2.]
 [7.7 2.8 6.7 2.]


```

[6.3 2.7 4.9 1.8]
[6.7 3.3 5.7 2.1]
[7.2 3.2 6.  1.8]
[6.2 2.8 4.8 1.8]
[6.1 3.  4.9 1.8]
[6.4 2.8 5.6 2.1]
[7.2 3.  5.8 1.6]
[7.4 2.8 6.1 1.9]
[7.9 3.8 6.4 2. ]
[6.4 2.8 5.6 2.2]
[6.3 2.8 5.1 1.5]
[6.1 2.6 5.6 1.4]
[7.7 3.  6.1 2.3]
[6.3 3.4 5.6 2.4]
[6.4 3.1 5.5 1.8]
[6.  3.  4.8 1.8]
[6.9 3.1 5.4 2.1]
[6.7 3.1 5.6 2.4]
[6.9 3.1 5.1 2.3]
[5.8 2.7 5.1 1.9]
[6.8 3.2 5.9 2.3]
[6.7 3.3 5.7 2.5]
[6.7 3.  5.2 2.3]
[6.3 2.5 5.  1.9]
[6.5 3.  5.2 2. ]
[6.2 3.4 5.4 2.3]
[5.9 3.  5.1 1.8]] <class 'numpy.ndarray'> (150, 4)

```

```
In [7]: set(iris)
```

```
Out[7]: {'DESCR', 'data', 'feature_names', 'filename', 'target', 'target_names'}
```

- The data is stored in a *numpy array* of 150 rows and 4 columns, each corresponding to the measurements of a flower.
- Then, we will import the headers of the data:

```
In [6]: header = iris['feature_names']
        print(header, type(header))
```

```
['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)'] <class 'list'>
```

- **Why do you think the data and the header are stored separately?**
- Afterwards, we will import the **class/target**:

```
In [8]: target = iris['target']
        print(target, type(target), target.shape)
```

```
[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
 2 2] <class 'numpy.ndarray'> (150,)
```

- The class/target is a *numpy array* which contains the **category** of each flowers.
- Each sample is labelled as 0, 1 or 2 instead of the iris type since the labels can be better used as numbers.
- A separate key called **target_names** contains the name corresponding to each numerical label.

```
In [9]: target_names = iris['target_names']
        print(target_names, type(target_names), target_names.shape)
```

```
['setosa' 'versicolor' 'virginica'] <class 'numpy.ndarray'> (3,)
```

- Finally, just in case you are interested, there is an entry containing the description of the dataset (a string):

```
In [10]: iris['DESCR']
```

```
Out[10]: '.. _iris_dataset:\n\nIris plants dataset\n-----\n\n**Data Set Character
```

1.2.2 Wrangling Data

- Accessing an individual entry of the dataset (along with its class/target):

```
In [16]: n = 89
        print(data[n], target[n])
```

```
[5.5 2.5 4.  1.3] 1
```

- Creating a table for each iris type (“manually”)

```
In [19]: setosa = data[0:50]
        print(setosa, setosa.shape)
```

```
[[5.1 3.5 1.4 0.2]
 [4.9 3.  1.4 0.2]
 [4.7 3.2 1.3 0.2]
 [4.6 3.1 1.5 0.2]
 [5.  3.6 1.4 0.2]
 [5.4 3.9 1.7 0.4]
 [4.6 3.4 1.4 0.3]
```

```

[5.  3.4 1.5 0.2]
[4.4 2.9 1.4 0.2]
[4.9 3.1 1.5 0.1]
[5.4 3.7 1.5 0.2]
[4.8 3.4 1.6 0.2]
[4.8 3.  1.4 0.1]
[4.3 3.  1.1 0.1]
[5.8 4.  1.2 0.2]
[5.7 4.4 1.5 0.4]
[5.4 3.9 1.3 0.4]
[5.1 3.5 1.4 0.3]
[5.7 3.8 1.7 0.3]
[5.1 3.8 1.5 0.3]
[5.4 3.4 1.7 0.2]
[5.1 3.7 1.5 0.4]
[4.6 3.6 1.  0.2]
[5.1 3.3 1.7 0.5]
[4.8 3.4 1.9 0.2]
[5.  3.  1.6 0.2]
[5.  3.4 1.6 0.4]
[5.2 3.5 1.5 0.2]
[5.2 3.4 1.4 0.2]
[4.7 3.2 1.6 0.2]
[4.8 3.1 1.6 0.2]
[5.4 3.4 1.5 0.4]
[5.2 4.1 1.5 0.1]
[5.5 4.2 1.4 0.2]
[4.9 3.1 1.5 0.2]
[5.  3.2 1.2 0.2]
[5.5 3.5 1.3 0.2]
[4.9 3.6 1.4 0.1]
[4.4 3.  1.3 0.2]
[5.1 3.4 1.5 0.2]
[5.  3.5 1.3 0.3]
[4.5 2.3 1.3 0.3]
[4.4 3.2 1.3 0.2]
[5.  3.5 1.6 0.6]
[5.1 3.8 1.9 0.4]
[4.8 3.  1.4 0.3]
[5.1 3.8 1.6 0.2]
[4.6 3.2 1.4 0.2]
[5.3 3.7 1.5 0.2]
[5.  3.3 1.4 0.2]] (50, 4)

```

```

In [21]: ## Use this cell to create and print versicolor
         ## and virginica (with the shape)
         versicolor = data[50:100]

```

```

virginica = data[100:150]
print(versicolor,versicolor.shape)
print(virginica,virginica.shape)

[[7.  3.2 4.7 1.4]
 [6.4 3.2 4.5 1.5]
 [6.9 3.1 4.9 1.5]
 [5.5 2.3 4.  1.3]
 [6.5 2.8 4.6 1.5]
 [5.7 2.8 4.5 1.3]
 [6.3 3.3 4.7 1.6]
 [4.9 2.4 3.3 1. ]
 [6.6 2.9 4.6 1.3]
 [5.2 2.7 3.9 1.4]
 [5.  2.  3.5 1. ]
 [5.9 3.  4.2 1.5]
 [6.  2.2 4.  1. ]
 [6.1 2.9 4.7 1.4]
 [5.6 2.9 3.6 1.3]
 [6.7 3.1 4.4 1.4]
 [5.6 3.  4.5 1.5]
 [5.8 2.7 4.1 1. ]
 [6.2 2.2 4.5 1.5]
 [5.6 2.5 3.9 1.1]
 [5.9 3.2 4.8 1.8]
 [6.1 2.8 4.  1.3]
 [6.3 2.5 4.9 1.5]
 [6.1 2.8 4.7 1.2]
 [6.4 2.9 4.3 1.3]
 [6.6 3.  4.4 1.4]
 [6.8 2.8 4.8 1.4]
 [6.7 3.  5.  1.7]
 [6.  2.9 4.5 1.5]
 [5.7 2.6 3.5 1. ]
 [5.5 2.4 3.8 1.1]
 [5.5 2.4 3.7 1. ]
 [5.8 2.7 3.9 1.2]
 [6.  2.7 5.1 1.6]
 [5.4 3.  4.5 1.5]
 [6.  3.4 4.5 1.6]
 [6.7 3.1 4.7 1.5]
 [6.3 2.3 4.4 1.3]
 [5.6 3.  4.1 1.3]
 [5.5 2.5 4.  1.3]
 [5.5 2.6 4.4 1.2]
 [6.1 3.  4.6 1.4]
 [5.8 2.6 4.  1.2]
 [5.  2.3 3.3 1. ]

```

[5.6 2.7 4.2 1.3]
 [5.7 3. 4.2 1.2]
 [5.7 2.9 4.2 1.3]
 [6.2 2.9 4.3 1.3]
 [5.1 2.5 3. 1.1]
 [5.7 2.8 4.1 1.3]] (50, 4)
 [[6.3 3.3 6. 2.5]
 [5.8 2.7 5.1 1.9]
 [7.1 3. 5.9 2.1]
 [6.3 2.9 5.6 1.8]
 [6.5 3. 5.8 2.2]
 [7.6 3. 6.6 2.1]
 [4.9 2.5 4.5 1.7]
 [7.3 2.9 6.3 1.8]
 [6.7 2.5 5.8 1.8]
 [7.2 3.6 6.1 2.5]
 [6.5 3.2 5.1 2.]
 [6.4 2.7 5.3 1.9]
 [6.8 3. 5.5 2.1]
 [5.7 2.5 5. 2.]
 [5.8 2.8 5.1 2.4]
 [6.4 3.2 5.3 2.3]
 [6.5 3. 5.5 1.8]
 [7.7 3.8 6.7 2.2]
 [7.7 2.6 6.9 2.3]
 [6. 2.2 5. 1.5]
 [6.9 3.2 5.7 2.3]
 [5.6 2.8 4.9 2.]
 [7.7 2.8 6.7 2.]
 [6.3 2.7 4.9 1.8]
 [6.7 3.3 5.7 2.1]
 [7.2 3.2 6. 1.8]
 [6.2 2.8 4.8 1.8]
 [6.1 3. 4.9 1.8]
 [6.4 2.8 5.6 2.1]
 [7.2 3. 5.8 1.6]
 [7.4 2.8 6.1 1.9]
 [7.9 3.8 6.4 2.]
 [6.4 2.8 5.6 2.2]
 [6.3 2.8 5.1 1.5]
 [6.1 2.6 5.6 1.4]
 [7.7 3. 6.1 2.3]
 [6.3 3.4 5.6 2.4]
 [6.4 3.1 5.5 1.8]
 [6. 3. 4.8 1.8]
 [6.9 3.1 5.4 2.1]
 [6.7 3.1 5.6 2.4]
 [6.9 3.1 5.1 2.3]

```

[5.8 2.7 5.1 1.9]
[6.8 3.2 5.9 2.3]
[6.7 3.3 5.7 2.5]
[6.7 3.  5.2 2.3]
[6.3 2.5 5.  1.9]
[6.5 3.  5.2 2. ]
[6.2 3.4 5.4 2.3]
[5.9 3.  5.1 1.8]] (50, 4)

```

- Creating a table for each iris type (“automatically”)

```

In [24]: ## In case that data is not in order or you don't want to count,
         ## we can use this alternative:
         import numpy as np
         setosa2 = data[np.where(target==0)]
         print(setosa, setosa.shape)

```

```

[[5.1 3.5 1.4 0.2]
 [4.9 3.  1.4 0.2]
 [4.7 3.2 1.3 0.2]
 [4.6 3.1 1.5 0.2]
 [5.  3.6 1.4 0.2]
 [5.4 3.9 1.7 0.4]
 [4.6 3.4 1.4 0.3]
 [5.  3.4 1.5 0.2]
 [4.4 2.9 1.4 0.2]
 [4.9 3.1 1.5 0.1]
 [5.4 3.7 1.5 0.2]
 [4.8 3.4 1.6 0.2]
 [4.8 3.  1.4 0.1]
 [4.3 3.  1.1 0.1]
 [5.8 4.  1.2 0.2]
 [5.7 4.4 1.5 0.4]
 [5.4 3.9 1.3 0.4]
 [5.1 3.5 1.4 0.3]
 [5.7 3.8 1.7 0.3]
 [5.1 3.8 1.5 0.3]
 [5.4 3.4 1.7 0.2]
 [5.1 3.7 1.5 0.4]
 [4.6 3.6 1.  0.2]
 [5.1 3.3 1.7 0.5]
 [4.8 3.4 1.9 0.2]
 [5.  3.  1.6 0.2]
 [5.  3.4 1.6 0.4]
 [5.2 3.5 1.5 0.2]
 [5.2 3.4 1.4 0.2]
 [4.7 3.2 1.6 0.2]

```

[4.8 3.1 1.6 0.2]
[5.4 3.4 1.5 0.4]
[5.2 4.1 1.5 0.1]
[5.5 4.2 1.4 0.2]
[4.9 3.1 1.5 0.2]
[5. 3.2 1.2 0.2]
[5.5 3.5 1.3 0.2]
[4.9 3.6 1.4 0.1]
[4.4 3. 1.3 0.2]
[5.1 3.4 1.5 0.2]
[5. 3.5 1.3 0.3]
[4.5 2.3 1.3 0.3]
[4.4 3.2 1.3 0.2]
[5. 3.5 1.6 0.6]
[5.1 3.8 1.9 0.4]
[4.8 3. 1.4 0.3]
[5.1 3.8 1.6 0.2]
[4.6 3.2 1.4 0.2]
[5.3 3.7 1.5 0.2]
[5. 3.3 1.4 0.2]] (50, 4)

```
In [25]: ## Verify that we get the same
         setosa == setosa2
```

[illegible]

[illegible]

- Creating a new table with “less” columns (by column number):

```
In [45]: ## creating a "reduced" table  
## with only the first two columns  
print(data[0])  
data_red1 = data[0,-3:]  
print(data_red1,data_red1.shape)
```

```
[5.1 3.5 1.4 0.2]
[3.5 1.4 0.2] (3,)
```

```
In [ ]: ## Use this cell to create a new dataset called data_red2
        ## with the last two columns
```

```
In [60]: ## Use this cell to create a new dataset called data_red3
         ## with the first and the third columns
         data[:,[0,2]][[6,7,8],:]
```

```
Out[60]: array([[4.6, 1.4],
                [5. , 1.5],
                [4.4, 1.4]])
```



```
In [61]: ## creating a "reduced" table with only the first column
col_0 = data[:,0]
print(col_0,col_0.shape)

[5.1 4.9 4.7 4.6 5.  5.4 4.6 5.  4.4 4.9 5.4 4.8 4.8 4.3 5.8 5.7 5.4 5.1
 5.7 5.1 5.4 5.1 4.6 5.1 4.8 5.  5.  5.2 5.2 4.7 4.8 5.4 5.2 5.5 4.9 5.
 5.5 4.9 4.4 5.1 5.  4.5 4.4 5.  5.1 4.8 5.1 4.6 5.3 5.  7.  6.4 6.9 5.5
 6.5 5.7 6.3 4.9 6.6 5.2 5.  5.9 6.  6.1 5.6 6.7 5.6 5.8 6.2 5.6 5.9 6.1
 6.3 6.1 6.4 6.6 6.8 6.7 6.  5.7 5.5 5.5 5.8 6.  5.4 6.  6.7 6.3 5.6 5.5
 5.5 6.1 5.8 5.  5.6 5.7 5.7 6.2 5.1 5.7 6.3 5.8 7.1 6.3 6.5 7.6 4.9 7.3
 6.7 7.2 6.5 6.4 6.8 5.7 5.8 6.4 6.5 7.7 7.7 6.  6.9 5.6 7.7 6.3 6.7 7.2
 6.2 6.1 6.4 7.2 7.4 7.9 6.4 6.3 6.1 7.7 6.3 6.4 6.  6.9 6.7 6.9 5.8 6.8
 6.7 6.7 6.3 6.5 6.2 5.9] (150,)
```

- Getting a column by it's name:

```
In [68]: print(header)
header.index('sepal length (cm)')

['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']
```

Out[68]: 0

```
In [ ]: sepal_length = data[:,header.index('sepal length (cm)')]
print(sepal_length,sepal_length.shape)

In [76]: pw = header.index('petal width (cm)')
sl = header.index('sepal length (cm)')
data[:,[pw,sl]]
```

Out[76]: array([[0.2, 5.1],
[0.2, 4.9],
[0.2, 4.7],
[0.2, 4.6],
[0.2, 5.],
[0.4, 5.4],
[0.3, 4.6],
[0.2, 5.],
[0.2, 4.4],
[0.1, 4.9],
[0.2, 5.4],
[0.2, 4.8],
[0.1, 4.8],
[0.1, 4.3],
[0.2, 5.8],
[0.4, 5.7],
[0.4, 5.4],

[0.3, 5.1],
[0.3, 5.7],
[0.3, 5.1],
[0.2, 5.4],
[0.4, 5.1],
[0.2, 4.6],
[0.5, 5.1],
[0.2, 4.8],
[0.2, 5.],
[0.4, 5.],
[0.2, 5.2],
[0.2, 5.2],
[0.2, 4.7],
[0.2, 4.8],
[0.4, 5.4],
[0.1, 5.2],
[0.2, 5.5],
[0.2, 4.9],
[0.2, 5.],
[0.2, 5.5],
[0.1, 4.9],
[0.2, 4.4],
[0.2, 5.1],
[0.3, 5.],
[0.3, 4.5],
[0.2, 4.4],
[0.6, 5.],
[0.4, 5.1],
[0.3, 4.8],
[0.2, 5.1],
[0.2, 4.6],
[0.2, 5.3],
[0.2, 5.],
[1.4, 7.],
[1.5, 6.4],
[1.5, 6.9],
[1.3, 5.5],
[1.5, 6.5],
[1.3, 5.7],
[1.6, 6.3],
[1. , 4.9],
[1.3, 6.6],
[1.4, 5.2],
[1. , 5.],
[1.5, 5.9],
[1. , 6.],
[1.4, 6.1],
[1.3, 5.6],

[1.4, 6.7],
[1.5, 5.6],
[1. , 5.8],
[1.5, 6.2],
[1.1, 5.6],
[1.8, 5.9],
[1.3, 6.1],
[1.5, 6.3],
[1.2, 6.1],
[1.3, 6.4],
[1.4, 6.6],
[1.4, 6.8],
[1.7, 6.7],
[1.5, 6.],
[1. , 5.7],
[1.1, 5.5],
[1. , 5.5],
[1.2, 5.8],
[1.6, 6.],
[1.5, 5.4],
[1.6, 6.],
[1.5, 6.7],
[1.3, 6.3],
[1.3, 5.6],
[1.3, 5.5],
[1.2, 5.5],
[1.4, 6.1],
[1.2, 5.8],
[1. , 5.],
[1.3, 5.6],
[1.2, 5.7],
[1.3, 5.7],
[1.3, 6.2],
[1.1, 5.1],
[1.3, 5.7],
[2.5, 6.3],
[1.9, 5.8],
[2.1, 7.1],
[1.8, 6.3],
[2.2, 6.5],
[2.1, 7.6],
[1.7, 4.9],
[1.8, 7.3],
[1.8, 6.7],
[2.5, 7.2],
[2. , 6.5],
[1.9, 6.4],
[2.1, 6.8],

```
[2. , 5.7],
[2.4, 5.8],
[2.3, 6.4],
[1.8, 6.5],
[2.2, 7.7],
[2.3, 7.7],
[1.5, 6. ],
[2.3, 6.9],
[2. , 5.6],
[2. , 7.7],
[1.8, 6.3],
[2.1, 6.7],
[1.8, 7.2],
[1.8, 6.2],
[1.8, 6.1],
[2.1, 6.4],
[1.6, 7.2],
[1.9, 7.4],
[2. , 7.9],
[2.2, 6.4],
[1.5, 6.3],
[1.4, 6.1],
[2.3, 7.7],
[2.4, 6.3],
[1.8, 6.4],
[1.8, 6. ],
[2.1, 6.9],
[2.4, 6.7],
[2.3, 6.9],
[1.9, 5.8],
[2.3, 6.8],
[2.5, 6.7],
[2.3, 6.7],
[1.9, 6.3],
[2. , 6.5],
[2.3, 6.2],
[1.8, 5.9]])
```

1.3 Importing YOUR data

- For the coursework output 2, you will need to import the data from a **.csv** file.
- For instance, the IRIS dataset would look something like this:
- Your datasets will have a **first column** with the id of each entry (**NOT** the same as the row index).
- Your dataset will have the class/target in the **last column**.

	A	B	C	D	E	F
1	flower_id	sepal_len	sepal_wid	petal_len	petal_wid	variety
2	45	5.1	3.5	1.4	0.2	0
3	88	4.9	3	1.4	0.2	0
4	100	4.7	3.2	1.3	0.2	0
5	133	4.6	3.1	1.5	0.2	0
6	160	5	3.6	1.4	0.2	0

dataset.png

- The **first row** contains the header.
- You need to find a pre-existing module that lets you import data from a csv file into a numpy array.
- Try to import the header in a different variable as the data.
- Since the classes/targets are numeric for all datasets, you can leave them on the same numpy array as the data.
- You don't need the target names, just work with the numbers!