

**Victor Hugo Avila Felipe - A01794425**



## **Data Analysis with Python**

**IBM: DA0101EN**

### **LEARNING OBJECTIVES**

In this course you will learn about:

- Data Acquisition
- How to Obtain Basic Insight From a Dataset
- Data Wrangling
- Exploratory Data Analysis
- Model Development
- Model Evaluation

## **Introduction to Data Analysis with Python**

- Problem requiring data analysis
- dataset to analyze in python
- overview of packages
- import and export data
- Basic insights

Can we estimate the price of used cars?

### **The problem**

why data analysis? data everywhere, helps discovery of information.

Tom wants to sell his car, but wants the best price. what affects the price?

### **Understand the data**

There are documentation and the .csv file.

The first attribute, "symboling", corresponds to the insurance risk level of a car. Cars are initially assigned a risk factor symbol associated with their price. Then, if an automobile is more risky, this symbol is adjusted by moving it up the scale. A value of +3 indicates that the auto is risky, -3 that it is probably pretty safe. The second attribute "normalized-losses" is the relative average loss payment per insured vehicle year. This value is normalized for all autos within a particular size classification (two-door small, station wagons, sports/speciality, etc...), and represents the average loss per car per year. The values range from 65 to 256. The other attributes are easy to understand.

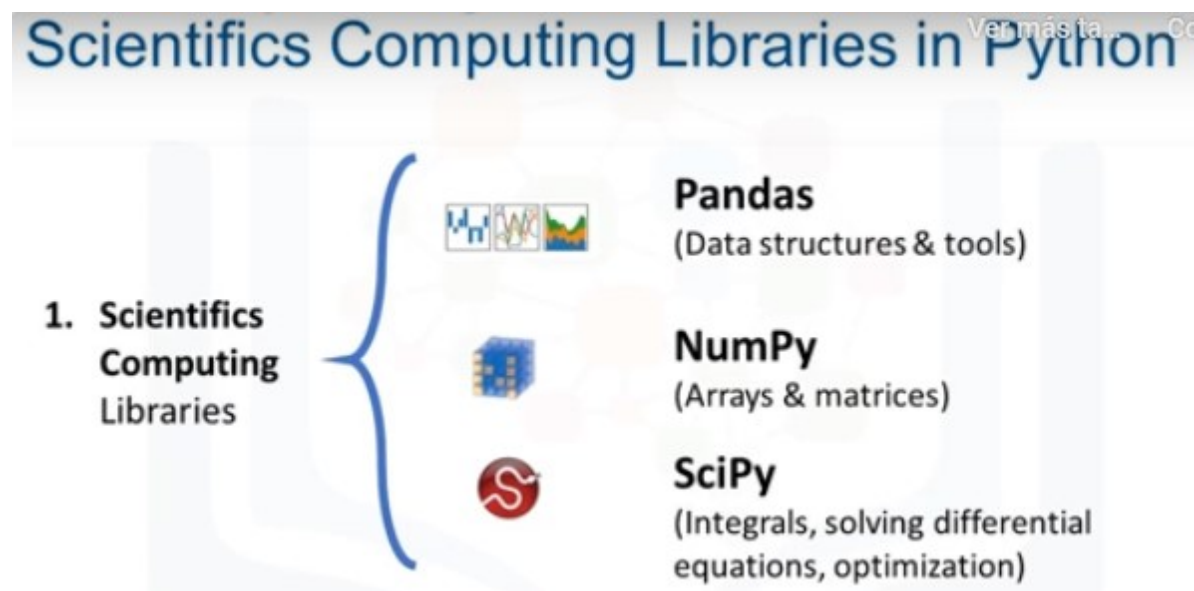
**Thus, the goal of this project is to predict "price" in terms of other car features.**

No.	Attribute name	attribute range	No.	Attribute name	attribute range
1	symboling	-3, -2, -1, 0, 1, 2, 3.	14	curb-weight	continuous from 1488 to 4066.
2	normalized-losses	continuous from 65 to 256.	15	engine-type	dohc, dohc, l, ohc, ohc, ohcv, rotor.
3	make	audi, bmw, etc.	16	num-of-cylinders	eight, five, four, six, three, twelve, two.
4	fuel-type	diesel, gas.	17	engine-size	continuous from 61 to 326.
5	aspiration	std, turbo.	18	fuel-system	1bbl, 2bbl, 4bbl, idi, mfi, mpfi, spdi, spfi.
6	num-of-doors	four, two.	19	bore	continuous from 2.54 to 3.94.
7	body-style	hardtop, wagon, etc.	20	stroke	continuous from 2.07 to 4.17.
8	drive-wheels	4wd, fwd, rwd.	21	compression-ratio	continuous from 7 to 23.
9	engine-location	front, rear.	22	horsepower	continuous from 48 to 288.
10	wheel-base	continuous from 86.6 to 120.9.	23	peak-rpm	continuous from 4150 to 6600.
11	length	continuous from 141.1 to 208.1.	24	city-mpg	continuous from 13 to 49.
12	width	continuous from 60.3 to 72.3.	25	highway-mpg	continuous from 16 to 54.
13	height	continuous from 47.8 to 59.8.	26	price	continuous from 5118 to 45400.

Target (Label)

## Python packages

A Python library is a collection of functions and methods that allow you to perform lots of actions without writing any code. The libraries usually contain built-in modules providing different functionalities, which you can use directly. And there are extensive libraries, offering a broad range of facilities.



## Visualization Libraries in Python

### 2. Visualization Libraries



#### Matplotlib

(plots & graphs, most popular)



#### Seaborn

(plots : heat maps, time series, violin plots)

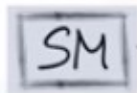
## Algorithmic Libraries in Python

### 3. Algorithmic libraries



#### Scikit-learn

(Machine Learning : regression, classification,... )



#### Statsmodels

(Explore data, estimate statistical models, and perform statistical tests.)

### importing and exporting data

Format and file path. `read_csv()` `df` `df.head(n)` `df.tail(n)`

Add headers `df.columns = headers` `headers = ["a","b","c"]`

export `df` to csv: `df.to_csv(path)`

csv, json, excel, sql

### Analyzing data

Check data type data distribution

object, float, int y datetime

- potential info and type mismatch
- compatibility with python methods

dataframe.dtypes

df.describe(include="all") -> count, mean, std deviation, min, 25%, 50%, 75%, max all ->  
UNIQUE, top, freq

## Data lab

- Data source: <https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data> ([https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data?utm\\_medium=Exinfluencer&utm\\_source=Exinfluencer&utm\\_content=000026UJ&utm\\_term=10006555&utm\\_id=NA-SkillsNetwork-Channel-SkillsNetworkCoursesIBMDeveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-0](https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data?utm_medium=Exinfluencer&utm_source=Exinfluencer&utm_content=000026UJ&utm_term=10006555&utm_id=NA-SkillsNetwork-Channel-SkillsNetworkCoursesIBMDeveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-0))
- Data type: csv

In [10]: *#install specific version of libraries used in Lab*

```
import sys
!{sys.executable} -m pip install pandas
!{sys.executable} -m pip install numpy
!{sys.executable} -m pip install matplotlib
!{sys.executable} -m pip install scipy
!{sys.executable} -m pip install seaborn
!{sys.executable} -m pip install ipywidgets
```

```
Requirement already satisfied: pandas in /home/flynn/anaconda3/lib/python3.9/
site-packages (1.4.2)
Requirement already satisfied: python-dateutil>=2.8.1 in /home/flynn/anaconda
3/lib/python3.9/site-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /home/flynn/anaconda3/lib/pyth
on3.9/site-packages (from pandas) (2021.3)
Requirement already satisfied: numpy>=1.18.5 in /home/flynn/anaconda3/lib/pyt
hon3.9/site-packages (from pandas) (1.21.5)
Requirement already satisfied: six>=1.5 in /home/flynn/anaconda3/lib/python3.
9/site-packages (from python-dateutil>=2.8.1->pandas) (1.16.0)
Requirement already satisfied: numpy in /home/flynn/anaconda3/lib/python3.9/s
ite-packages (1.21.5)
Requirement already satisfied: matplotlib in /home/flynn/anaconda3/lib/python
3.9/site-packages (3.5.1)
Requirement already satisfied: kiwisolver>=1.0.1 in /home/flynn/anaconda3/lib
/python3.9/site-packages (from matplotlib) (1.3.2)
Requirement already satisfied: numpy>=1.17 in /home/flynn/anaconda3/lib/pytho
n3.9/site-packages (from matplotlib) (1.21.5)
Requirement already satisfied: fonttools>=4.22.0 in /home/flynn/anaconda3/lib
/python3.9/site-packages (from matplotlib) (4.25.0)
```

In [13]:

```
Input In [13]
  pip install pandas
    ^
SyntaxError: invalid syntax
```

```
In [11]: # import pandas library
import pandas as pd
```

In [12]: *#This function will download the dataset into your browser*

```
from pyodide.http import pyfetch

async def download(url, filename):
    response = await pyfetch(url)
    if response.status == 200:
        with open(filename, "wb") as f:
```

```
-----
ModuleNotFoundError                                Traceback (most recent call last)
Input In [12], in <cell line: 3>()
      1 #This function will download the dataset into your browser
----> 3 from pyodide.http import pyfetch
      5 async def download(url, filename):
      6     response = await pyfetch(url)

ModuleNotFoundError: No module named 'pyodide'
```

## Read Data

We use `pandas.read_csv()` function to read the csv file. In the brackets, we put the file path along with a quotation mark so that pandas will read the file into a dataframe from that address. The file path can be either an URL or your local file address.

Because the data does not include headers, we can add an argument `headers = None` inside the `read_csv()` method so that pandas will not automatically set the first row as a header.

You can also assign the dataset to any variable you create.

```
In [ ]: path = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM"
```

```
In [ ]: #you will need to download the dataset; if you are running locally, please com
await download(path, "auto.csv")
```

This dataset was hosted on IBM Cloud object. Click [HERE \(https://cocl.us\)](https://cocl.us)

[/DA101EN\\_object\\_storage?utm\\_medium=Exinfluencer&utm\\_source=Exinfluencer&utm\\_content=000026UJ&utm\\_term=10006555&utm\\_id=NA-SkillsNetwork-Channel-SkillsNetworkCoursesIBMDDeveloperSkillsNetworkDA0101ENSkillsNetwork20235326-2021-01-01](https://www.skillsnetwork.co.uk/courses/bmdeveloper/skillsnetworkda0101enskillsnetwork20235326-2021-01-01)  
for free storage

## Import pandas library

import pandas as pd

## Read the online file by the URL provides above, and assign it to variable "df"

```
df = pd.read_csv(path, header=None)
```

After reading the dataset, we can use the `dataframe.head(n)` method to check the top `n` rows of the dataframe, where `n` is an integer. Contrary to `dataframe.head(n)`, `dataframe.tail(n)` will show you the bottom `n` rows of the dataframe.

```
In [ ]: # show the first 5 rows using dataframe.head() method
print("The first 5 rows of the dataframe")
```

```
In [ ]: print("The last 10 rows of the dataframe\n")
```

```
In [ ]: # create headers list
headers = ["symboling", "normalized-losses", "make", "fuel-type", "aspiration", "n",
           "drive-wheels", "engine-location", "wheel-base", "length", "width", "heig",
           "num-of-cylinders", "engine-size", "fuel-system", "bore", "stroke", "comp",
           "peak-rpm", "city-mpg", "highway-mpg", "price"]
```

```
In [ ]: df.columns = headers
```

```
In [ ]: df1=df.replace('?',np.NaN)
```

```
In [ ]: df=df1.dropna(subset=["price"], axis=0)
```

```
In [ ]:
```

## Save Dataset

Correspondingly, Pandas enables us to save the dataset to csv. By using the `dataframe.to_csv()` method, you can add the file path and name along with quotation marks in the brackets.

For example, if you would save the dataframe **df** as **automobile.csv** to your local machine, you may use the syntax below, where `index = False` means the row names will not be written.

In [ ]:

## Read/Save Other Data Formats

Data Formate	Read	Save
csv	pd.read_csv()	df.to_csv()
json	pd.read_json()	df.to_json()
excel	pd.read_excel()	df.to_excel()
hdf	pd.read_hdf()	df.to_hdf()
sql	pd.read_sql()	df.to_sql()
...	...	...

In [ ]:

In [ ]:

```
# check the data type of data frame "df" by .dtypes
```

In [ ]:

In [ ]:

In [ ]:

```
# describe all the columns in "df"
```

In [ ]:

In [ ]:

In [ ]:

```
# Look at the info of "df"
```

### Question 1

1/1 point (graded)

What does CSV stand for?

☒ Comma-separated values

☐ Car sold values

☐ Car state values

☐ None of the above


Save

Submit

You have used 1 of 2 attempts

### Question 2

0/1 point (graded)

[Return to Question 1](#)
[Return to Question 3](#)
[Return to Question 4](#)
[Return to Question 5](#)

In the data set, which of the following represents an attribute or feature?

☐ Row

☐ Column

☒ Each element in the dataset



Submit

You have used 2 of 2 attempts

### Question 3

1/1 point (graded)

What is the name of what we want to predict?

☒ Target

☐ Feature

☐ Dataframe



Save

Submit

You have used 1 of 2 attempts

### Question 4

1/1 point (graded)

What is the command to display the first five rows of a dataframe `df` ?

☒ `df.head()`

☐ `df.tail()`



Submit

You have used 1 of 1 attempt

### Question 5

1/1 point (graded)

What command do you use to get the data type of each row of the dataframe `df` ?

☒ `df.dtypes`

☐ `df.head()`

☐ `df.tail()`



Save

Submit

You have used 1 of 2 attempts

### Question 6

1/1 point (graded)

How do you get a statistical summary of a dataframe `df` ?

☒ `df.describe()`



☐ `df.head()`☐ `df.tail()`[Save](#)**Submit**

You have used 1 of 2 attempts

### Question 7

1/1 point (graded)

If you use the method `describe()` without changing any of the arguments, you will get a statistical summary of all the columns of type "object".

☒ False☐ True**Submit**

You have used 1 of 1 attempt

In [ ]:

In [ ]:

In [ ]:

In [ ]: