**Maestría en Inteligencia Artificial Aplicada**

**Materia: Ciencia y analítica de datos**

**Profesor: Jobish Vallikavungal Devassia**

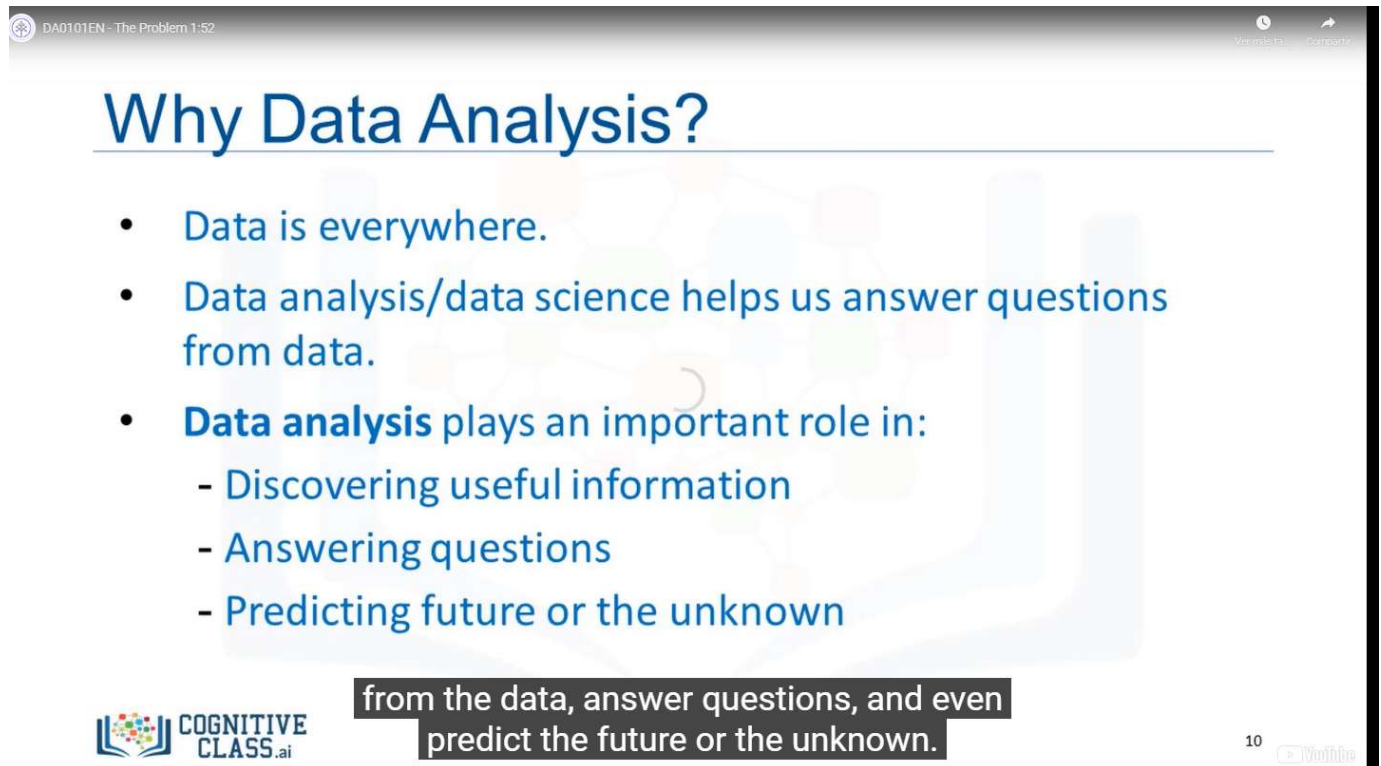**Semana 4: Data Analysis with Python (IBM) : Módulo 1**

**Integrante:**

- **Norma de los Ángeles García López - A01228278**

**Fecha de entrega: 11 de octubre del 2022**

# ▾ MÓDULO 1

Introduction to Data Analysis with Python



*Problem: Estimate how much money can Tom get from the car?*

Some Hypothesis questions

Atributos del dataset



The target value is "price" which is the one that we are going to try to estimate with out model, the other variables are the predictors

*Python Packages for Data Science*

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.

# Scientifics Computing Libraries in Python

**1. Scientifics Computing Libraries**

**Pandas** (Data structures & tools)

**NumPy** (Arrays & matrices)

**SciPy** (Integrals, solving differential equations, optimization)

SciPy includes functions for some advanced math problems, as listed on this slide, as

COGNITIVE CLASS.ai

5

## ▾ *Data visualization Methods*

Using data visualization methods is the best way to communicate with others, showing them meaningful results of analysis. These libraries enable you to create graphs, charts and maps.

# Visualization Libraries in Python

**2. Visualization Libraries**

**Matplotlib** (plots & graphs, most popular)

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

It's very easy to generate various plots such as heat maps, time series, and violin

COGNITIVE CLASS.ai

8

## ▾ *Algorithmic Libraries in Python*

With Machine Learning algorithms, we're able to develop a model using our dataset, and obtain predictions.



## ▾ *Importing and Exporting Data in Python*

Data acquisition is a process of loading and reading data into notebook from various sources. To read any data using Python's pandas package, there are two important factors to consider:

**Format is the way data is encoded.** We can usually tell different encoding schemes by looking at the ending of the file name. Some common encodings are csv, json, xlsx, hdf and so forth.

**The (file) path tells us where the data is stored.** Usually it is stored either on the computer we are using, or online on the internet.



```
import pandas as pd

url = "https://archive.ics.uci.edu/ml/machine-learningdatabases/autos/imports-85.data"

df = pd.read_csv(url)
```

# Importing a CSV without a header

```python
import pandas as pd

url = "https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data"

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

DA0101EN Importing and Exporting Data in Python v2

# Printing the dataframe in Python

- **df** prints the entire dataframe (not recommended for large datasets)
- **df.head(n)** to show the first *n* rows of data frame.
- **df.tail(n)** shows the bottom *n* rows of data frame.

df.head()

Header

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 3 | ? | alfa-romero | gas | std | two | convertible | rwd | front | 88.6 | ... | 130 | mpfi | 3.47 | 2.68 | 9.0 | 111 | 5000 | 21 | 27 | 13495 |
| 1 | 3 | ? | alfa-romero | gas | std | two | convertible | rwd | front | 88.6 | ... | 130 | mpfi | 3.47 | 2.68 | 9.0 | 111 | 5000 | 21 | 27 | 16500 |
| 2 | 1 | ? | alfa-romero | gas | std | two | hatchback | rwd | front | 94.5 | ... | 152 | mpfi | 2.68 | 3.47 | 9.0 | 154 | 5000 | 19 | 26 | 16500 |
| 3 | 2 | 164 | audi | gas | std | four | sedan | fwd | front | 99.8 | ... | 109 | mpfi | 3.19 | 3.40 | 10.0 | 102 | 5500 | 24 | 30 | 13950 |
| 4 | 2 | 164 | audi | gas | std | four | sedan | 4wd | front | 99.4 | ... | 136 | mpfi | 3.19 | 3.40 | 8.0 | 115 | 5500 | 18 | 22 | 17450 |

COGNITIVE CLASS

23

DA0101EN Importing and Exporting Data in Python v2

# Adding headers

- Replace default header (by `df.columns = headers`)

```
headers = ["symboling","normalized-losses","make","fuel-type","aspiration","num-of-doors","body-style",
"drive-wheels","engine-location","wheel-base","length","width","height","curb-weight","engine-type",
"num-of-cylinders","engine-size","fuel-system","bore","stroke","compression-ratio","horsepower","peak-
rpm","city-mpg","highway-mpg","price"]
```

```
df.columns=headers
```

```
df.head(5)
```

| Data Format | Read | Save |
|---|---|---|
| csv | pd.read_csv() | df.to_csv() |
| json | pd.read_json() | df.to_json() |
| Excel | pd.read_excel() | df.to_excel() |
| sql | pd.read_sql() | df.to_sql() |

▾ *Getting Started Analyzing Data in Python*

# Basic insights from the data

- Understand your data before you begin any analysis

- Should check:
    - Data Types
    - Data Distribution

- Locate potential issues with the data

# Basic Insights of Dataset - Data Types

| Pandas Type | Native Python Type | Description |
|---|---|---|
| object | string | numbers and strings |
| int64 | int | Numeric characters |
| float64 | float | Numeric characters with decimals |
| datetime64, timedelta[ns] | N/A (but see the datetime module in Python's standard library) | time data. |

## Why check data types?

- potential info and type mismatch

- compatibility with python methods

When the "dtype" method is applied to the data set, the datatype of each column is returned in a Series. A good data scientist's intuition tells us that most of the data types make sense. The make of cars, for example, are names, so this information should be of type object.

# Basic Insights of Dataset - Data Types

- In pandas, we use `dataframe.dtypes` to check data types

`df.dtypes`

| | |
|---|---|
| symboling | int64 |
| normalized-losses | object |
| make | object |
| fuel-type | object |
| aspiration | object |
| num-of-doors | object |
| body-style | object |
| drive-wheels | object |
| engine-location | object |
| wheel-base | float64 |
| length | float64 |
| width | float64 |
| height | float64 |
| curb-weight | int64 |
| engine-type | object |
| num-of-cylinders | object |
| engine-size | int64 |
| fuel-system | object |
| bore | object |
| stroke | object |
| compression-ratio | float64 |
| horsepower | object |
| peak-rpm | object |
| city-mpg | int64 |
| highway-mpg | int64 |
| price | object |
| dtype: object | |

Now we would like to check the statistical summary of each column to learn about the distribution of data in each column. The statistical metrics can tell the data scientist if there are mathematical issues that may exist, such as extreme outliers and large deviations. The data scientist may have to address these issues later. To get the quick statistics, we use the describe method. It returns the number of terms in the column as "count", average column value as "mean", column standard deviation as "std", the maximum and minimum values, as well as the boundary of each of the quartiles.

# dataframe.describe(include="all")

- Provides full summary statistics

`df.describe(include="all")`

| | symboling | normalized-losses | make | fuel-type | aspiration | num-of-doors | body-style | drive-wheels | engine-location | wheel-base | ... | engine-size | fuel-system | bore | stroke |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 205.000000 | 205 | 205 | 205 | 205 | 205 | 205 | 205 | 205 | 205.000000 | ... | 205.000000 | 205 | 205 | 205 |
| unique | NaN | 52 | 22 | 2 | 2 | 3 | 5 | 3 | 2 | NaN | ... | NaN | 8 | 39 | 37 |
| top | NaN | ? | toyota | gas | std | four | sedan | fwd | front | NaN | ... | NaN | mpfi | 3.62 | 3.40 |
| freq | NaN | 41 | 32 | 185 | 168 | 114 | 96 | 120 | 202 | NaN | ... | NaN | 94 | 23 | 20 |
| mean | 0.834146 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | 98.756585 | ... | 126.907317 | NaN | NaN | NaN |
| std | 1.245307 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | 6.021776 | ... | 41.642693 | NaN | NaN | NaN |
| min | -2.000000 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | 86.600000 | ... | 61.000000 | NaN | NaN | NaN |
| 25% | 0.000000 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | 94.500000 | ... | 97.000000 | NaN | NaN | NaN |
| 50% | 1.000000 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | 97.000000 | ... | 120.000000 | NaN | NaN | NaN |
| 75% | 2.000000 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | 102.400000 | ... | 141.000000 | NaN | NaN | NaN |
| max | 3.000000 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | 120.900000 | ... | 326.000000 | NaN | NaN | NaN |

21

"Unique" is the number of distinct objects in the column, "top" is the most frequently occurring object, and "freq" is the number of times the top object appears in the column. Some values in the table are shown here as "NaN", which stands for "not a number". This is because that particular statistical metric cannot be calculated for that specific column data type.



Another method you can use to check your dataset is the dataframe.info function. This function shows the top 30 rows and bottom 30 rows of the dataframe.

## ▾ Graded Review Questions

Question 1

1/1 punto (calificado)

What does CSV stand for?

- ⦿ Comma-separated values
- ◯ Car sold values
- ◯ Car state values
- ◯ None of the above

✔

Guardar

**Enviar**   Ha realizado 1 de 2 intentos

✔ Correcto (1/1 punto)

Question 2

1 punto posible (calificable)

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

- ◯ Row
- ◯ Column
- ◯ Each element in the dataset

Guardar

Enviar   Ha realizado 0 de 2 intentos

Question 3

---

## Question 3

1/1 punto (calificado)

What is the name of what we want to predict?

- ⦿ Target
- ◯ Feature
- ◯ Dataframe

✔

Guardar

**Enviar**   Ha realizado 1 de 2 intentos

✔ Correcto (1/1 punto)

## Question 4

1/1 punto (calificado)

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

- ⦿ `df.head()`
- ◯ `df.tail()`

✔

Enviar   Ha realizado 1 de 1 intento

Question 5

1/1 punto (calificado)

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

- ◉ `df.dtypes`
- ○ `df.head()`
- ○ `df.tail()`

✔

Guardar

[ Enviar ]　Ha realizado 1 de 2 intentos

✔ Correcto (1/1 punto)

Question 6

1/1 punto (calificado)

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

- ◉ `df.describe()`
- ○ `df.head()`
- ○ `df.tail()`

✔

Guardar

[ Enviar ]　Ha realizado 1 de 2 intentos
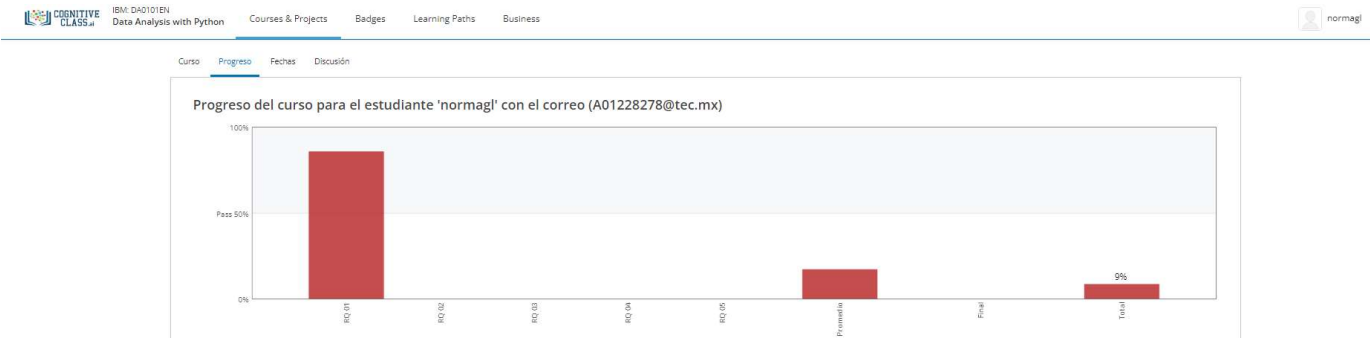
✔ Correcto (1/1 punto)

Question 7

1/1 punto (calificado)

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

✔

[ Enviar ]　Ha realizado 1 de 1 intento

✔ Correcto (1/1 punto)

▾ PROGRESO

COGNITIVE CLASS.ai    IBM: DA0101EN
Data Analysis with Python    Courses & Projects    Badges    Learning Paths    Business                    normagl

Curso    Progreso    Fechas    Discusión

Progreso del curso para el estudiante 'normagl' con el correo (A01228278@tec.mx)



Productos de pago de Colab  -  Cancelar contratos

COGNITIVE CLASS.ai    IBM: DA0101EN
Data Analysis with Python    Courses & Projects    Badges    Learning Paths    Business

Curso    Progreso    Fechas    Discusión