12:19 PM Fri 14 Jan •••











Not Trusted Python 3.8 o







Introduction Notebook

Estimated time needed: 10 minutes

Objectives

After completing this lab you will be able to:

- · Acquire data in various ways
- · Obtain insights from data with Pandas library

Table of Contents

- 1. Data Acquisition
- 2. Basic Insight of Dataset

Data Acquisition

There are various formats for a dataset: .csv, .json, .xlsx etc. The dataset can be stored in different places, on your local machine or sometimes online.

In this section, you will learn how to load a dataset into our Jupyter Notebook.

In our case, the Automobile Dataset is an online source, and it is in a CSV (comma separated value) format. Let's use this dataset as an example to practice data reading.

- Data source: https://archive.ics.uci.edu/ml/machine-learning-databases/autos/imports-85.data
- Data type: csv

The Pandas Library is a useful tool that enables us to read various datasets into a dataframe; our Jupyter notebook platforms have a built-in **Pandas Library** so that all we need to do is import Pandas without installing.

```
In [1]: *#install specific version of libraries used in lab
! mamba install pandas==1.3.3 -y
! mamba install numpy=1.21.2 -y

/usr/bin/sh: mamba: command not found
/usr/bin/sh: mamba: command not found

In [2]: *# import pandas library
import pandas as pd
import numpy as np
```

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.

This dataset was hosted on IBM Cloud object. Click HERE for free storage.

```
In [3]: ** Import pandas library
import pandas as pd

# Read the online file by the URL provides above, and assign it to variable "df"
other_path = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-DA0101EN-St
df = pd.read_csv(other_path, header=None)
```



jp-tok.dataplatform.cloud.ibm.com 🔒 🕞 AA ŌŌ

Not Trusted Python 3.8 o Insert Cell Kernel Help View ↑ U Run • C → Format Markdown

In [4]: other_path

Out[4]: 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-DA0101EN-SkillsNetwo rk/labs/Data%20files/auto.csv'

In [5]: **▶** df

Out [5]: 0 2 16 17 18 19 20 21 std two convertible rwd front 88.6 ... 130 mpfi 3.47 2.68 9.0 111 5000 21 27 13495 ? alfa-romero gas std two convertible rwd front 88.6 ... 130 mpfi 3.47 2.68 9.0 111 5000 21 27 16500 ? alfa-romero gas ? alfa-romero std two hatchback rwd front 94.5 ... 152 mpfi 2.68 3.47 9.0 154 5000 19 26 16500 **3** 2 164 std four sedan fwd front 99.8 ... 109 mpfi 3.19 3.40 10.0 102 5500 24 30 13950 audi gas sedan 4wd front 99.4 ... 136 mpfi 3.19 3.40 8.0 115 5500 18 22 17450 **4** 2 164 std four gas audi sedan rwd front 109.1 ... 141 mpfi 3.78 3.15 9.5 114 5400 23 28 16845 **200** -1 95 volvo std four sedan rwd front 109.1 ... 141 mpfi 3.78 3.15 8.7 160 5300 19 25 19045 **201** -1 95 gas turbo four volvo 202 -1 95 std four sedan rwd front 109.1 ... 173 mpfi 3.58 2.87 8.8 134 5500 18 23 21485 volvo gas sedan rwd front 109.1 ... 145 idi 3.01 3.40 23.0 106 4800 26 27 22470 95 volvo diesel turbo four **203** -1 **204** -1 95 gas turbo four sedan rwd front 109.1 ... 141 mpfi 3.78 3.15 9.5 114 5400 19 25 22625

205 rows × 26 columns

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

The first 5 rows of the dataframe

volvo

Out[6]: 0 9 ... 16 17 18 19 ? alfa-romero gas std two convertible rwd front 88.6 ... 130 mpfi 3.47 2.68 9.0 111 5000 21 27 13495 ? alfa-romero gas std two convertible rwd front 88.6 ... 130 mpfi 3.47 2.68 9.0 111 5000 21 27 16500 ? alfa-romero gas std two hatchback rwd front 94.5 ... 152 mpfi 2.68 3.47 9.0 154 5000 19 26 16500 audi gas std four sedan fwd front 99.8 ... 109 mpfi 3.19 3.40 10.0 102 5500 24 30 13950 **3** 2 164 **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

5 rows × 26 columns

Question #1:

Check the bottom 10 rows of data frame "df".

In [8]: * # Write your code below and press Shift+Enter to execute print("Bottom 10 rows of dataframe df are ") df.tail(10)

Bottom 10 rows of dataframe df are

Out[8]: 8 9 ... 16 17 18 19 20 21 std four wagon rwd front 104.3 ... 141 mpfi 3.78 3.15 9.5 114 5400 23 28 13415 **195** -1 74 volvo gas **196** -2 103 volvo std four sedan rwd front 104.3 ... 141 mpfi 3.78 3.15 9.5 114 5400 24 28 15985 **197** -1 74 volvo std four wagon rwd front 104.3 ... 141 mpfi 3.78 3.15 9.5 114 5400 24 28 16515 gas turbo four sedan rwd front 104.3 ... 130 mpfi 3.62 3.15 7.5 162 5100 17 22 18420 **198** -2 103 volvo **199** -1 74 volvo gas turbo four wagon rwd front 104.3 ... 130 mpfi 3.62 3.15 7.5 162 5100 17 22 18950 **200** -1 95 volvo gas std four sedan rwd front 109.1 ... 141 mpfi 3.78 3.15 9.5 114 5400 23 28 16845 gas turbo four sedan rwd front 109.1 ... 141 mpfi 3.78 3.15 8.7 160 5300 19 25 19045 gas std four sedan rwd front 109.1 ... 173 mpfi 3.58 2.87 8.8 134 5500 18 23 21485 95 volvo 203 -1 95 volvo diesel turbo four sedan rwd front 109.1 ... 145 idi 3.01 3.40 23.0 106 4800 26 27 22470 204 -1 95 volvo gas turbo four sedan rwd front 109.1 ... 141 mpfi 3.78 3.15 9.5 114 5400 19 25 22625

10 rows × 26 columns

12:19 PM Fri 14 Jan **?** 22% **□** $\bullet \bullet \bullet$





jp-tok.dataplatform.cloud.ibm.com 🔒 🗁 💍







Not Trusted Python 3.8 o

```
Edit View Insert Cell Kernel Help
```

Add Headers

Take a look at our dataset. Pandas automatically set the header with an integer starting from 0.

To better describe our data, we can introduce a header. This information is available at: https://archive.ics.uci.edu/ml/datasets/Automobile.

Thus, we have to add headers manually.

First, we create a list "headers" that include all column names in order. Then, we use dataframe.columns = headers to replace the headers with the list we created.

```
In [9]: ▶ # create headers list
          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"]
          print("headers\n", 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', 'highwa y-mpg', 'price']

We replace headers and recheck our dataframe:

```
In [10]:  df.columns = headers
          df.head(10)
```

Out[10]:	symboling	normalized- losses	make	fuel- type	aspiration	num- of- doors			engine- location		engine- " size	fuel- system	bore	stroke	compression- ratio	horsepowe
	0 3	?	alfa- romero	gas	std	two	convertible	rwd	front	88.6 .	130	mpfi	3.47	2.68	9.0	111
	1 3	?	alfa- romero	gas	std	two	convertible	rwd	front	88.6 .	130	mpfi	3.47	2.68	9.0	111
	2 1	?	alfa- romero	gas	std	two	hatchback	rwd	front	94.5 .	152	mpfi	2.68	3.47	9.0	154
	3 2	164	audi	gas	std	four	sedan	fwd	front	99.8 .	109	mpfi	3.19	3.40	10.0	102
	4 2	164	audi	gas	std	four	sedan	4wd	front	99.4 .	136	mpfi	3.19	3.40	8.0	115
	-				-4-1	4		£	£ ±	00.0	100	t:	0.40	0.40	0.5	440

```
In [11]: • df.columns
```

```
Out[11]: Index(['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'],
               dtype='object')
```

We need to replace the "?" symbol with NaN so the dropna() can remove the missing values:

We can drop missing values along the column "price" as follows:

```
In [13]:  df=df1.dropna(subset=["price"], axis=0)
          df.head(20)
```

4111	1044 (20)															
out[13]:	symboling	normalized- losses	make	fuel- type	aspiration	num- of- doors			engine- location	wheel- base	engine- size	fuel- system	bore	stroke	compression- ratio	horsepo
	0 3	NaN	alfa- romero	gas	std	two	convertible	rwd	front	88.6 .	130	mpfi	3.47	2.68	9.0	
	1 3	NaN	alfa- romero	gas	std	two	convertible	rwd	front	88.6 .	130	mpfi	3.47	2.68	9.0	
	2 1	NaN	alfa- romero	gas	std	two	hatchback	rwd	front	94.5 .	152	mpfi	2.68	3.47	9.0	
	3 2	164	audi	gas	std	four	sedan	fwd	front	99.8 .	109	mpfi	3.19	3.40	10.0	
	4 2	164	audi	gas	std	four	sedan	4wd	front	99.4 .	136	mpfi	3.19	3.40	8.0	
	-	NI-NI	I!		-4-1	4		£	f	00.0	100	£	0.10	0.40	0.5	

12:19 PM Fri 14 Jan **?** 22% **□** $\bullet \bullet \bullet$



AA

jp-tok.dataplatform.cloud.ibm.com 🔒 🕃 💍











```
Not Trusted Python 3.8 o
 Edit View Insert Cell Kernel Help
 ⊕ ♣ 🖟 1 ♠ ♦ Run • ♦ Format Markdown →
In [15]: • df.shape
```

Question #2:

Out[15]: (201, 26)

Find the name of the columns of the dataframe.

```
In [16]: ** # Write your code below and press Shift+Enter to execute
          df.columns
    Out[16]: Index(['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'],
```

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.

```
df.to_csv("automobile.csv", index=False)
```

We can also read and save other file formats. We can use similar functions like pd.read_csv() and df.to_csv() for other data formats. The functions are listed in the following table:

Read/Save Other Data Formats

dtype='object')

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

Basic Insight of Dataset

After reading data into Pandas dataframe, it is time for us to explore the dataset.

float64

int64

object

obiect

There are several ways to obtain essential insights of the data to help us better understand our dataset.

Data Types

height

curb-weight

engine-type

num-of-cvlinders

Data has a variety of types.

The main types stored in Pandas dataframes are object, float, int, bool and datetime64. In order to better learn about each attribute, it is always good for us to know the data type of each column. In Pandas:

```
In [17]: ▶ df.dtypes
    Out[17]: symboling
                                      int64
             normalized-losses
                                    object
             make
                                     object
             fuel-type
                                     object
             aspiration
                                     object
             num-of-doors
                                     object
             body-style
                                     object
             drive-wheels
                                     object
                                     object
             engine-location
                                    float64
             wheel-base
                                    float64
              length
             width
                                    float64
```

12:19 PM Fri 14 Jan **?** 22% **■** $\bullet \bullet \bullet$



AA

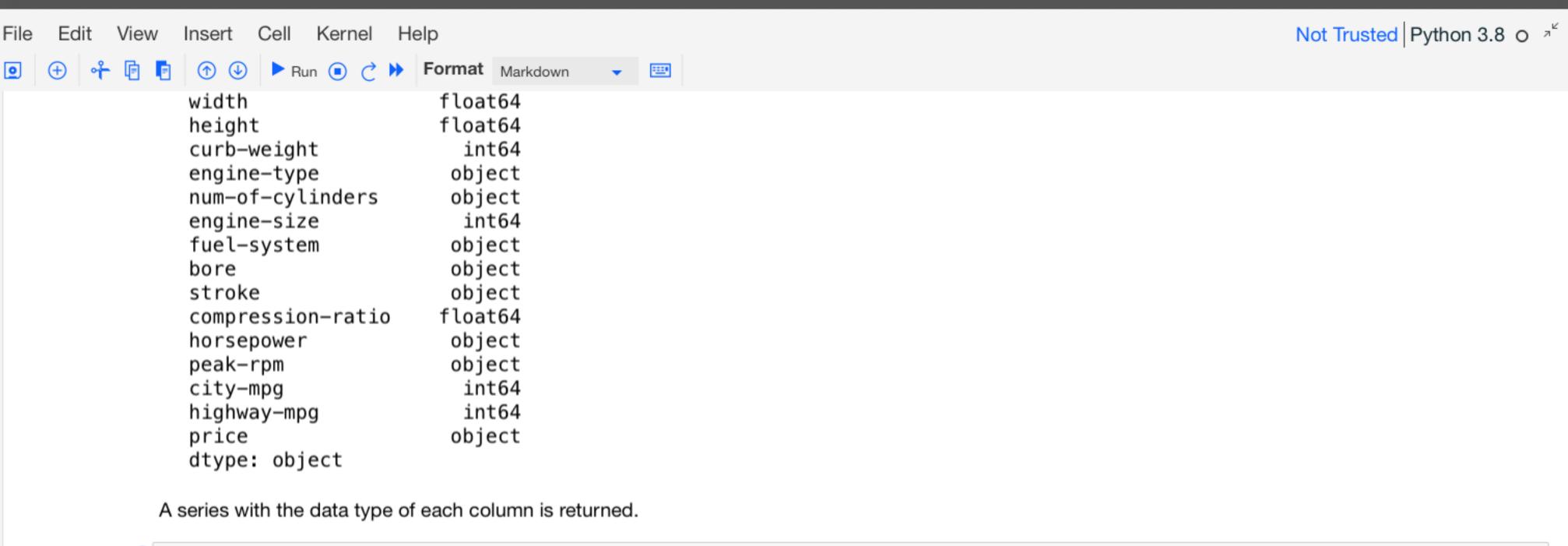
jp-tok.dataplatform.cloud.ibm.com 🔒 🕃 💍











In [18]: ** # check the data type of data frame "df" by .dtypes print(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 float64 length width float64 height float64 curb-weight int64 engine-type object num-of-cylinders object engine-size int64 fuel-system object bore object stroke object float64 compression-ratio object horsepower peak-rpm object int64 city-mpg int64 highway-mpg object price dtype: object

As shown above, it is clear to see that the data type of "symboling" and "curb-weight" are int64, "normalized-losses" is object, and "wheel-base" is float64, etc.

These data types can be changed; we will learn how to accomplish this in a later module.

Describe

If we would like to get a statistical summary of each column e.g. count, column mean value, column standard deviation, etc., we use the describe method:

dataframe.describe()

This method will provide various summary statistics, excluding NaN (Not a Number) values.

In [19]: df.describe()

_			
0u	tΙ	19	1:

	symboling	wheel-base	length	width	height	curb-weight	engine-size	compression-ratio	city-mpg	highway-mpg
count	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000	201.000000
mean	0.840796	98.797015	174.200995	65.889055	53.766667	2555.666667	126.875622	10.164279	25.179104	30.686567
std	1.254802	6.066366	12.322175	2.101471	2.447822	517.296727	41.546834	4.004965	6.423220	6.815150
min	-2.000000	86.600000	141.100000	60.300000	47.800000	1488.000000	61.000000	7.000000	13.000000	16.000000
25%	0.000000	94.500000	166.800000	64.100000	52.000000	2169.000000	98.000000	8.600000	19.000000	25.000000
50%	1.000000	97.000000	173.200000	65.500000	54.100000	2414.000000	120.000000	9.000000	24.000000	30.000000
75%	2.000000	102.400000	183.500000	66.600000	55.500000	2926.000000	141.000000	9.400000	30.000000	34.000000
max	3.000000	120.900000	208.100000	72.000000	59.800000	4066.000000	326.000000	23.000000	49.000000	54.000000

12:19 PM Fri 14 Jan **?** 22% **□** ...

Ringer

AA

↑ U Run • ↑ Format Markdown

jp-tok.da

5

Not Trusted Python 3.8 o Insert Cell Kernel Help View

This shows the statistical summary of all numeric-typed (int, float) columns.

For example, the attribute "symboling" has 205 counts, the mean value of this column is 0.83, the standard deviation is 1.25, the minimum value is -2, 25th percentile is 0, 50th percentile is 1, 75th percentile is 2, and the maximum value is 3.

However, what if we would also like to check all the columns including those that are of type object?

You can add an argument include = "all" inside the bracket. Let's try it again.

In [20]: • # describe all the columns in "df" df.describe(include = "all")

Out[20]:

	symboling	normalized- losses	make	fuel- type	aspiration	num- of- doors	body- style	drive- wheels	engine- location	wheel- base	 engine- size	fuel- system	bore	stroke	compression- ratio	hor
count	201.000000	164	201	201	201	199	201	201	201	201.000000	 201.000000	201	197	197	201.000000	
unique	NaN	51	22	2	2	2	5	3	2	NaN	 NaN	8	38	36	NaN	
top	NaN	161	toyota	gas	std	four	sedan	fwd	front	NaN	 NaN	mpfi	3.62	3.40	NaN	
freq	NaN	11	32	181	165	113	94	118	198	NaN	 NaN	92	23	19	NaN	
mean	0.840796	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	98.797015	 126.875622	NaN	NaN	NaN	10.164279	
std	1.254802	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	6.066366	 41.546834	NaN	NaN	NaN	4.004965	
min	-2.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	86.600000	 61.000000	NaN	NaN	NaN	7.000000	
25%	0.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	94.500000	 98.000000	NaN	NaN	NaN	8.600000	
50%	1.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	97.000000	 120.000000	NaN	NaN	NaN	9.000000	
75%	2.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	102.400000	 141.000000	NaN	NaN	NaN	9.400000	
max	3.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	120.900000	 326.000000	NaN	NaN	NaN	23.000000	

11 rows × 26 columns

Now it provides the statistical summary of all the columns, including object-typed attributes.

We can now see how many unique values there, which one is the top value and the frequency of top value in the object-typed columns.

Some values in the table above show as "NaN". This is because those numbers are not available regarding a particular column type.

Question #3:

You can select the columns of a dataframe by indicating the name of each column. For example, you can select the three columns as follows:

dataframe[[' column 1 ',column 2', 'column 3']]

Where "column" is the name of the column, you can apply the method ".describe()" to get the statistics of those columns as follows:

dataframe[[' column 1 ',column 2', 'column 3']].describe()

Apply the method to ".describe()" to the columns 'length' and 'compression-ratio'.

In [21]: • # Write your code below and press Shift+Enter to execute df[['length','compression-ratio']].describe()

Out[21]:

	length	compression-ratio
count	201.000000	201.000000
mean	174.200995	10.164279
std	12.322175	4.004965
min	141.100000	7.000000
25%	166.800000	8.600000
50%	173.200000	9.000000
75%	183.500000	9.400000
max	208.100000	23.000000

Info

Another method you can use to check your dataset is:

dataframe.info()

It provides a concise summary of your DataFrame.

This method prints information about a DataFrame including the index dtype and columns, non-null values and memory usage.

In [22]: • # look at the info of "df"



Edit View

AΑ

Insert Cell Kernel Help

jp-tok.dataplatform.cloud.ibm.com 🔒 🗁 💍









Not Trusted Python 3.8 o



```
₩ ===
In [22]: • # look at the info of "df"
          df.info()
             <class 'pandas.core.frame.DataFrame'>
             Int64Index: 201 entries, 0 to 204
             Data columns (total 26 columns):
                  Column
                                     Non-Null Count
                                                    Dtype
                  symboling
                                                     int64
                                     201 non-null
                  normalized-losses
                                     164 non-null
                                                     object
                  make
                                     201 non-null
                                                     object
                                     201 non-null
                                                     object
                  fuel-type
                                     201 non-null
                                                     object
                  aspiration
                  num-of-doors
                                     199 non-null
                                                     object
                  body-style
                                     201 non-null
                                                     object
                  drive-wheels
                                     201 non-null
                                                     object
                  engine-location
                                     201 non-null
                                                     object
                  wheel-base
                                     201 non-null
                                                     float64
              10
                  length
                                     201 non-null
                                                     float64
                                                     float64
              11
                  width
                                     201 non-null
              12
                                                     float64
                  height
                                     201 non-null
              13
                                                     int64
                  curb-weight
                                     201 non-null
                                     201 non-null
              14
                  engine-type
                                                     object
                  num-of-cylinders
                                     201 non-null
                                                     object
                  engine-size
                                     201 non-null
                                                     int64
              16
                                     201 non-null
              17
                  fuel-system
                                                     object
                                     197 non-null
              18
                                                     object
                  bore
                                     197 non-null
                  stroke
                                                     object
                                                     float64
                  compression-ratio
                                     201 non-null
                                                     object
                                     199 non-null
              21 horsepower
              22 peak-rpm
                                     199 non-null
                                                     object
                 city-mpg
```

Excellent! You have just completed the Introduction Notebook!

201 non-null

201 non-null

201 non-null

dtypes: float64(5), int64(5), object(16)

int64

int64

object

Thank you for completing this lab!

memory usage: 42.4+ KB

24 highway-mpg

25 price

Author

Joseph Santarcangelo

Other Contributors

Mahdi Noorian PhD

Bahare Talayian

Eric Xiao

Steven Dong

Parizad

Hima Vasudevan

Fiorella Wenver

Yi Yao.

Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2020-10-30	2.3	Lakshmi	Changed URL of the csv
2020-09-22	2.2	Nayef	Added replace() method to remove '?'
2020-09-09	2.1	Lakshmi	Made changes in info method of dataframe
2020-08-27	2.0	Lavanya	Moved lab to course repo in GitLab