



Introduction to Pandas in Python

Estimated time needed: 15 minutes

Objectives

After completing this lab you will be able to:

- Use Pandas to access and view data

Table of Contents

- About the Dataset
- Introduction of Pandas
- Viewing Data and Accessing Data
- Quiz on DataFrame

About the Dataset

The table has one row for each album and several columns.

- artist: Name of the artist
- album: Name of the album
- released_year: Year the album was released
- length_min_sec: Length of the album (hours,minutes,seconds)
- genre: Genre of the album
- music_recording_sales_millions: Music recording sales (millions in USD) on [SONG//DATABASE]
- claimed_sales_millions: Album's claimed sales (millions in USD) on [SONG//DATABASE]
- date_released: Date on which the album was released
- soundtrack: Indicates if the album is the movie soundtrack (Y or N)
- rating_of_friends: Indicates the rating from your friends from 1 to 10

You can see the dataset here:

```
# Load the data
# Read the file
# Print first five rows of the data frame
```

Introduction of Pandas

In [1]: `# Dependency needed to install file`

`!pip install xlrd`

Requirement already satisfied: xlrd in c:\python\lib\site-packages (2.0.1)

In [2]: `# Import required library`

`import pandas as pd`

After the import command, we now have access to a large number of pre-built classes and functions. This assumes the library is installed; in our lab environment all the necessary libraries are installed. One way pandas allows you to work with data is a data frame. Let's go through the process to go from a comma separated values (.csv) file to a data frame. This variable `csv_path` stores the path of the .csv, that is used as an argument to the `read_csv` function. The result is stored in the object `df`, this is a common short form used for a variable referring to a Pandas data frame.

In [3]: `# Read data from CSV file`

`csv_path = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-PY0101EN-SkillsNetwork/labs/M02/Lab1/musics.csv'`

We can use the method `head()` to examine the first five rows of a data frame:

In [4]: `# Print first five rows of the data frame`

`df.head()`

	Artist	Album	Released	Length	Genre	Music Recording Sales (millions)	Claimed Sales (millions)	Released.1	Soundtrack	Rating
0	Michael Jackson	Thriller	1982	0:42:19	pop, rock, R&B	46.0	65	30-Nov-82	NaN	10.0
1	AC/DC	Back in Black	1980	0:42:11	hard rock	26.1	50	25-Jul-80	NaN	9.5
2	Pink Floyd	The Dark Side of the Moon	1973	0:42:49	progressive rock	24.2	45	01-Mar-73	NaN	9.0
3	Whitney Houston	The Bodyguard	1992	0:57:44	R&B, soul, pop	27.4	44	17-Nov-92	Y	8.5
4	Meat Loaf	Bat Out of Hell	1977	0:46:33	hard rock, progressive rock	20.6	43	21-Oct-77	NaN	8.0

We use the path of the excel file and the function `read_excel`. The result is a data frame as before:

In [5]: `# Read data from Excel File and print the first five rows`

`xlsx_path = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-PY0101EN/Chapter%204/Datasets/TopSellingAlbums.xlsx'`

`df = pd.read_excel(xlsx_path)`

We can use the method `head()` to examine the first five rows of a data frame:

In [6]: `# Access to the column Length`

`x = df[['Length']]`

`x`

Out[6]: `Length`

`0 0:42:19`

`1 0:42:11`

`2 0:42:49`

`3 0:57:44`

`4 0:46:33`

`5 0:43:08`

`6 0:15:54`

`7 0:40:01`

The process is shown in the figure:

`x=df[['Length']]`

Artist	Album	Released	Length	Genre	Music Recording Sales (millions)	Claimed Sales (millions)	Released.1	Soundtrack	Rating
0 Michael Jackson	Thriller	1982	0:42:19	pop, rock, R&B	46.0	65	30-Nov-82	NaN	10.0
1 AC/DC	Back in Black	1980	0:42:11	hard rock	26.1	50	25-Jul-80	NaN	9.5
2 Pink Floyd	The Dark Side of the Moon	1973	0:42:49	progressive rock	24.2	45	01-Mar-73	NaN	9.0
3 Whitney Houston	The Bodyguard	1992	0:57:44	R&B, soul, pop	27.4	44	17-Nov-92	Y	8.5
4 Meat Loaf	Bat Out of Hell	1977	0:46:33	hard rock, progressive rock	20.6	43	21-Oct-77	NaN	8.0
5 Eagles	Their Greatest Hits (1971-1975)	1976	0:43:08	rock, soft rock, folk rock	32.2	42	17-Feb-76	NaN	7.5
6 Bee Gees	Saturday Night Fever	1977	1:15:54	disco	20.6	40	15-Nov-77	Y	7.0
7 Fleetwood Mac	Rumours	1977	0:40:01	soft rock	27.9	40	04-Feb-77	NaN	6.5

The process is shown in the figure:

`y=df[['Artist','Length','Genre']]`

Artist	Album	Released	Length	Genre	Music Recording Sales (millions)	Claimed Sales (millions)	Released.1	Soundtrack	Rating
0 Michael Jackson	Thriller	1982	0:42:19	pop, rock, R&B	46.0	65	30-Nov-82	NaN	10.0
1 AC/DC	Back in Black	1980	0:42:11	hard rock	26.1	50	25-Jul-80	NaN	9.5
2 Pink Floyd	The Dark Side of the Moon	1973	0:42:49	progressive rock	24.2	45	01-Mar-73	NaN	9.0
3 Whitney Houston	The Bodyguard	1992	0:57:44	R&B, soul, pop	27.4	44	17-Nov-92	Y	8.5
4 Meat Loaf	Bat Out of Hell	1977	0:46:33	hard rock, progressive rock	20.6	43	21-Oct-77	NaN	8.0
5 Eagles	Their Greatest Hits (1971-1975)	1976	0:43:08	rock, soft rock, folk rock	32.2	42	17-Feb-76	NaN	7.5
6 Bee Gees	Saturday Night Fever	1977	1:15:54	disco	20.6	40	15-Nov-77	Y	7.0
7 Fleetwood Mac	Rumours	1977	0:40:01	soft rock	27.9	40	04-Feb-77	NaN	6.5

One way to access unique elements is the `iloc` method, where you can access the 1st row and the 1st column as follows:

In [10]: `# Access the value on the first row and the first column`

`df.iloc[0, 0]`

'Michael Jackson'

You can access the 2nd row and the 1st column as follows:

In [11]: `# Access the value on the second row and the first column`

`df.iloc[1, 0]`

'AC/DC'

You can access the 1st row and the 3rd column as follows:

In [12]: `# Access the value on the first row and the third column`

`df.iloc[0, 2]`

1982

This is shown in the following image

Artist	Album	Released	Length	Genre	Music Recording Sales (millions)	Claimed Sales (millions)	Released.1	Soundtrack	Rating
0 Michael Jackson	Thriller	1982	0:42:19	pop, rock, R&B	46.0	65	30-Nov-82	NaN	10.0
1 AC/DC	Back in Black	1980	0:42:11	hard rock	26.1	50	25-Jul-80	NaN	9.5
2 Pink Floyd	The Dark Side of the Moon	1973	0:42:49	progressive rock	24.2	45	01-Mar-73	NaN	9.0
3 Whitney Houston	The Bodyguard	1992	0:57:44	R&B, soul, pop	27.4	44	17-Nov-92	Y	8.5
4 Meat Loaf	Bat Out of Hell	1977	0:46:33	hard rock, progressive rock	20.6	43	21-Oct-77	NaN	8.0
5 Eagles	Their Greatest Hits (1971-1975)	1976	0:43:08	rock, soft rock, folk rock	32.2	42	17-Feb-76	NaN	7.5
6 Bee Gees	Saturday Night Fever	1977	1:15:54	disco	20.6	40	15-Nov-77	Y	7.0
7 Fleetwood Mac	Rumours	1977	0:40:01	soft rock	27.9	40	04-Feb-77	NaN	6.5

The process is shown in the figure:

`y=df[['Artist','Length','Genre']]`

Artist	Album	Released	Length	Genre	Music Recording Sales (millions)	Claimed Sales (millions)	Released.1	Soundtrack	Rating