

Python Exercises

Part 1:

Video Game Sales

We have a file with video game sales data with the following columns:

- Ranking: Ranking of the sale of the video game worldwide
- Name: Name of the video game
- Platform: Platform on which it was launched (Wii, PS4, PC ... etc)
- Year: Launch year
- Genre: Video game genre
- Publisher: Company that published the game
- NA_Sales - Sales in North America (in millions)
- EU_Sales - Sales in Europe (in millions)
- JP_Sales - Sales in Japan (in millions)
- Other_Sales - Sales in the rest of the world (in millions)
- Global_Sales - Worldwide Sales (in millions)

Loading the dataset:

```
In [1]: import pandas as pd
from scipy import stats

games = pd.read_csv('vg-sales.csv')

games.head()
```

```
Out[1]:
```

Rank	Name	Platform	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales
0	1	Wii Sports	Wii 2006.0	Sports	Nintendo	41.49	29.02	3.77	8.46	82.74
1	2	Super Mario Bros.	NES 1985.0	Platform	Nintendo	29.08	3.58	6.81	0.77	40.24
2	3	Mario Kart Wii	Wii 2008.0	Racing	Nintendo	15.85	12.88	3.79	3.31	35.82
3	4	Wii Sports Resort	Wii 2009.0	Sports	Nintendo	15.75	11.01	3.28	2.96	33.00
4	5	Pokemon Red/Pokemon Blue	GB 1996.0	Role-Playing	Nintendo	11.27	8.89	10.22	1.00	31.37

Finding the video game that has generated the most revenue globally

```
In [2]: games.sort_values(['Global_Sales'], ascending=False).head(1)['Name']

Out[2]:
```

Wii Sports
Name: Name, dtype: object

Finding the platform with the most average global sales

```
In [3]: plataformas = games.groupby('Platform')
plataformas.mean().sort_values('Global_Sales', ascending=False).head(1)
```

```
Out[3]:
```

Platform	Rank	Year	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales
GB	3392.030612	1995.958763	1.166531	0.487959	0.868571	0.083673	2.606633

Finding the publisher with the most global sales

```
In [4]: empresas = games.groupby('Publisher')
empresas.sum().sort_values('Global_Sales', ascending=False).head(1)
```

```
Out[4]:
```

Publisher	Rank	Year	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales
Nintendo	2714164	1394666.0	816.87	418.74	455.42	95.33	1786.56

Discovering which year had the most global sales

```
In [5]: años = games.groupby('Year')
años.sum().sort_values('Global_Sales', ascending=False).head(1)
```

```
Out[5]:
```

Year	Rank	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales
2008.0	1208807	351.44	184.4	60.26	82.39	678.9

Finding the top selling game for each platform

```
In [6]: años = games.groupby('Platform', 'Name', 'Global_Sales')
años.max().sort_values('Global_Sales', ascending=False).groupby('Platform').head(1)
```

```
Out[6]:
```

Platform	Name	Global_Sales	Rank	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales
Wii	Wii Sports	82.74	1	2006.0	Sports	Nintendo	41.49	29.02	3.77	8.46
NES	Super Mario Bros.	40.24	2	1985.0	Platform	Nintendo	29.08	3.58	6.81	0.77
GB	Pokemon Red/Pokemon Blue	31.37	5	1996.0	Role-Playing	Nintendo	11.27	8.89	10.22	1.00
DS	New Super Mario Bros.	30.01	7	2006.0	Platform	Nintendo	11.38	9.23	6.50	2.90
X360	Kinect Adventures!	21.82	16	2010.0	Misc	Microsoft Game Studios	14.97	4.94	0.24	1.67
PS2	Grand Theft Auto V	21.81	17	2013.0	Action	Take-Two Interactive	7.01	9.27	0.97	4.14
PS3	Grand Theft Auto San Andreas	20.40	18	2004.0	Action	Take-Two Interactive	9.43	4.40	0.41	10.57
SNES	Super Mario World	20.61	19	1990.0	Platform	Nintendo	12.78	3.75	3.54	0.55
GBA	Ruby/Pokemon Sapphire	15.85	26	2002.0	Role-Playing	Nintendo	6.06	3.90	5.38	0.50
3DS	Pokemon X/Pokemon Y	14.35	33	2013.0	Role-Playing	Nintendo	5.17	4.05	4.34	0.79
PS4	Call of Duty: Black Ops 3	14.24	34	2015.0	Shooter	Activision	5.77	5.81	0.35	2.31
N64	Super Mario 64	11.89	47	1996.0	Platform	Nintendo	6.91	2.85	1.91	0.23
PS	Gran Turismo	10.95	53	1997.0	Racing	Sony Computer Entertainment	4.02	3.87	2.54	0.52
XB	Halo 2	8.49	80	2004.0	Shooter	Microsoft Game Studios	6.82	1.53	0.05	0.08
PC	The Sims	8.11	84	2009.0	Simulation	Electronic Arts	0.98	6.42	0.00	0.71
PC60	Pac-Man	7.81	90	1982.0	Puzzle	Atari	7.28	0.45	0.00	0.08
PSP	Grand Theft Auto: Liberty City Stories	7.72	91	2005.0	Action	Take-Two Interactive	2.90	2.83	0.24	1.75
GC	Call of Duty: Black Ops 3	7.70	102	2015.0	Shooter	Activision	4.52	2.09	0.01	0.67
XONE	Super Smash Bros. Melee	7.07	108	2001.0	Fighting	Nintendo	4.41	1.04	1.39	0.22
WiiU	Mario Kart 8	6.96	109	2014.0	Racing	Nintendo	3.13	2.07	1.27	0.49
GEN	Sonic the Hedgehog 2	6.03	144	1992.0	Platform	Sega	4.47	1.20	0.16	0.19
DC	Sonic Adventure	2.42	638	1998.0	Platform	Sega	1.26	0.61	0.46	0.08
PSV	Minecraft	2.25	715	2014.0	Misc	Sony Computer Entertainment Europe	0.28	0.79	0.87	0.32
SCD	Virtua Fighter 2	1.93	890	1995.0	Fighting	Sega	0.34	0.26	1.30	0.03
SAT	Sonic CD	1.50	1263	1993.0	Platform	Sega	1.00	0.36	0.09	0.05
WS	Final Fantasy	0.51	3933	2000.0	Role-Playing	Squaresoft	0.00	0.00	0.51	0.00
NG	Samurai Shodown II	0.25	6683	1994.0	Fighting	SNK	0.00	0.00	0.25	0.00
TG16	Pokeyunsei	0.14	9225	1995.0	Adventure	NEC	0.00	0.00	0.14	0.00
DO	Doukenuuts	0.06	12627	1995.0	Adventure	Konami Digital Entertainment	0.00	0.00	0.06	0.00
GG	Sonic the Hedgehog 2 (GBA)	0.04	13527	1992.0	Platform	Sega	0.00	0.00	0.04	0.00
PCFX	Blue Breaker/Ken Vlorio/Hokemito	0.03	14559	1996.0	Role-Playing	NEC	0.00	0.00	0.03	0.00

Finding the top seller for each genre

```
In [7]: años = games.groupby('Genre', 'Name', 'Global_Sales')
años.max().sort_values('Global_Sales', ascending=False).groupby('Genre').head(1)
```

```
Out[7]:
```

Genre	Name	Global_Sales	Rank	Platform	Year	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales
Sports	Wii Sports	82.74	1	Wii 2006.0	Nintendo	41.49	29.02	3.77	8.46	
Platform	Super Mario Bros.	40.24	2	NES 1985.0	Nintendo	29.08	3.58	6.81	0.77	
Racing	Mario Kart Wii	35.82	3	Wii 2008.0	Nintendo	15.85	12.88	3.79	3.31	
Role-Playing	Pokemon Red/Pokemon Blue	31.37	5	GB 1996.0	Nintendo	11.27	8.89	10.22	1.00	
Puzzle	Tetris	30.26	6	GB 1989.0	Nintendo	23.20	2.26	4.22	0.58	
Misc	Wii Play	29.02	8	Wii 2006.0	Nintendo	14.03	9.20	2.93	2.85	
Shooter	Duck Hunt	28.31	10	NES 1984.0	Nintendo	26.93	0.63	0.28	0.47	
Simulation	Nintendogs	24.76	11	DS 2005.0	Nintendo	9.07	11.00	1.93	2.75	
Action	Grand Theft Auto V	21.81	17	PS3 2013.0	Take-Two Interactive	7.01	9.27	0.97	4.14	
Fighting	Super Smash Bros. Brawl	13.04	40	Wii 2008.0	Nintendo	6.75	2.61	2.96	1.02	
Adventure	Super Mario Land 2: 6 Golden Coins	11.18	51	GB 1992.0	Nintendo	6.16	2.04	2.69	0.29	
Strategy	Pokemon Stadium	5.45	166	N64 1999.0	Nintendo	3.18	1.24	0.94	0.09	

Creating a function that groups video game release dates by decades, and shows how have evolved their global sales

```
In [8]: import pandas as pd
from scipy import stats

def decades(row):
    year = row['Year']
    if 1979 <= year < 1989:
        return "1980s"
    elif 1989 <= year < 2000:
        return "1990s"
    elif 1999 <= year < 2010:
        return "2000s"
    elif 2009 <= year < 2018:
        return "2010s"
    else:
        return "Others"

games['Decade'] = games.apply(decades, axis=1)
games.head(20)
```

```
Out[8]:
```

Rank	Name	Platform	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales	Decade
0	1	Wii Sports	Wii 2006.0	Sports	Nintendo	41.49	29.02	3.77	8.46	82.74	2000s
1	2	Super Mario Bros.	NES 1985.0	Platform	Nintendo	29.08	3.58	6.81	0.77	40.24	1980s
2	3	Mario Kart Wii	Wii 2008.0	Racing	Nintendo	15.85	12.88	3.79	3.31	35.82	2000s
3	4	Wii Sports Resort	Wii 2009.0	Sports	Nintendo	15.75	11.01	3.28	2.96	33.00	2000s
4	5	Pokemon Red/Pokemon Blue	GB 1996.0	Role-Playing	Nintendo	11.27	8.89	10.22	1.00	31.37	1990s
5	6	Tetris	GB 1989.0	Puzzle	Nintendo	23.20	2.26	4.22	0.58	30.26	1980s
6	7	New Super Mario Bros.	DS 2006.0	Platform	Nintendo	11.38	9.23	6.50	2.90	30.01	2000s
7	8	Wii Play	Wii 2006.0	Misc	Nintendo	14.03	9.20	2.93	2.85	29.02	2000s
8	9	New Super Mario Bros. Wii	Wii 2009.0	Platform	Nintendo	14.59	7.06	4.70	2.26	28.62	2000s
9	10	Duck Hunt	NES 1984.0	Shooter	Nintendo	26.93	0.63	0.28	0.47	28.31	1980s
10	11	Nintendogs	DS 2005.0	Simulation	Nintendo	9.07	11.00	1.93	2.75	24.76	2000s
11	12	Mario Kart DS	DS 2005.0	Racing	Nintendo	9.81	7.57	4.13	1.92	23.42	2000s
12	13	Pokemon Gold/Pokemon Silver	GB 1999.0	Role-Playing	Nintendo	9.00	6.18	7.20	0.71	23.10	1990s
13	14	Wii Fit	Wii 2007.0	Sports	Nintendo	8.94	8.03	3.60	2.15	22.72	2000s
14	15	Wii Fit Plus	Wii 2009.0	Sports	Nintendo	9.09	8.59	2.53	1.79	22.00	2000s
15	16	Kinect Adventures!	X360 2010.0	Misc	Microsoft Game Studios	14.97	4.94	0.24	1.67	21.82	2010s
16	17	Grand Theft Auto V	PS3 2013.0	Action	Take-Two Interactive	7.01	9.27	0.97	4.14	21.40	2010s
17	18	Grand Theft Auto San Andreas	PS2 2004.0	Action	Take-Two Interactive	9.43	0.40	0.41	10.57	20.81	2000s
18	19	Super Mario World	SNES 1990.0	Platform	Nintendo	12.78	3.75	3.54	0.55	20.61	1990s
19	20	Brain Age: Train Your Brain in Minutes a Day	DS 2005.0	Misc	Nintendo	4.75	9.26	4.16	2.05	20.22	2000s

Part 2

5000 IMD Movies (The Internet Movie Database)

In this project we will focus on the exploratory data analysis of a 5000 movie dataset obtained from the Internet Movie Database (<https://www.imdb.com>).

```
In [9]: import pandas as pd
import numpy as np
import datetime

df = pd.read_csv('movie_metadata.csv')
```

```
Out[9]:
```

df.dtypes

```
Out[11]:
```

	color	director_name	num_critic_for_reviews	duration	director_facebook_likes	actor_3_facebook_likes	actor_2_name	actor_1_facebook_likes
0	Color	James Cameron	723.0	178.0	0.0	855.0	Joel David Moore	1000.0
1	Color	Gore Verbinski	302.0	169.0	563.0	1000.0	Orlando Bloom	4000.0
2	Color	Sam Mendes	602.0	148.0	0.0	161.0	Rory Kinnear	1100.0
3	Color	Christopher Nolan	813.0	164.0	22000.0	23000.0	Christian Bale	2700.0
4	NaN	Doug Walker	NaN	NaN	131.0	NaN	Rob Walker	131.0

5 rows x 28 columns

Finding the difference between budget and total box office per director on an annual basis by adding a new column called diff_gross.

```
In [13]: limpio = df.copy()
limpio['diff_gross'] = limpio['gross'] - limpio['budget']
limpio.head()
```

```
Out[14]:
```

color	director_name	num_critic_for_reviews	duration	director_facebook_likes	actor_3_facebook_likes	actor_2_name	actor_1_facebook_likes	
0	Color	James Cameron	723.0	178.0	0.0	855.0	Joel David Moore	1000.0
1	Color	Gore Verbinski	302.0	169.0	563.0	1000.0	Orlando Bloom	4000.0
2	Color	Sam Mendes	602.0	148.0	0.0	161.0	Rory Kinnear	1100.0
3	Color	Christopher Nolan	813.0	164.0	22000.0	23000.0	Christian Bale	2700.0
5	Color	Andrew Stanton	462.0	132.0	475.0	530.0	Samantha Morton	640.0

5 rows x 29 columns

```
In [15]: directors = limpio.groupby('director_name')
directors.sum().sort_values('diff_gross', ascending=False).head()
```

```
Out[15]:
```

director_name	num_critic_for_reviews	duration	director_facebook_likes	actor_3_facebook_likes	actor_1_facebook_likes	gross	num_voted
Steven Spielberg	6526.0	3429.0	350000.0	42157.0	271942.0	4.114233e+09	8
George Lucas	1345.0	655.0	0.0	9929.0	62000.0	1.741418e+09	2
James Cameron	1878.0	1098.0	0.0	4303.0	38780.0	1.730876e+09	3
Jos Whedon	2317.0	606.0	0.0	57173.0	80000.0	1.618708e+09	2
Chris Columbus	1567.0	1398.0	0.0	17183.0	161870.0	1.618708e+09	1

Finding which directors have obtained the highest number of critic reviews

```
In [16]: directors = limpio.groupby('director_name')
directors.sum().sort_values('num_critic_for_reviews', ascending=False).head()
```

```
Out[16]:
```

director_name	num_critic_for_reviews	duration	director_facebook_likes	actor_3_facebook_likes	actor_1_facebook_likes	gross	num_voted
Steven Spielberg	6526.0	3429.0	350000.0	42157.0	271942.0	4.114233e+09	8
Ridley Scott	4930.0	2419.0	0.0	10167.0	187713.0	1.377722e+09	4
Tim Burton	4200.0	1759.0	20800.0	46100.0	328705.0	2.071275e+09	3
Clint Eastwood	4172.0	250.0	30400.0	6753.0	266874.0	1.378321e+09	3
Christopher Nolan	4090.0	1122.0	17600.0	93698.0	150716.0	1.813228e+09	8

IMDB is known for the disproportion between the evaluation of professional critics and that of the users of the platform. High values mean that professional critics value the film more and vice versa. I will be finding this ratio and adding it in a new column called critic_ratio.

```
In [17]: limpio['critic_ratio'] = limpio['num_critic_for_reviews'] / limpio['num_user_for_reviews']
limpio.head()
```

```
Out[17]:
```

color	director_name	num_critic_for_reviews	duration	director_facebook_likes	actor_3_facebook_likes	actor_2_name	actor_1_facebook_likes	
0	Color	James Cameron	723.0	178.0	0.0	855.0	Joel David Moore	1000.0
1	Color	Gore Verbinski	302.0	169.0	563.0	1000.0	Orlando Bloom	4000.0
2	Color	Sam Mendes	602.0	148.0	0.0	161.0	Rory Kinnear	1100.0
3	Color	Christopher Nolan	813.0	164.0	22000.0	23000.0	Christian Bale	2700.0
5	Color	Andrew Stanton	462.0	132.0	475.0	530.0	Samantha Morton	640.0

5 rows x 30 columns