

I T P 4 4 9

Data Exploration

Using Pandas and DataFrames

Lecture 6



More with NumPy, Pandas

```
[[19 19 16 10]
 [17 11 15 15]
 [12 10 19 10]
 [10 16 13 16]
 [11 13 13 15]
 [15 10 10 16]
 [12 19 14 19]
 [10 16 11 12]]
```

```
[[19 19]
 [17 11]
 [12 10]
 [10 16]]
```

```
[[11 13]
 [15 10]
 [12 19]
 [10 16]]
```

```
[[16 10]
 [15 15]
 [19 10]
 [13 16]]
```

```
[[13 15]
 [10 16]
 [14 19]
 [11 12]]
```

Do the following:

Split a (8, 4) shaped matrix of random integers, ranging from 10 to 20, into 4 arrays with shapes (4, 2)

The screenshot shows a PyCharm IDE window with the following code in the editor:

```
1 import numpy as np
2
3 x = np.random.randint(low=10, high=20, size=(8, 4))
4
5 [xL, xR] = np.hsplit(x, 2)
6 [xLT, xLB] = np.vsplit(xL, 2)
7 [xRT, xRB] = np.vsplit(xR, 2)
8
9 print(x, '\n\n', xLT, '\n\n', xLB, '\n\n', xRT, '\n\n', xRB)
```

The code imports numpy, generates an 8x4 array 'x' with integer values from 10 to 20, and then splits it into four sub-arrays: 'xL' and 'xR' (horizontal split), and 'xLT', 'xLB', 'xRT', and 'xRB' (vertical split). The final line prints all these arrays with newlines between them.

`hsplit` splits the array horizontally, i.e side by side.

`vsplit` splits the array vertically, i.e top and bottom

```
Run: in_class_coding ×  C:\Users\Reza\Desktop\ITP449_Fall2020\class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Cl  
[  
    [[10 18 19 11]  
     [19 11 14 10]  
     [14 19 13 15]  
     [14 11 11 11]  
     [10 14 18 14]  
     [19 19 14 17]  
     [13 11 13 18]]
```

1: Project

2: Structure

3: Favorites

4: Run

Pr... in_class_coding.py

```
1 import numpy as np
2
3 a = np.array([5, 3, -12, np.nan, 6, np.nan, 45, 1])
4
5 b = a[np.isnan(a) == False]
# Alternative
6 c = a[~np.isnan(a)]
7
8
9 print('a =', a)
10 print('b =', b)
11 print('c =', c)
12
```

nan = not a number – missing data

~ tilde is the complement (inverse) operator. It is a bitwise inverse.

Run: in_class_coding

```
C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Cl
a = [ 5.  3. -12. nan  6. nan 45.  1.]
b = [ 5.  3. -12.  6. 45.  1.]
c = [ 5.  3. -12.  6. 45.  1.]
```

Process finished with exit code 0

5: TODO

6: Problems

Terminal

Python Console

Event Log

Do the following:

Create a Series of length 20 that contains random integers between 1 and 10. Then print the number of integers in the Series that are greater than or equal to 5.

ITP449_Fall2020 > Class > In Class Coding >  in_class_coding.py

 in_class_coding ▾



1: Project

7: Structure

★ 2: Favorites

in_class_coding.py

The screenshot shows the PyCharm IDE interface. On the left, the project structure is displayed under the 'ITP449_Fall2020' root folder. The 'Class' folder contains sub-folders like '.idea', 'Files', 'In Class Coding', and 'in_class_code'. A 'venv' folder is also present. The 'External Libraries' and 'Scratches and Console' sections are visible at the bottom. The main code editor window on the right contains the following Python script:

```
1 import pandas as pd
2 import numpy as np
3
4 a = pd.Series(np.random.randint(low=1, high=11, size=20))
5
6 print(a.values)
7 print('There are', len(a[a.values >= 5]), 'numbers greater than or equal to 5')
8
9
10
```

Run: in_class_coding >



```
C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Cl
[ 5  5  7  6  6  5  1  2  4  7  4  5  6  6 10  3  1  2  2  3]
There are 11 numbers greater than or equal to 5
Process finished with exit code 0
```

4: Run

TODO

6: Problem

> Terminal

Python Consol

Event Log

Do the following:

Create the DataFrame shown below and filter to show only the rows for desserts that are worth the calories.

	dessert	worthIt
0	cupcake	False
1	donut	True
2	cake	True
3	ice cream	True

1: Project

2: Structure

2: Favorites

Pr... - in_class_coding.py ×

```
1 import pandas as pd
2
3 data = {'dessert': ['cupcake', 'donut', 'cake', 'ice cream'],
4         'worthIt': [False, True, True, True]}
5
6 df = pd.DataFrame(data)
7
8 print(df[df['worthIt'].values])
```

⚠ 1 ✘ 131 ⌂

Run: in_class_coding ×

C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/CL

dessert worthIt

1	donut	True
2	cake	True
3	ice cream	True

Process finished with exit code 0

Run

TODO

⚠ 6: Problems

Terminal

Python Console

Event Log

Do the following:

Create a DataFrame with the following data on the highest grossing movies:

	ReleaseYear		Movie	Expense	Profit
0	2009		Avatar	328262000	1499182307
1	2019		Avengers: Endgame	475560000	1010189087
2	2015	Star Wars Ep VII: The Force Awakens	Awakens	381704000	843123707
3	2013		Frozen	245904000	797747501
4	2005	HP and the Goblet of Fire		208064000	789505652

Pr... in_class_coding.py x

1 import pandas as pd
2
3 data = {'ReleaseYear': [2009, 2019, 2015, 2013, 2005],
4 'Movie': ['Avatar', 'Avengers: Endgame', 'Star Wars Ep VII: The Force Awakens',
5 'Frozen', 'HP and the Goblet of Fire'],
6 'Expense': [328262000, 475560000, 381704000, 245904000, 208064000],
7 'Profit': [1499182307, 1010189087, 843123707, 797747501, 789505652]}
8
9 df = pd.DataFrame(data)
10 pd.set_option('display.max_columns', None)
11 print(df)

Run: in_class_coding x

	ReleaseYear	Movie	Expense	Profit
0	2009	Avatar	328262000	1499182307
1	2019	Avengers: Endgame	475560000	1010189087
2	2015	Star Wars Ep VII: The Force Awakens	381704000	843123707
3	2013	Frozen	245904000	797747501
4	2005	HP and the Goblet of Fire	208064000	789505652

Selection

Do the following:

Select the column *Movie* from the DataFrame.

1: Project

Pr... + × | in_class_coding.py x

- ITP449_Fall20
 - Class
 - .idea
 - Files
 - In Class
 - in_class_coding.py
 - venv lib
- External Libr
- Scratches and

2: Structure

```
1 import pandas as pd
2
3 data = {'ReleaseYear': [2009, 2019, 2015, 2013, 2005],
4         'Movie': ['Avatar', 'Avengers: Endgame', 'Start Wars Ep VII: The Force Awakens',
5                   'Frozen', 'HP and the Goblet of Fire'],
6         'Expense': [328262000, 475560000, 381704000, 245904000, 208064000],
7         'Profit': [1499182307, 1010189087, 843123707, 797747501, 789505652]}
8
9 df = pd.DataFrame(data)
10 pd.set_option('display.max_columns', None)
11 # Single column selection. All of the following are equivalent
12 print(df['Movie'])
13 print(df.Movie)
14 print(df.loc[:, 'Movie'])
15 print(df.iloc[:, 1])
```

2: Favorites

Run: in_class_coding x

▶	0 Avatar
↑	1 Avengers: Endgame
↓	2 Start Wars Ep VII: The Force Awakens
⤵	3 Frozen
⤶	4 HP and the Goblet of Fire
»	Name: Movie, dtype: object

3: Favorites

Do the following:

Select the columns *Expense* and *Profit* from the DataFrame.

1: Project

	Pr...	+	-	⚙
✓	ITP449_Fall2020	C		

- Class
- .idea
- Files
- In Class Coding
 - in_class_coding.py
- venv library

2: Structure

>	External Libraries
⌚	Scratches and Cons

3: Favorites

▶	↑	Expense	Profit
■	↓	0 328262000	1499182307
—	⟳	1 475560000	1010189087
—	⟲	2 381704000	843123707
»	»	3 245904000	797747501
»	»	4 208064000	789505652

Filtering

Do the following:

Create a condition that determines whether *ReleaseYear* is greater than or equal to 2010.

1: Project

Pr... + | - in_class_coding.py x

- ITP449_Fall2020 C
 - Class
 - .idea
 - Files
 - In Class Coding
 - in_class_coding.py
 - venv library
- External Libraries
- Scratches and Cons

2: Structure

3: Favorites

4: Run

Run

Run

```
1 import pandas as pd
2
3 data = {'ReleaseYear': [2009, 2019, 2015, 2013, 2005],
4         'Movie': ['Avatar', 'Avengers: Endgame', 'Star Wars Ep VII: The Force Awakens',
5                   'Frozen', 'HP and the Goblet of Fire'],
6         'Expense': [328262000, 475560000, 381704000, 245904000, 208064000],
7         'Profit': [1499182307, 1010189087, 843123707, 797747501, 789505652]}
8
9 df = pd.DataFrame(data)
10 pd.set_option('display.max_columns', None)
11 # Row filtering. ReleaseYear >= 2010
12 print(df['ReleaseYear'] >= 2010)
13 print(df.ReleaseYear >= 2010)
14
```

5: TODO

TODO

TODO

6: Problems

Problems

Problems

7: Terminal

Terminal

Terminal

8: Python Console

Python Console

Python Console

Event Log

Do the following:

Display the DataFrame where *ReleaseYear* is greater than or equal to 2010.

in_class_coding.py

```
1 import pandas as pd
2
3 data = {'ReleaseYear': [2009, 2019, 2015, 2013, 2005],
4         'Movie': ['Avatar', 'Avengers: Endgame', 'Star Wars Ep VII: The Force Awakens',
5                   'Frozen', 'HP and the Goblet of Fire'],
6         'Expense': [328262000, 475560000, 381704000, 245904000, 208064000],
7         'Profit': [1499182307, 1010189087, 843123707, 797747501, 789505652]}
8
9 df = pd.DataFrame(data)
10 pd.set_option('display.max_columns', None)
11 # DataFrame selection with row filtering. ReleaseYear >= 2010
12 print(df[df['ReleaseYear'] >= 2010])
13 print(df[df.ReleaseYear >= 2010])
14 print(df.loc[df['ReleaseYear'] >= 2010, :])
15 print(df.loc[df.ReleaseYear >= 2010, :])
```

A 1 X 137

Run: in_class_coding

	ReleaseYear	Movie	Expense	Profit
1	2019	Avengers: Endgame	475560000	1010189087
2	2015	Star Wars Ep VII: The Force Awakens	381704000	843123707
3	2013	Frozen	245904000	797747501

Run

TODO

Problems

Terminal

Python Console

Event Log

Do the following:

Display the *Profit* column where *ReleaseYear* is greater than or equal to 2010.

1: Project

Pr... + | - in_class_coding.py x

- ITP449_Fall2020 C
 - Class
 - .idea
 - Files
 - In Class Coding
 - in_class_coding.py
 - venv library
- External Libraries
- Scratches and Cons

2: Structure

Do the following:

Display columns *Profit* and *Expense* where *ReleaseYear* is greater than or equal to 2010.

1: Project

	Pr...	+	÷	⚙
✓	ITP449_Fall2020	C		
	Class			
>	.idea			
>	Files			
✓	In Class Coding			
>	in_class_coding.py			
>	venv library			
>	External Libraries			
⌚	Scratches and Cons			

2: Structure

	Run:	in_class_coding	×	⚙	-
▶		Profit	Expense		
↑		1	1010189087	475560000	
↓		2	843123707	381704000	
⤵		3	797747501	245904000	
»		Process finished with exit code 0			
»					

1: Project

	Pr...	+	÷	⚙
✓	ITP449_Fall2020	C		

- Class
- .idea
- Files
- In Class Coding
 - in_class_coding.py
- venv library

2: Structure

	External Libraries
>	Scratches and Cons

Run: in_class_coding

	Profit	Expense
1	1010189087	475560000
2	843123707	381704000
3	797747501	245904000

> Run

TODO

6: Problems

Terminal

Python Console

Event Log

Arithmetic

Do the following:

Calculate the *Profit* in millions of dollars: divide the *Profit* column by 1,000,000

1: Project

Pr... + | - in_class_coding.py x

- ITP449_Fall2020 C
 - Class
 - .idea
 - Files
 - In Class Coding
 - in_class_coding.py
 - venv library
- External Libraries
- Scratches and Cons

2: Structure

3: Favorites

```
import pandas as pd
data = {'ReleaseYear': [2009, 2019, 2015, 2013, 2005],
        'Movie': ['Avatar', 'Avengers: Endgame', 'Star Wars Ep VII: The Force Awakens',
                  'Frozen', 'HP and the Goblet of Fire'],
        'Expense': [328262000, 475560000, 381704000, 245904000, 208064000],
        'Profit': [1499182307, 1010189087, 843123707, 797747501, 789505652]}

df = pd.DataFrame(data)
pd.set_option('display.max_columns', None)
# Calculate Profit in millions. Divide by 1,000,000
print(df.Profit / 1000000)
print(df['Profit'] / 1000000)
```

Run:

Run: in_class_coding x

Run

Up

Down

Left

Right

Down

Up

>>

>>>

Do the following:

Calculate the age of each movie from the *ReleaseYear*

1: Project 1: Structure

1: in_class_coding.py

```
1 import pandas as pd
2
3 data = {'ReleaseYear': [2009, 2019, 2015, 2013, 2005],
4         'Movie': ['Avatar', 'Avengers: Endgame', 'Star Wars Ep VII: The Force Awakens',
5                   'Frozen', 'HP and the Goblet of Fire'],
6         'Expense': [328262000, 475560000, 381704000, 245904000, 208064000],
7         'Profit': [1499182307, 1010189087, 843123707, 797747501, 789505652]}
8
9 df = pd.DataFrame(data)
10 pd.set_option('display.max_columns', None)
11 # Calculate the age of each movie: 2020 - ReleaseYear
12 print(2020 - df.ReleaseYear)
13 print(2020 - df['ReleaseYear'])
```

2: Favorites

Run: in_class_coding

▶	0	11
↑	1	1
↓	2	5
⤵	3	7
⤶	4	15
»»	Name: ReleaseYear, dtype: int64	

2: Run TODO 6: Problems Terminal Python Console Event Log

17:1 CRLF UTF-8 4 spaces Python 3.8 (Class)

Do the following:

Calculate *Expense* as a percentage of *Profit*

1: Project

Pr... + | - in_class_coding.py x

- ITP449_Fall2020 C
 - Class
 - .idea
 - Files
 - In Class Coding
 - in_class_coding.py
 - venv library
- External Libraries
- Scratches and Cons

```
1 import pandas as pd
2
3 data = {'ReleaseYear': [2009, 2019, 2015, 2013, 2005],
4         'Movie': ['Avatar', 'Avengers: Endgame', 'Star Wars Ep VII: The Force Awakens',
5                   'Frozen', 'HP and the Goblet of Fire'],
6         'Expense': [328262000, 475560000, 381704000, 245904000, 208064000],
7         'Profit': [1499182307, 1010189087, 843123707, 797747501, 789505652]}
8
9 df = pd.DataFrame(data)
10 pd.set_option('display.max_columns', None)
11 # Calculate Expense as percentage of Profit
12 print(df.Expense / df.Profit * 100)
13 print(df['Expense'] / df['Profit'] * 100)
14
```

A 1 137

2: Structure

Run: in_class_coding x



	0	21.896070
	1	47.076335
	2	45.272597
	3	30.824791
	4	26.353706
		dtype: float64

2: Favorites

3: Favorites

▶ 4: Run TODO 6: Problems Terminal Python Console Event Log

Aggregation

Do the following:

Calculate the total *Expense* for all the movies in the DataFrame

ITP449_Fall2020 > Class > In Class Coding > in_class_coding.py

 in_class_coding ▾



1: Project

Z: Structure

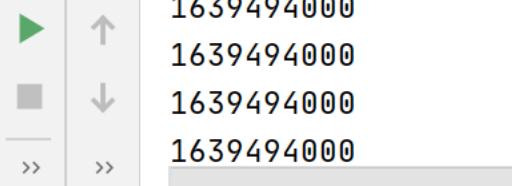
2: Favorites

```
import pandas as pd
import numpy as np

data = {'ReleaseYear': [2009, 2019, 2015, 2013, 2005],
        'Movie': ['Avatar', 'Avengers: Endgame', 'Star Wars Ep VII: The Force Awakens',
                  'Frozen', 'HP and the Goblet of Fire'],
        'Expense': [328262000, 475560000, 381704000, 245904000, 208064000],
        'Profit': [1499182307, 1010189087, 843123707, 797747501, 789505652]}

df = pd.DataFrame(data)
pd.set_option('display.max_columns', None)
# aggregation of DataFrames
print(np.sum(df.Expense))
print(np.sum(df['Expense']))
print(df.Expense.sum())
print(df['Expense'].sum())
```

Run: in_class_coding >



► 4: Run

TODO

6: Problem

>_ Terminal

Python Console

Event Log

Do the following:

Calculate the mean *Expense* and *Profit* for all the movies in the DataFrame

1: Project

Pr... + | - in_class_coding.py x
ITP449_Fall2020 C
Class
.idea
Files
In Class Coding
in_class_coding.py
venv library
External Libraries
Scratches and Cons

2: Structure

3: Favorites

Modifying DataFrames

Do the following:

Add a column of 1's to the DataFrame

1: Project ITP449_Fall2020 > Class > In Class Coding > in_class_coding.py

1: Structure Run: in_class_coding

1: Favorites

```
9 df = pd.DataFrame(data)
10 pd.set_option('display.max_columns', None)
11 # modifying DataFrames
12 df['Ones'] = 1
13 print(df)
```

	ReleaseYear	Movie	Expense	Profit	\
0	2009	Avatar	328262000	1499182307	
1	2019	Avengers: Endgame	475560000	1010189087	
2	2015	Star Wars Ep VII: The Force Awakens	381704000	843123707	
3	2013	Frozen	245904000	797747501	
4	2005	HP and the Goblet of Fire	208064000	789505652	

	Ones
0	1
1	1
2	1
3	1
4	1

4: Run TODO 6: Problems Terminal Python Console Event Log

11:15 CRLF UTF-8 4 spaces Python 3.8 (Class)

Do the following:

Add a column titled *IsOld* of boolean type that is True for any movie with a release date older than 2015

1: Project ITP449_Fall2020 > Class > In Class Coding > in_class_coding.py

1: Structure

1: Favorites

in_class_coding.py

```
9 df = pd.DataFrame(data)
10 pd.set_option('display.max_columns', None)
11 # modifying DataFrames
12 df['IsOld'] = df.ReleaseYear < 2015
13 print(df)
```

Run: in_class_coding

	ReleaseYear	Movie	Expense	Profit	\
0	2009	Avatar	328262000	1499182307	
1	2019	Avengers: Endgame	475560000	1010189087	
2	2015	Star Wars Ep VII: The Force Awakens	381704000	843123707	
3	2013	Frozen	245904000	797747501	
4	2005	HP and the Goblet of Fire	208064000	789505652	

	IsOld
0	True
1	False
2	False
3	True
4	True

4: Run TODO 6: Problems Terminal Python Console Event Log

11:23 CRLF UTF-8 4 spaces Python 3.8 (Class)

Do the following:

Add a column titled *ExpensePercent* of float type
that is equal to the *Expense* as a percentage of
Profit

1: Project ITP449_Fall2020 > Class > In Class Coding > in_class_coding.py

1: Structure Run: in_class_coding

```
df = pd.DataFrame(data)
pd.set_option('display.max_columns', None)
# modifying DataFrames
df['ExpensePercent'] = df.Expense / df.Profit * 100
print(df)
```

Z: Favorites

	ReleaseYear	Movie	Expense	Profit	\
0	2009	Avatar	328262000	1499182307	
1	2019	Avengers: Endgame	475560000	1010189087	
2	2015	Star Wars Ep VII: The Force Awakens	381704000	843123707	
3	2013	Frozen	245904000	797747501	
4	2005	HP and the Goblet of Fire	208064000	789505652	

	ExpensePercent
0	21.896070
1	47.076335
2	45.272597
3	30.824791
4	26.353706

4: Run TODO 6: Problems Terminal Python Console Event Log

17:1 CRLF UTF-8 4 spaces Python 3.8 (Class)

Do the following:

Delete the *ExpensePercent* column from the DataFrame

1: Project

	Pr...	+	÷	⚙	-
✓	ITP449_Fall2020	C			

- Class
- .idea
- Files
- In Class Coding
 - in_class_coding.py
- venv library

External Libraries

Scratches and Cons

2: Structure

	Pr...	+	÷	⚙	-
✓	ITP449_Fall2020	C			

- Class
- .idea
- Files
- In Class Coding
 - in_class_coding.py
- venv library

External Libraries

Scratches and Cons

2: Favorites

	Run:	in_class_coding	×	⚙	-
▶	ReleaseYear		Movie	Expense	Profit

- 0 2009
- 1 2019
- 2 2015
- 3 2013
- 4 2005

	Run	TODO	Problems	Terminal	Python Console	Event Log
▶	4: Run	5: TODO	6: Problems	7: Terminal	8: Python Console	9: Event Log

22:1 CRLF UTF-8 4 spaces Python 3.8 (Class)

Combine methods for analysis

Do the following:

Calculate the total *Expense* and *Profit* for movies released before 2015

1: Project

Pr... + | - in_class_coding.py x

- ITP449_Fall2020 C
 - Class
 - .idea
 - Files
 - In Class Coding
 - in_class_coding.py
 - venv library
- External Libraries
- Scratches and Cons

2: Structure

3: Favorites

4: Run

Run: in_class_coding x

```
C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Class/In Class Cod  
782230000  
3086435460
```

>>> Process finished with exit code 0

Run TODO Problems Terminal Python Console Event Log

1: Project

Pr... + | - in_class_coding.py x

- ITP449_Fall2020 C
 - Class
 - .idea
 - Files
 - In Class Coding
 - in_class_coding.py
 - venv library
- External Libraries
- Scratches and Cons

2: Structure

```
import pandas as pd
import numpy as np

data = {'ReleaseYear': [2009, 2019, 2015, 2013, 2005],
        'Movie': ['Avatar', 'Avengers: Endgame', 'Star Wars Ep VII: The Force Awakens',
                  'Frozen', 'HP and the Goblet of Fire'],
        'Expense': [328262000, 475560000, 381704000, 245904000, 208064000],
        'Profit': [1499182307, 1010189087, 843123707, 797747501, 789505652]}

df = pd.DataFrame(data)
pd.set_option('display.max_columns', None)
print(np.sum(df.loc[df['ReleaseYear'] < 2015, ['Expense', 'Profit']]))

dfOld = df[df.ReleaseYear < 2015]
print(dfOld.Expense.sum())
print(np.sum(dfOld['Profit']))
```

A 1 137

2: Favorites

Run: in_class_coding x

	Expense	782230000
	Profit	3086435460
	dtype: int64	
	782230000	
	3086435460	



3: Favorites

▶ 4: Run TODO 6: Problems Terminal Python Console

Event Log

Do the following:

Display the *movies* where the *Expense* is less than 30% of *Profit*

1: Project

Pr... + | - in_class_coding.py x

- ITP449_Fall2020 C
 - Class
 - .idea
 - Files
 - In Class Coding
 - in_class_coding.py
 - venv library
- External Libraries
- Scratches and Cons

2: Structure

3: Favorites

Do the following:

Display the *Movie* and *ReleaseYear* columns
sorted by the *ReleaseYear* from newest to oldest

1: Project

2: Structure

2: Favorites

Pr... - in_class_coding.py x

```
1 import pandas as pd
2
3 data = {'ReleaseYear': [2009, 2019, 2015, 2013, 2005],
4         'Movie': ['Avatar', 'Avengers: Endgame', 'Star Wars Ep VII: The Force Awakens',
5                   'Frozen', 'HP and the Goblet of Fire'],
6         'Expense': [328262000, 475560000, 381704000, 245904000, 208064000],
7         'Profit': [1499182307, 1010189087, 843123707, 797747501, 789505652]}
8
9 df = pd.DataFrame(data)
10 pd.set_option('display.max_columns', None)
11
12 print(df[['Movie', 'ReleaseYear']].sort_values(by='ReleaseYear', ascending=False))
13
```

Run: in_class_coding x



	Movie	ReleaseYear
1	Avengers: Endgame	2019
2	Star Wars Ep VII: The Force Awakens	2015
3	Frozen	2013
0	Avatar	2009
4	HP and the Goblet of Fire	2005

Run

TODO

Problems

Terminal

Python Console

Event Log

NumPy and Pandas

In Action

Description

This dataset contains a list of video games with sales greater than 100,000 copies. It was generated by a scrape of vgchartz.com.

Fields include

- Rank - Ranking of overall sales
- Name - The game's name
- Platform - Platform of the game's release (i.e. PC, PS4, etc.)
- Year - Year of the game's release
- Genre - Genre of the game
- Publisher - Publisher of 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 - Total worldwide sales.

The script to scrape the data is available at <https://github.com/GregorUT/vgchartzScrape>. It is based on BeautifulSoup using Python.

There are 16,598 records. 2 records were dropped due to incomplete information.

1: Project

Pr...	+	-
ITP449_Fall2020	C:\	
Class		
.idea		
Files		
In Class Coding		
in_class_coding.py		
venv	library r	
External Libraries		

2: Structure

Run: in_class_coding



```
Rank           Name Platform  Year      Genre Publisher \
0   1    Wii Sports     Wii 2006.0  Sports  Nintendo
1   2 Super Mario Bros.    NES 1985.0 Platform  Nintendo
2   3      Mario Kart Wii     Wii 2008.0  Racing  Nintendo
3   4    Wii Sports Resort     Wii 2009.0  Sports  Nintendo
4   5 Pokemon Red/Pokemon Blue     GB 1996.0 Role-Playing  Nintendo
```

```
NA_Sales EU_Sales JP_Sales Other_Sales Global_Sales
0   41.49    29.02    3.77      8.46    82.74
1   29.08    3.58    6.81      0.77    40.24
2   15.85   12.88    3.79      3.31    35.82
3   15.75   11.01    3.28      2.96    33.00
4   11.27    8.89   10.22      1.00    31.37
```

2: Favorites



▶ 4: Run

TODO

6: Problems

Terminal

Python Console

Event Log

1: Project

2: Structure

2: Favorites

 Pr... in_class_coding.py x
ITP449_Fall2020 C:\
 Class
 .idea
 Files
 In Class Coding
 in_class_coding.py
 venv library r
External Libraries

Run: in_class_coding x

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 16598 entries, 0 to 16597
Data columns (total 11 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Rank        16598 non-null   int64  
 1   Name         16598 non-null   object  
 2   Platform     16598 non-null   object  
 3   Year         16327 non-null   float64 
 4   Genre        16598 non-null   object  
 5   Publisher    16540 non-null   object  
 6   NA_Sales     16598 non-null   float64 
 7   EU_Sales     16598 non-null   float64 
 8   JP_Sales     16598 non-null   float64 
 9   Other_Sales  16598 non-null   float64 
 10  Global_Sales 16598 non-null   float64
```

Run

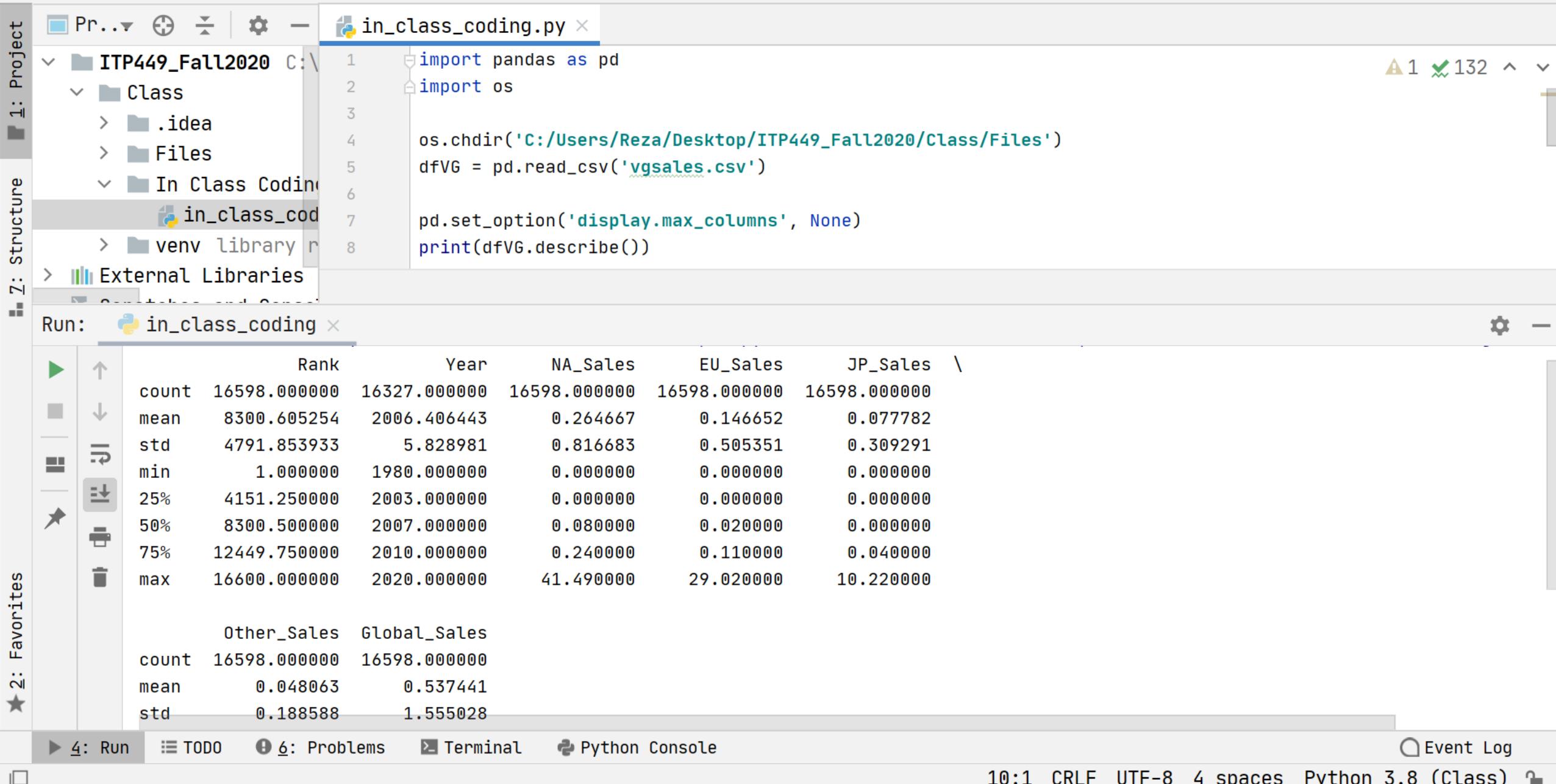
TODO

6: Problems

Terminal

Python Console

Event Log



Do the following:

Print the number of rows and number of columns in the dataset.

1: Project

2: Structure

2: Favorites

Pr... - in_class_coding.py ×

ITP449_Fall2020 C:\

- Class
- .idea
- Files
- In Class Coding
 - in_class_coding.py
- venv library r

External Libraries

Scratches and Console

Run: in_class_coding ×

C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Class/In Class Coding/in_class_coding.py"

Number of rows and columns: (16598, 11)

Number of rows: 16598

Number of columns: 11

Process finished with exit code 0



▶ 4: Run

TODO

6: Problems

Terminal

Python Console

Event Log

Do the following:

Print the average sales for Global, North America, European Union and Japan for games with global sales of at least 1 million.

1: Project

Pr... + - | in_class_coding.py x
ITP449_Fall2020 C:\
 Class
 .idea
 Files
 In Class Coding
 in_class_coding.py
 venv library r
External Libraries
Scratches and Conso

2: Structure

3: Favorites

4: Run

5: TODO

6: Problems

7: Terminal

8: Python Console

9: Event Log

10: File

11: Edit

12: View

13: Navigate

14: Code

15: Refactor

16: Run

17: Tools

18: VCS

19: Window

20: Help

21: File

22: Edit

23: View

24: Navigate

25: Code

26: Refactor

27: Run

28: Tools

29: VCS

30: Window

31: Help

32: File

33: Edit

34: View

35: Navigate

36: Code

37: Refactor

38: Run

39: Tools

40: VCS

41: Window

42: Help

43: File

44: Edit

45: View

46: Navigate

47: Code

48: Refactor

49: Run

50: Tools

51: VCS

52: Window

53: Help

54: File

55: Edit

56: View

57: Navigate

58: Code

59: Refactor

60: Run

61: Tools

62: VCS

63: Window

64: Help

65: File

66: Edit

67: View

68: Navigate

69: Code

70: Refactor

71: Run

72: Tools

73: VCS

74: Window

75: Help

76: File

77: Edit

78: View

79: Navigate

80: Code

81: Refactor

82: Run

83: Tools

84: VCS

85: Window

86: Help

87: File

88: Edit

89: View

90: Navigate

91: Code

92: Refactor

93: Run

94: Tools

95: VCS

96: Window

97: Help

98: File

99: Edit

100: View

101: Navigate

102: Code

103: Refactor

104: Run

105: Tools

106: VCS

107: Window

108: Help

109: File

110: Edit

111: View

112: Navigate

113: Code

114: Refactor

115: Run

116: Tools

117: VCS

118: Window

119: Help

120: File

121: Edit

122: View

123: Navigate

124: Code

125: Refactor

126: Run

127: Tools

128: VCS

129: Window

130: Help

131: File

132: Edit

133: View

134: Navigate

135: Code

136: Refactor

137: Run

138: Tools

139: VCS

140: Window

141: Help

142: File

143: Edit

144: View

145: Navigate

146: Code

147: Refactor

148: Run

149: Tools

150: VCS

151: Window

152: Help

153: File

154: Edit

155: View

156: Navigate

157: Code

158: Refactor

159: Run

160: Tools

161: VCS

162: Window

163: Help

164: File

165: Edit

166: View

167: Navigate

168: Code

169: Refactor

170: Run

171: Tools

172: VCS

173: Window

174: Help

175: File

176: Edit

177: View

178: Navigate

179: Code

180: Refactor

181: Run

182: Tools

183: VCS

184: Window

185: Help

186: File

187: Edit

188: View

189: Navigate

190: Code

191: Refactor

192: Run

193: Tools

194: VCS

195: Window

196: Help

197: File

198: Edit

199: View

200: Navigate

201: Code

202: Refactor

203: Run

204: Tools

205: VCS

206: Window

207: Help

208: File

209: Edit

210: View

211: Navigate

212: Code

213: Refactor

214: Run

215: Tools

216: VCS

217: Window

218: Help

219: File

220: Edit

221: View

222: Navigate

223: Code

224: Refactor

225: Run

226: Tools

227: VCS

228: Window

229: Help

230: File

231: Edit

232: View

233: Navigate

234: Code

235: Refactor

236: Run

237: Tools

238: VCS

239: Window

240: Help

241: File

242: Edit

243: View

244: Navigate

245: Code

246: Refactor

247: Run

248: Tools

249: VCS

250: Window

251: Help

252: File

253: Edit

254: View

255: Navigate

256: Code

257: Refactor

258: Run

259: Tools

260: VCS

261: Window

262: Help

263: File

264: Edit

265: View

266: Navigate

267: Code

268: Refactor

269: Run

270: Tools

271: VCS

272: Window

273: Help

274: File

275: Edit

276: View

277: Navigate

278: Code

279: Refactor

280: Run

281: Tools

282: VCS

283: Window

284: Help

285: File

286: Edit

287: View

288: Navigate

289: Code

290: Refactor

291: Run

292: Tools

293: VCS

294: Window

295: Help

296: File

297: Edit

298: View

299: Navigate

300: Code

301: Refactor

302: Run

Do the following:

Print all the game genres.

```
['Sports' 'Platform' 'Racing' 'Role-Playing' 'Puzzle' 'Misc' 'Shooter'  
 'Simulation' 'Action' 'Fighting' 'Adventure' 'Strategy']
```

1: Project

ITP449_Fall2020 C:\
 Class
 .idea
 Files
 In Class Coding
 in_class_coding.py
 venv library r
External Libraries
Scratches and Conso

Run: in_class_coding

```
C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Class/In Class Coding/in_class_coding.py"  
['Sports' 'Platform' 'Racing' 'Role-Playing' 'Puzzle' 'Misc' 'Shooter'  
 'Simulation' 'Action' 'Fighting' 'Adventure' 'Strategy']
```

```
Process finished with exit code 0
```

2: Favorites



Run

TODO

Problems

Terminal

Python Console

Event Log

Do the following:

Print the total global sales for games of the Racing and Sports genres.

Global sales of Racing games: 732.04

Global sales of Sports games: 1330.93

1: Project

2: Structure

3: Favorites

Pr... in_class_coding.py

```
1 import pandas as pd
2 import numpy as np
3 import os
4
5 os.chdir('C:/Users/Reza/Desktop/ITP449_Fall2020/Class/Files')
6 dfVG = pd.read_csv('vgsales.csv')
7
8 print('Global sales of "Racing" games:', np.sum(dfVG.Global_Sales[dfVG['Genre'] == 'Racing']))
9 print('Global sales of "Sports" games:', np.sum(dfVG.Global_Sales[dfVG['Genre'] == 'Sports']))
10
```

Run: in_class_coding



C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Cl"

Global sales of "Racing" games: 732.04

Global sales of "Sports" games: 1330.93

Process finished with exit code 0

4: Run

TODO

Problems

Terminal

Python Console

Event Log

Do the following:

Print the top five highest grossing games in North America:

	Name	Publisher	NA_Sales
0	Wii Sports	Nintendo	41.49
1	Super Mario Bros.	Nintendo	29.08
9	Duck Hunt	Nintendo	26.93
5	Tetris	Nintendo	23.20
2	Mario Kart Wii	Nintendo	15.85

1: Project 1: Structure

in_class_coding.py

```
1 import pandas as pd
2 import os
3
4 os.chdir('C:/Users/Reza/Desktop/ITP449_Fall2020/Class/Files')
5 dfVG = pd.read_csv('vgsales.csv')
6
7 print(dfVG[['Name', 'Publisher', 'NA_Sales']].sort_values(by='NA_Sales', ascending=False).head())
8
```

Run: in_class_coding

```
C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Class/in_class_coding.py"
          Name Publisher  NA_Sales
0        Wii Sports   Nintendo    41.49
1  Super Mario Bros.   Nintendo    29.08
9        Duck Hunt   Nintendo    26.93
5         Tetris   Nintendo    23.20
2  Mario Kart Wii   Nintendo    15.85
```

Process finished with exit code 0

Do the following:

Print the top five Shooter games from Nintendo.

		Name	Platform	Global_Sales
9		Duck Hunt	NES	28.31
84		GoldenEye 007	N64	8.09
206	Link's Crossbow Training		Wii	5.00
235		Splatoon	WiiU	4.57
296		Star Fox 64	N64	4.03

1: Project 1: Structure

```
in_class_coding.py ×
1 import pandas as pd
2 import os
3
4 os.chdir('C:/Users/Reza/Desktop/ITP449_Fall2020/Class/Files')
5 dfVG = pd.read_csv('vgsales.csv')
6
7 shooter = dfVG[dfVG['Genre'] == 'Shooter']
8
9 print(shooter.loc[shooter['Publisher'] == 'Nintendo', ['Name', 'Platform', 'Global_Sales']].head())
```

Run: in_class_coding ×

		Name	Platform	Global_Sales
9		Duck Hunt	NES	28.31
84		GoldenEye 007	N64	8.09
206		Link's Crossbow Training	Wii	5.00
235		Splatoon	WiiU	4.57
296		Star Fox 64	N64	4.03

Process finished with exit code 0

Do the following:

Print all the sales totals for each genre.

Genre	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales
Action	877.83	525.00	159.95	187.38	1751.18
Adventure	105.80	64.13	52.07	16.81	239.04
Fighting	223.59	101.32	87.35	36.68	448.91
Misc	410.24	215.98	107.76	75.32	809.96
Platform	447.05	201.63	130.77	51.59	831.37
Puzzle	123.78	50.78	57.31	12.55	244.95
Racing	359.42	238.39	56.69	77.27	732.04
Role-Playing	327.28	188.06	352.31	59.61	927.37
Shooter	582.60	313.27	38.28	102.69	1037.37
Simulation	183.31	113.38	63.70	31.52	392.20
Sports	683.35	376.85	135.37	134.97	1330.93
Strategy	68.70	45.34	49.46	11.36	175.12

1: Project in_class_coding.py ×

1 import pandas as pd
2 import os
3
4 os.chdir('C:/Users/Reza/Desktop/ITP449_Fall2020/Class/Files')
5 dfVG = pd.read_csv('vgsales.csv')
6
7 print(dfVG.groupby(['Genre']).sum().loc[:, 'NA_Sales':'Global_Sales'])

A 1 134

2: Structure 1: Project in_class_coding.py ×

Run: in_class_coding ×

Genre	NA_Sales	EU_Sales	JP_Sales	Other_Sales	Global_Sales
Action	877.83	525.00	159.95	187.38	1751.18
Adventure	105.80	64.13	52.07	16.81	239.04
Fighting	223.59	101.32	87.35	36.68	448.91
Misc	410.24	215.98	107.76	75.32	809.96
Platform	447.05	201.63	130.77	51.59	831.37
Puzzle	123.78	50.78	57.31	12.55	244.95
Racing	359.42	238.39	56.69	77.27	732.04
Role-Playing	327.28	188.06	352.31	59.61	927.37
Shooter	582.60	313.27	38.28	102.69	1037.37
Simulation	183.31	113.38	63.70	31.52	392.20
Sports	683.35	376.85	135.37	134.97	1330.93
Strategy	68.70	45.34	49.46	11.36	175.12