

## DV LAB : Pandas Indexing and Selection

Dinesh Kumar K

225229108

### Simple Series and DataFrames

#### Import necessary modules

```
In [4]: import pandas as pd
```

#### Create a Series to store Temperature values for 1 week

```
In [5]: temp_trichy = pd.Series([40.2,39.8,36.3,39.1,41.3,32.9,36.6])
```

#### show temperature values

```
In [6]: temp_trichy
```

```
Out[6]: 0    40.2  
        1    39.8  
        2    36.3  
        3    39.1  
        4    41.3  
        5    32.9  
        6    36.6  
        dtype: float64
```

#### What is the weather on 2nd day?

```
In [7]: temp_trichy[1]
```

```
Out[7]: 39.8
```

#### Find all days and temperatures where temperature over 40.0 degree Celsius

```
In [8]: temp_trichy[[0,4]]
```

```
Out[8]: 0    40.2  
        4    41.3  
        dtype: float64
```

#### Find only day, not temperature where temperature over 40.0 degree Celsius

```
In [13]: temp_trichy[temp_trichy>=40].index
```

```
Out[13]: Int64Index([0, 4], dtype='int64')
```

### Create a Dataframe for student details from List

```
In [16]: students=[['DS01','Rex','1msc'],['Ds02','peter','2msc'],['cs01','ann','3bsc']]  
df_stud=pd.DataFrame(students,columns=['rollno','name','class'])
```

#### show df\_stud dataframe

```
In [17]: df_stud
```

```
Out[17]:
```

	rollno	name	class
0	DS01	Rex	1msc
1	Ds02	peter	2msc
2	cs01	ann	3bsc

Display all column names of df\_stud

```
In [14]: df_stud.columns
Out[14]: Index(['rollno', 'name', 'class'], dtype='object')
```

Add a new column "address" with values ['Delhi', 'Bangalore', 'Chennai'] to df\_stud

```
In [19]: df_stud['address']=['Delhi','Bangalore','Chennai']
In [20]: df_stud
Out[20]:
```

	rollno	name	class	address
0	DS01	Rex	1msc	Delhi
1	Ds02	peter	2msc	Bangalore
2	cs01	ann	3bsc	Chennai

Create a Dataframe for Phone book from Dictionary

```
In [34]: phonebook={'rex':[9942002764,'rex@abc.com'], 'sam':[9932176542,'sam@xyz.com'], 'peter':[9865323645, 'ann@bhc.edu.com']}
df_phonebook=pd.DataFrame.from_dict(phonebook,orient='index')
In [35]: df_phonebook
Out[35]:
```

	0	1
rex	9942002764	rex@abc.com
sam	9932176542	sam@xyz.com
peter	9865323645	ann@bhc.edu.com

Exploratory Data Analysis on Video Game Review Dataset

Import ign.csv dataset

```
In [36]: reviews=pd.read_csv("ign.csv")
```

Show top-5 rows

```
In [37]: reviews.head()
Out[37]:
```

	Unnamed: 0	score_phrase	title	url	platform	score	genre	editors_choice	release_year	release_month	release_day
0	0	Amazing	LittleBigPlanet PS Vita	/games/littlebigplanet-vita/vita-98907	PlayStation Vita	9.0	Platformer	Y	2012	9	12
1	1	Amazing	LittleBigPlanet PS Vita -- Marvel Super Hero E...	/games/littlebigplanet-ps-vita-marvel-super-he...	PlayStation Vita	9.0	Platformer	Y	2012	9	12
2	2	Great	Splice: Tree of Life	/games/splice/ipad-141070	iPad	8.5	Puzzle	N	2012	9	12
3	3	Great	NHL 13	/games/nhl-13/xbox-360-128182	Xbox 360	8.5	Sports	N	2012	9	11
4	4	Great	NHL 13	/games/nhl-13/ps3-128181	PlayStation 3	8.5	Sports	N	2012	9	11

Show bottom 3 rows

```
In [38]: reviews.tail(3)
```

Out[38]:

	Unnamed: 0	score_phrase	title	url	platform	score	genre	editors_choice	release_year	release_month	release_day
18622	18622	Mediocre	Star Ocean: Integrity and Faithlessness	/games/star-ocean-5/ps4-20035681	PlayStation 4	5.8	RPG	N	2016	6	28
18623	18623	Masterpiece	Inside	/games/inside-playdead/xbox-one-121435	Xbox One	10.0	Adventure	Y	2016	6	28
18624	18624	Masterpiece	Inside	/games/inside-playdead/pc-20055740	PC	10.0	Adventure	Y	2016	6	28

How many rows and columns here?

```
In [39]: reviews.shape
```

Out[39]: (18625, 11)

What are the datatypes?

```
In [40]: reviews.dtypes
```

Out[40]:

Unnamed: 0	int64
score_phrase	object
title	object
url	object
platform	object
score	float64
genre	object
editors_choice	object
release_year	int64
release_month	int64
release_day	int64
dtype:	object

Selecting Columns

```
In [41]: reviews['title'].tail()
```

Out[41]:

18620	Tokyo Mirage Sessions #FE
18621	LEGO Star Wars: The Force Awakens
18622	Star Ocean: Integrity and Faithlessness
18623	Inside
18624	Inside

Name: title, dtype: object

Select multiple columns, title and genre and print head

```
In [42]: reviews[['title', 'genre']].head(10)
```

Out[42]:

	title	genre
0	LittleBigPlanet PS Vita	Platformer
1	LittleBigPlanet PS Vita -- Marvel Super Hero E...	Platformer
2	Splice: Tree of Life	Puzzle
3	NHL 13	Sports
4	NHL 13	Sports
5	Total War Battles: Shogun	Strategy
6	Double Dragon: Neon	Fighting
7	Guild Wars 2	RPG
8	Double Dragon: Neon	Fighting
9	Total War Battles: Shogun	Strategy

Selection Using Position

Select top-5 rows and all columns, same as head() using iloc

```
In [43]: reviews.iloc[0:5]
```

Out[43]:

	Unnamed: 0	score_phrase	title	url	platform	score	genre	editors_choice	release_year	release_month	release_day
0	0	Amazing	LittleBigPlanet PS Vita	/games/littlebigplanet-vita/vita-98907	PlayStation Vita	9.0	Platformer	Y	2012	9	12
1	1	Amazing	LittleBigPlanet PS Vita -- Marvel Super Hero E...	/games/littlebigplanet-ps-vita-marvel-super-he...	PlayStation Vita	9.0	Platformer	Y	2012	9	12
2	2	Great	Splice: Tree of Life	/games/splice/ipad-141070	iPad	8.5	Puzzle	N	2012	9	12
3	3	Great	NHL 13	/games/nhl-13/xbox-360-128182	Xbox 360	8.5	Sports	N	2012	9	11
4	4	Great	NHL 13	/games/nhl-13/ps3-128181	PlayStation 3	8.5	Sports	N	2012	9	11

Select rows from position 5 onwards, and columns from position 5 onwards.

```
In [46]: reviews.iloc[4:,4:].head()
```

Out[46]:

	platform	score	genre	editors_choice	release_year	release_month	release_day
4	PlayStation 3	8.5	Sports	N	2012	9	11
5	Macintosh	7.0	Strategy	N	2012	9	11
6	Xbox 360	3.0	Fighting	N	2012	9	11
7	PC	9.0	RPG	Y	2012	9	11
8	PlayStation 3	3.0	Fighting	N	2012	9	11

Select the first column, and all of the rows for the column

```
In [47]: reviews.iloc[:,0].head()
```

Out[47]:

0	0
1	1
2	2
3	3
4	4

Name: Unnamed: 0, dtype: int64

the 10th row, and all of the columns for that row.

```
In [48]: reviews.iloc[9,:]
```

Out[48]:

Unnamed: 0	9
score_phrase	Good
title	Total War Battles: Shogun
url	/games/total-war-battles-shogun/pc-142564
platform	PC
score	7
genre	Strategy
editors_choice	N
release_year	2012
release_month	9
release_day	11

Name: 9, dtype: object

First column is not useful. So remove it

```
In [49]: reviews=reviews.drop("Unnamed: 0",axis=1)
```

### Selection using Row and Column Labels

```
In [50]: df_stud
```

```
Out[50]:
```

	rollno	name	class	address
0	DS01	Rex	1msc	Delhi
1	Ds02	peter	2msc	Bangalore
2	cs01	ann	3bsc	Chennai

### Print all names using loc

```
In [51]: df_stud.loc[:, 'name']
```

```
Out[51]:
```

```
0    Rex
1  peter
2    ann
Name: name, dtype: object
```

### Let us come back to our reviews. Display the first five rows of reviews using the loc method

```
In [53]: reviews.loc[:5, :]
```

```
Out[53]:
```

	score_phrase	title	url	platform	score	genre	editors_choice	release_year	release_month	release_day
0	Amazing	LittleBigPlanet PS Vita	/games/littlebigplanet-vita/vita-98907	PlayStation Vita	9.0	Platformer	Y	2012	9	12
1	Amazing	LittleBigPlanet PS Vita -- Marvel Super Hero E...	/games/littlebigplanet-ps-vita-marvel-super-he...	PlayStation Vita	9.0	Platformer	Y	2012	9	12
2	Great	Splice: Tree of Life	/games/splice/ipad-141070	iPad	8.5	Puzzle	N	2012	9	12
3	Great	NHL 13	/games/nhl-13/xbox-360-128182	Xbox 360	8.5	Sports	N	2012	9	11
4	Great	NHL 13	/games/nhl-13/ps3-128181	PlayStation 3	8.5	Sports	N	2012	9	11
5	Good	Total War Battles: Shogun	/games/total-war-battles-shogun/mac-142565	Macintosh	7.0	Strategy	N	2012	9	11

### Select score\_phrase column using loc and print head

```
In [54]: reviews.loc[:4, 'score_phrase']
```

```
Out[54]:
```

```
0    Amazing
1    Amazing
2     Great
3     Great
4     Great
Name: score_phrase, dtype: object
```

### Print top 10 values of column label "score\_phrase"

```
In [55]: reviews.loc[:9, 'score_phrase']
```

```
Out[55]:
```

```
0    Amazing
1    Amazing
2     Great
3     Great
4     Great
5     Good
6    Awful
7    Amazing
8    Awful
9     Good
Name: score_phrase, dtype: object
```

### Select from reviews of rows from 5 to 15

```
In [56]: some_reviews=reviews.loc[5:15, :]
```

```
In [57]: some_reviews.head()
```

```
Out[57]:
```

	score_phrase	title	url	platform	score	genre	editors_choice	release_year	release_month	release_day
5	Good	Total War Battles: Shogun	/games/total-war-battles-shogun/mac-142565	Macintosh	7.0	Strategy	N	2012	9	11
6	Awful	Double Dragon: Neon	/games/double-dragon-neon/xbox-360-131320	Xbox 360	3.0	Fighting	N	2012	9	11
7	Amazing	Guild Wars 2	/games/guild-wars-2/pc-896298	PC	9.0	RPG	Y	2012	9	11
8	Awful	Double Dragon: Neon	/games/double-dragon-neon/ps3-131321	PlayStation 3	3.0	Fighting	N	2012	9	11
9	Good	Total War Battles: Shogun	/games/total-war-battles-shogun/pc-142564	PC	7.0	Strategy	N	2012	9	11

### Select score of first 3 rows some\_reviews

```
In [58]: some_reviews.loc[:, 'score'].head(3)
```

```
Out[58]:
```

```
5    7.0
6    3.0
7    9.0
Name: score, dtype: float64
```

### Select "score", "genre", and "release\_year" columns from reviews dataframe and print head

```
In [59]: reviews.loc[:, ['score', 'genre', 'release_year']].head()
```

```
Out[59]:
```

	score	genre	release_year
0	9.0	Platformer	2012
1	9.0	Platformer	2012
2	8.5	Puzzle	2012
3	8.5	Sports	2012
4	8.5	Sports	2012

### What is the datatype of "score" column?

```
In [60]: a=reviews.loc[:, 'score']
type(a)
```

```
Out[60]: pandas.core.series.Series
```

## Aggregate Columns

### Find average value of score column in reviews dataframe

```
In [61]: reviews.score.mean()
```

```
Out[61]: 6.950459060402666
```

### Find average value of all numeric columns

```
In [62]: reviews.mean()
```

```
Out[62]: score          6.950459
release_year    2006.515329
release_month      7.138470
release_day      15.603866
dtype: float64
```

### Find average value for each numeric column

```
In [63]: reviews.mean()
```

```
Out[63]: score          6.950459
release_year    2006.515329
release_month      7.138470
release_day      15.603866
dtype: float64
```

Find average value for each row containing numeric values and print head

```
In [64]: reviews.mean(axis=1).head()

Out[64]: 0    510.500
         1    510.500
         2    510.375
         3    510.125
         4    510.125
dtype: float64
```

Find lowest, highest, median, standard deviation of score column of reviews dataframe

show median of "score" column of reviews dataframe

```
In [65]: reviews.score.median()

Out[65]: 7.3
```

show minimum of "score" column of reviews dataframe

```
In [66]: a=reviews.score
         min(a)

Out[66]: 0.5
```

show maximum of "score" column of reviews dataframe

```
In [67]: max(a)

Out[67]: 10.0
```

show standard deviation of "score" column of reviews dataframe

```
In [68]: reviews['score'].std()

Out[68]: 1.7117358608045874
```

How many non-null values in "score" column of reviews dataframe?

```
In [69]: reviews['score'].notnull().sum()

Out[69]: 18625
```

Show the summary of reviews dataframe

```
In [70]: reviews.describe()

Out[70]:
```

	score	release_year	release_month	release_day
count	18625.000000	18625.000000	18625.000000	18625.000000
mean	6.950459	2006.515329	7.13847	15.603866
std	1.711736	4.587529	3.47671	8.690128
min	0.500000	1970.000000	1.00000	1.000000
25%	6.000000	2003.000000	4.00000	8.000000
50%	7.300000	2007.000000	8.00000	16.000000
75%	8.200000	2010.000000	10.00000	23.000000
max	10.000000	2016.000000	12.00000	31.000000

Check if review score has any correlation with other columns of reviews

```
In [71]: reviews.corr()
```

Out[71]:

	score	release_year	release_month	release_day
score	1.000000	0.062716	0.007632	0.020079
release_year	0.062716	1.000000	-0.115515	0.016867
release_month	0.007632	-0.115515	1.000000	-0.067964
release_day	0.020079	0.016867	-0.067964	1.000000

Math Operations on DF columns

Divide the values of "score" column in reviews dataframe by 2. There will be too many values, so just print head

```
In [72]: (reviews.score/2).head()
```

Out[72]:

0	4.50
1	4.50
2	4.25
3	4.25
4	4.25

Name: score, dtype: float64

Boolean Indexing in Pandas

Select all video games whose review score > 7, call it score\_filter

```
In [73]: score_filter=(reviews.score>7)
```

Print head of score\_filter

```
In [74]: score_filter.head()
```

Out[74]:

0	True
1	True
2	True
3	True
4	True

Name: score, dtype: bool

Select all rows for score\_filter column and print its head

```
In [75]: filtered_reviews=reviews[score_filter]
```

```
In [76]: filtered_reviews.head()
```

Out[76]:

	score_phrase	title	url	platform	score	genre	editors_choice	release_year	release_month	release_day
0	Amazing	LittleBigPlanet PS Vita	/games/littlebigplanet-vita/vita-98907	PlayStation Vita	9.0	Platformer	Y	2012	9	12
1	Amazing	LittleBigPlanet PS Vita -- Marvel Super Hero E...	/games/littlebigplanet-ps-vita-marvel-super-he...	PlayStation Vita	9.0	Platformer	Y	2012	9	12
2	Great	Splice: Tree of Life	/games/splice/ipad-141070	iPad	8.5	Puzzle	N	2012	9	12
3	Great	NHL 13	/games/nhl-13/xbox-360-128182	Xbox 360	8.5	Sports	N	2012	9	11
4	Great	NHL 13	/games/nhl-13/ps3-128181	PlayStation 3	8.5	Sports	N	2012	9	11

Show the size of filtered\_reviews

```
In [77]: filtered_reviews.shape
```

Out[77]:

(9800, 10)
------------



### Show top 10 "title" from filtered\_reviews

```
In [78]: (filtered_reviews.title).head(10)
```

```
Out[78]: 0          LittleBigPlanet PS Vita
1  LittleBigPlanet PS Vita -- Marvel Super Hero E...
2          Splice: Tree of Life
3          NHL 13
4          NHL 13
7          Guild Wars 2
10         Tekken Tag Tournament 2
11         Tekken Tag Tournament 2
13         Mark of the Ninja
14         Mark of the Ninja
Name: title, dtype: object
```

### Find games released for the Xbox One platform that have a score of more than 7

### FIND CREATE A FILTER, CALLED XBOX\_ONE\_FILTER FOR THE CONDITIONS

```
In [79]: xbox_one_filter = (reviews["score"] > 7) & (reviews["platform"] == "Xbox One")
```

### SELECT THOSE ROWS FROM REVIEWS OF XBOX\_ONE\_FILTER AND PRINT HEAD

```
In [80]: filtered_reviews2 = reviews[xbox_one_filter]
         filtered_reviews2.head()
```

```
Out[80]:
```

	score_phrase	title	url	platform	score	genre	editors_choice	release_year	release_month	release_day
17137	Amazing	Gone Home	/games/gone-home/xbox-one-20014361	Xbox One	9.5	Simulation	Y	2013	8	15
17197	Amazing	Rayman Legends	/games/rayman-legends/xbox-one-20008449	Xbox One	9.5	Platformer	Y	2013	8	26
17295	Amazing	LEGO Marvel Super Heroes	/games/lego-marvel-super-heroes/xbox-one-20000826	Xbox One	9.0	Action	Y	2013	10	22
17313	Great	Dead Rising 3	/games/dead-rising-3/xbox-one-124306	Xbox One	8.3	Action	N	2013	11	18
17317	Great	Killer Instinct	/games/killer-instinct-2013/xbox-one-20000538	Xbox One	8.4	Fighting	N	2013	11	18

### WHAT IS THE SIZE OF FILTERED\_REVIEWS 2

```
In [81]: filtered_reviews2.shape
```

```
Out[81]: (140, 10)
```

### SELECT ALL VIDEO GAMES WHICH ARE 'ACTION'

```
In [82]: action_reviews = reviews[reviews.genre == 'Action']
```

```
In [83]: action_reviews.head()
```

```
Out[83]:
```

	score_phrase	title	url	platform	score	genre	editors_choice	release_year	release_month	release_day
17	Great	Avengers Initiative	/games/avengers-initiative/iphone-141579	iPhone	8.0	Action	N	2012	9	5
34	Good	War of the Roses	/games/war-of-the-roses-140577/pc-115849	PC	7.3	Action	N	2012	10	3
45	Amazing	Bad Piggies	/games/bad-piggies/iphone-141455	iPhone	9.2	Action	Y	2012	10	1
49	Okay	Demon's Score	/games/demons-score/iphone-118050	iPhone	6.9	Action	N	2012	9	27
69	Great	Hotline Miami	/games/hotline-miami/pc-139657	PC	8.8	Action	Y	2012	10	26

```
In [84]: action_reviews.shape
```

```
Out[84]: (3797, 10)
```

### PLOT REVIEW RATINGS OF TWO PLAY STATIONS AND COMPARE WHICH ONEHAS MORE RATINGS?

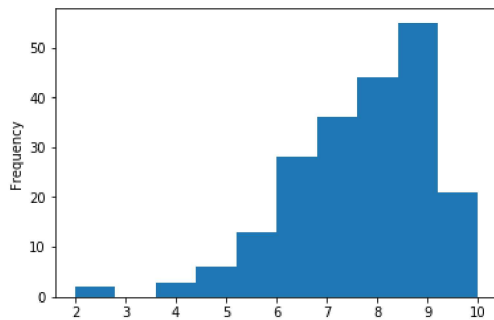
Now that we know how to filter, We can create plots to observe the review distribution for the Xbox one vs the review distributionfor the PlayStation 4 . This will help us figure out which console has better games. We can do this via a histogram, which will plot the frequencies for different score

ranges

**PLOT HISTOGRAM FOR THE FREQUENCIES OF DIFFERENT SCORE RANGES OF XBOX ONE PLATFORM**

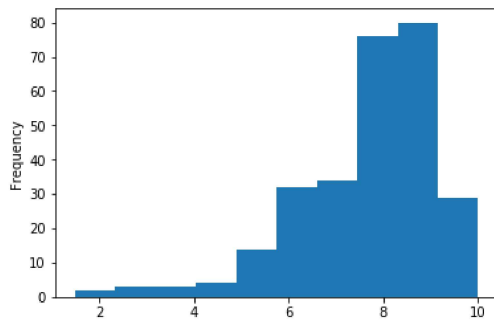
```
In [87]: import matplotlib.pyplot as plt  
reviews[reviews["platform"] == "Xbox One"]["score"].plot(kind="hist")
```

```
Out[87]: <matplotlib.axes._subplots.AxesSubplot at 0x233e69544e0>
```

**PLOT HISTOGRAM FOR FREQUENCIES OF THE SCORE OF PLAY STATION 4 PLATFORM**

```
In [89]: reviews[reviews["platform"] == "PlayStation 4"]["score"].plot(kind="hist")
```

```
Out[89]: <matplotlib.axes._subplots.AxesSubplot at 0x233e698c2b0>
```



```
In [ ]:
```