## ∨ STC Jawwy

```
1 """
2 Here we install libraries that are not installed by default
3 Example:  pyslsb
4 Feel free to add any library you are planning to use.
5 """
6 !pip install pyxlsb
```

> Requirement already satisfied: pyxlsb in /usr/local/lib/python3.10/dist-packages (1.0.10)

```
1 !pip install Openpyxl
```

> Requirement already satisfied: Openpyxl in /usr/local/lib/python3.10/dist-packages (3.1.2)
> Requirement already satisfied: et-xmlfile in /usr/local/lib/python3.10/dist-packages (from Openpyxl) (1.1.0)

```
1 # Import the required libraries
2 """
3 Please feel free to import any required libraries as per your needs
4 """
5 import pandas as pd      # provides high-performance, easy to use structures and data analysis tools
6 import pyxlsb            # Excel extention to read xlsb files (the input file)
7 import numpy as np       # provides fast mathematical computation on arrays and matrices
8 from datetime import date
```

## ∨ Jawwy dataset

The dataset consists of meta details about the movies and tv shows as genre. Also details about Users activities, spent duration and if watching in High definition or standard definition. You have to analyse this dataset to find top insights, findings and to solve the four tasks assigned to you.

```
1 dataframe = pd.read_excel("/content/sample_data/stc TV Data Set_T1.xlsx")
2 # Please make a copy of dataset if you are going to work directly and make changes on the dataset
3 # you can use   df=dataframe.copy()
```

```
1 # check the data shape
2 dataframe.shape
```

> (1048575, 13)

```
1 #df=dataframe.copy()
```

```
1 # display the first 5 rows
2 dataframe.head()
```

| | Column1 | date_ | user_id_maped | program_name | duration_seconds | program_class | seas |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 2017-05-27 | 26138 | 100 treets | 40 | MOVIE | |
| 1 | 3 | 2017-05-21 | 7946 | Moana | 17 | MOVIE | |
| 2 | 4 | 2017-08-10 | 7418 | The Mermaid Princess | 8 | MOVIE | |
| 3 | 5 | 2017-07-26 | 19307 | The Mermaid Princess | 76 | MOVIE | |
| 4 | 7 | 2017- | 15860 | Churchill | 87 | MOVIE | |

```
1 # Data Preprocessing on the input data
2 dataframe = dataframe.drop([1,1])        # dropping the index column
```

```
1 dataframe['program_name'] = dataframe['program_name'].str.strip()  # trim spaces in movies names to avoid misspellings in input data
2 #dataframe['date_'] = pd.to_datetime(dataframe['date_'], unit='D', origin='1899-12-30')  # read date column as date data type
3 dataframe[['duration_seconds', 'season','episode','series_title','hd']] = dataframe[['duration_seconds', 'season','episode','series_
4 dataframe[['user_id_maped', 'program_name','program_class','program_desc','program_genre','original_name']] = dataframe[['user_id_ma
```

```
1 # display the dataset after applying data types
2 dataframe.head()
```

| | Column1 | date_ | user_id_maped | program_name | duration_seconds | program_class | seas |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 2017-05-27 | 26138 | 100 treets | 40 | MOVIE | |
| 2 | 4 | 2017-08-10 | 7418 | The Mermaid Princess | 8 | MOVIE | |
| 3 | 5 | 2017-07-26 | 19307 | The Mermaid Princess | 76 | MOVIE | |
| 4 | 7 | 2017-07-07 | 15860 | Churchill | 87 | MOVIE | |
| 5 | 8 | 2017-08-19 | 20775 | Beavis And Butt-Head Do America | 3 | MOVIE | |

```
1 # describe the numeric values in the dataset
2 dataframe.describe()
```
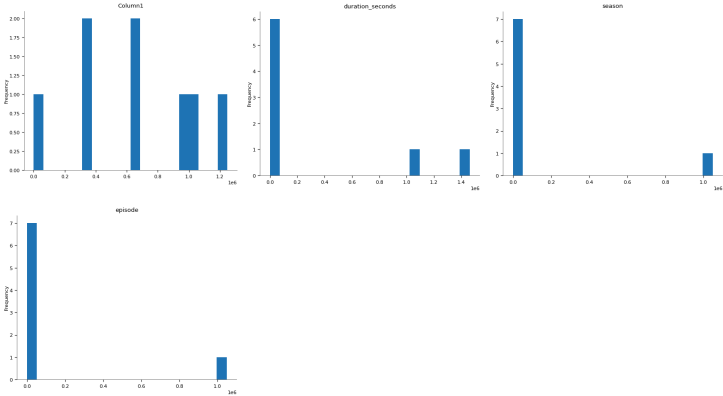
1 to 8 of 8 entries    Filter

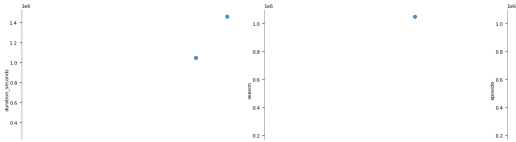| index | Column1 | duration_seconds | season | episode | s |
|---|---|---|---|---|---|
| count | 1048574.0 | 1048574.0 | 1048574.0 | 1048574.0 | |
| mean | 628173.5990669233 | 1230.9583701293375 | 1.342139896659654 | 6.157958331982292 | 0.0120! |
| std | 359703.70509784657 | 6821.061038505865 | 2.1040959772571464 | 12.22015904257431 | 0.1091! |
| min | 1.0 | 2.0 | 0.0 | 0.0 | |
| 25% | 318067.5 | 52.0 | 0.0 | 0.0 | |
| 50% | 630355.5 | 119.0 | 1.0 | 1.0 | |
| 75% | 939822.75 | 1328.0 | 1.0 | 9.0 | |
| max | 1247852.0 | 1461329.0 | 23.0 | 282.0 | |

Show 25 ∨ per page

Like what you see? Visit the [data table notebook](#) to learn more about interactive tables.

**Distributions**

**2-d distributions**

```
1 # check if any column has null value in the dataset
2 dataframe.isnull().any()
```

Column1         False
date_           False

```
user_id_maped        False
program_name         False
duration_seconds     False
program_class        False
season               False
episode              False
program_desc         False
program_genre        False
series_title         False
hd                   False
original_name        False
dtype: bool
```

## ⌄ Task 1

You are required to work on task one to study and HD flag for available dataset

```
1 # make a copy of the dataframe for working on task 1
2 df=dataframe.copy()
```

```
 1 # Here we try to get the most watched movies (Total Views / Total Users Views / Total watch time)
 2 # For series we concatenated the Session episode to differentiate between episodes
 3 grouped=df.copy()
 4 grouped.loc[grouped['program_class'] == 'SERIES/EPISODES', 'program_name'] = grouped['program_name']+'_SE'+grouped['season'].astype(
 5 grouped = grouped.groupby(['program_name','program_class'])\
 6 .agg({'user_id_maped': [('co1', 'nunique'),('co2', 'count')],\
 7      'duration_seconds': [('co3', 'sum')] }).reset_index()
 8 grouped.columns = ['program_name','program_class','No of Users who Watched', 'No of watches', 'Total watch time in seconds']
 9 grouped['Total watch time in houres']=grouped['Total watch time in seconds']/3600
10 grouped = grouped.drop(columns=['Total watch time in seconds'])
11 grouped = grouped.sort_values(by=['Total watch time in houres', 'No of watches','No of Users who Watched'], ascending=False).reset_i
12
```

```
1 # show the result
2 grouped.head(35)
```

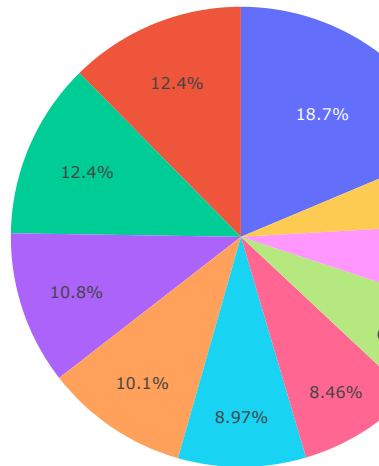| | program_name | program_class | No of Users who Watched | No of watches | Total watch time in houres |
|---|---|---|---|---|---|
| 0 | The Boss Baby | MOVIE | 3389 | 24047 | 2961.350833 |
| 1 | The Amazing pider-Man | MOVIE | 1011 | 2877 | 1966.119167 |
| 2 | The Expendables | MOVIE | 853 | 2119 | 1961.159444 |
| 3 | Moana | MOVIE | 2173 | 8080 | 1706.172222 |
| 4 | Trolls | MOVIE | 2613 | 13793 | 1601.023056 |
| 5 | Bean | MOVIE | 949 | 3617 | 1423.955000 |
| 6 | The murfs | MOVIE | 867 | 3132 | 1342.141111 |
| 7 | Hotel Transylvania | MOVIE | 491 | 1947 | 1096.533611 |
| 8 | Cloudy With a Chance of Meatballs | MOVIE | 683 | 2076 | 948.674722 |
| 9 | The Man With The Iron Fists | MOVIE | 707 | 2505 | 859.626389 |
| 10 | Salt | MOVIE | 563 | 1082 | 767.392778 |
| 11 | Unbroken | MOVIE | 625 | 1429 | 763.078333 |
| 12 | ParaNorman | MOVIE | 614 | 1746 | 747.065556 |
| 13 | Youm Maloosh Lazma | MOVIE | 1131 | 2278 | 718.109722 |
| 14 | Ferdinand | MOVIE | 1278 | 6817 | 714.223056 |
| 15 | White Chicks | MOVIE | 307 | 916 | 711.840833 |
| 16 | Jurassic Park | MOVIE | 504 | 1192 | 693.394444 |
| 17 | The November Man | MOVIE | 494 | 1219 | 679.492222 |
| 18 | Total Recall | MOVIE | 587 | 1108 | 661.820000 |
| 19 | Robin Hood | MOVIE | 588 | 1209 | 643.935000 |
| 20 | Public Enemies | MOVIE | 368 | 716 | 634.035000 |
| 21 | Daddy Day Camp | MOVIE | 263 | 647 | 625.338333 |
| 22 | Oblivion | MOVIE | 790 | 1678 | 609.391111 |
| 23 | Blitz | MOVIE | 562 | 1200 | 570.521944 |
| 24 | War for the Planet of the Apes | MOVIE | 879 | 2028 | 567.597778 |
| 25 | Inside Man | MOVIE | 532 | 1567 | 560.386111 |
| 26 | Bad Boys | MOVIE | 438 | 871 | 559.277500 |
| 27 | Easy A | MOVIE | 513 | 990 | 557.068611 |
| 28 | Battleship | MOVIE | 634 | 1324 | 552.857222 |
| 29 | Baywatch | MOVIE | 2062 | 7436 | 548.995556 |
| 30 | Police tory | MOVIE | 409 | 737 | 520.077222 |

Next steps:   Generate code with `grouped`    ◉ View recommended plots

```
1 # we import Visualization libraries
2 # you can ignore and use any other graphing libraries
3 import matplotlib.pyplot as plt # a comprehensive library for creating static, animated, and interactive visualizations
4 import plotly #a graphing library makes interactive, publication-quality graphs. Examples of how to make line plots, scatter plots,
5 import plotly.express as px
6 import plotly.graph_objects as go
7 from plotly.subplots import make_subplots
```

```
1 # plot top 10 Programs
2 fig = px.pie(grouped.head(10), values='Total watch time in hours', names='program_name',\
3          hover_data=['program_class'],title='top 10 programs in total watch time in houres')
4 fig.show()
```

top 10 programs in total watch time in houres



```
1  # Here we try to study the customer experience against Program class
2  grouped=df.copy()
3  grouped = grouped.groupby('program_class')\
4  .agg({'user_id_maped': [('co1', 'nunique'),('co2', 'count')],\
5        'duration_seconds': [('co3', 'sum')] }).reset_index()
6  grouped.columns = ['program_class','No of Users who Watched', 'No of watches', 'Total watch time in seconds']
7  grouped['Total watch time in houres']=grouped['Total watch time in seconds']/3600
8  grouped = grouped.drop(columns=['Total watch time in seconds'])
9  grouped = grouped.sort_values(by=['Total watch time in houres', 'No of watches','No of Users who Watched'], ascending=False).reset_i
10
```

```
1  # show the result
2  grouped.head()
```
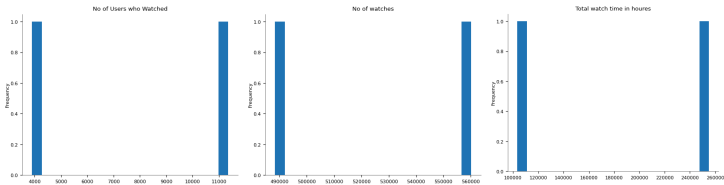
1 to 2 of 2 entries    Filter    □    ?

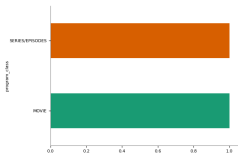| index | program_class | No of Users who Watched | No of watches | Total watch time in houres |
|---|---|---|---|---|
| 0 | SERIES/EPISODES | 3901 | 560174 | 255097.7875 |
| 1 | MOVIE | 11355 | 488400 | 103444.14083333334 |

Show 25 ⌄ per page

📊

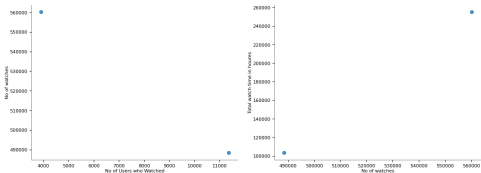Like what you see? Visit the [data table notebook](#) to learn more about interactive tables.
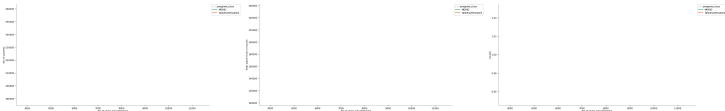
### Distributions



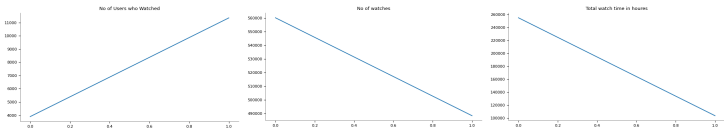### Categorical distributions



### 2-d distributions



### Time series



### Values



### Faceted distributions

<string>:5: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14

             <string>:5: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14

 <string>:5: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14



- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
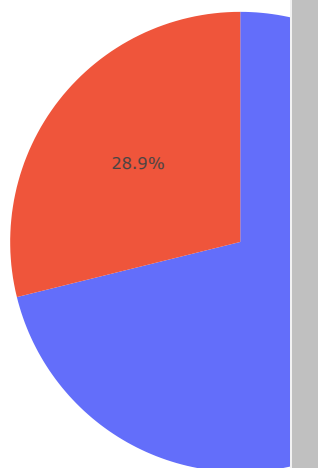
Next steps:    [ Generate code with  grouped ]    [ ◯ View recommended plots ]

```
1  # plot the total watch time against total number of users and report your findings
2  fig = px.pie(grouped, values='Total watch time in houres', names='program_class',\
3                hover_data=['program_class'],title='Total duration spent by program_class')
4  fig2 = px.pie(grouped, values='No of Users who Watched', names='program_class',\
5                hover_data=['program_class'],title='Total Users watching by program_class')
6
7  fig.update_traces(sort=False)
8  fig2.update_traces(sort=False)
9  fig.show()
10 fig2.show()
```

Total duration spent by program_class



Total Users watching by program_class