

In []: `# DATA CLEANING - YOUTUBE CHANNEL DATASET`

In [1]: `import pandas as pd`

In [27]: `ytube = pd.read_csv('D:\\\\24 - Machine_Learning\\download files\\top-5000-youtube-channels.csv')
ytube`

Out[27]:

	Rank	Grade	Channel name	Video Uploads	Subscribers	Video views
	0	1st	A++	Zee TV	82757	18752951
	1	2nd	A++	T-Series	12661	61196302
	2	3rd	A++	Cocomelon - Nursery Rhymes	373	19238251
	3	4th	A++	SET India	27323	31180559
	4	5th	A++	WWE	36756	32852346

	4995	4,996th	B+	Uras Benlioğlu	706	2072942
	4996	4,997th	B+	HI-TECH MUSIC LTD	797	1055091
	4997	4,998th	B+	Mastersaint	110	3265735
	4998	4,999th	B+	Bruce McIntosh	3475	32990
	4999	5,000th	B+	SehatAQUA	254	21172

5000 rows × 6 columns

In [28]: `ytube.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5000 entries, 0 to 4999
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Rank         5000 non-null   object
1   Grade        5000 non-null   object
2   Channel name 5000 non-null   object
3   Video Uploads 5000 non-null   object
4   Subscribers   5000 non-null   object
5   Video Views   5000 non-null   int64
dtypes: int64(1), object(5)
memory usage: 234.5+ KB
```

In [29]: `ytube['Rank'] = ytube['Rank'].str[0:-2].str.replace('.', '').astype('int')
Here we remove the last 2 character from the data of the rank column`

In [30]: `ytube`

Out[30]:

	Rank	Grade	Channel name	Video Uploads	Subscribers	Video views
	0	1	A++	Zee TV	82757	18752951
	1	2	A++	T-Series	12661	61196302
	2	3	A++	Cocomelon - Nursery Rhymes	373	19238251
	3	4	A++	SET India	27323	31180559
	4	5	A++	WWE	36756	32852346

	4995	4996	B+	Uras Benlioğlu	706	2072942
	4996	4997	B+	HI-TECH MUSIC LTD	797	1055091
	4997	4998	B+	Mastersaint	110	3265735
	4998	4999	B+	Bruce McIntosh	3475	32990
	4999	5000	B+	SehatAQUA	254	21172

5000 rows × 6 columns

In [31]: `ytube.dropna()
Here we are removing the Nan values`

Out[31]:

	Rank	Grade	Channel name	Video Uploads	Subscribers	Video views
	0	1	A++	Zee TV	82757	18752951
	1	2	A++	T-Series	12661	61196302
	2	3	A++	Cocomelon - Nursery Rhymes	373	19238251
	3	4	A++	SET India	27323	31180559
	4	5	A++	WWE	36756	32852346

	4995	4996	B+	Uras Benlioğlu	706	2072942
	4996	4997	B+	HI-TECH MUSIC LTD	797	1055091
	4997	4998	B+	Mastersaint	110	3265735
	4998	4999	B+	Bruce McIntosh	3475	32990
	4999	5000	B+	SehatAQUA	254	21172

5000 rows × 6 columns

In [32]: `ytube['Grade'].unique()`

Out[32]: `array(['A++ ', 'A+ ', 'A ', '\xa0 ', 'A- ', 'B+ '], dtype=object)`

In [33]: `mask1 = ytube[ytube['Subscribers'].str.contains('-')].index
ytube.drop(labels = mask1, axis = 0, inplace = True)`

In [34]: `ytube['Subscribers'] = ytube['Subscribers'].astype('int')`

In [35]: `mask1`

Out[35]: `Int64Index([17, 108, 115, 142, 143, 152, 156, 175, 180, 189,
...,
4892, 4893, 4895, 4912, 4936, 4941, 4948, 4956, 4961, 4990],
dtype='int64', length=387)`

In [36]: `ytube.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 4613 entries, 0 to 4999
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Rank         4613 non-null   int32
1   Grade        4613 non-null   object
2   Channel name 4613 non-null   object
3   Video Uploads 4613 non-null   object
4   Subscribers   4613 non-null   int32
5   Video views   4613 non-null   int64
dtypes: int32(2), int64(1), object(3)
memory usage: 216.2+ KB
```

In [37]: `mask2 = ytube[ytube['Video Uploads'].str.contains('-')].index
ytube.drop(labels = mask2, axis = 0, inplace = True)`

In [38]: `ytube['Video Uploads'] = ytube['Video Uploads'].astype('int')`

In [39]: `mask2`

Out[39]: `Int64Index([2323, 3072, 4898], dtype='int64')`

In [40]: `ytube.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 4610 entries, 0 to 4999
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Rank         4610 non-null   int32
1   Grade        4610 non-null   object
2   Channel name 4610 non-null   object
3   Video Uploads 4610 non-null   int32
4   Subscribers   4610 non-null   int32
5   Video views   4610 non-null   int64
dtypes: int32(3), int64(1), object(2)
memory usage: 198.1+ KB
```

In [41]: `ytube['Grade'].unique()`

Out[41]: `array(['A++ ', 'A+ ', 'A ', 'A- ', 'B+ '], dtype=object)`

In [42]: `channel_map = {'A++ ':5, 'A+ ':4, 'A ':3, 'A- ':2, 'B+ ':1}`

In [43]: `ytube['Grade'] = ytube['Grade'].map(channel_map)`

In [44]: `ytube`

Out[44]:

	Rank	Grade	Channel name	Video Uploads	Subscribers	Video views
	0	1	5	Zee TV	82757	18752951
	1	2	5	T-Series	12661	61196302
	2	3	5	Cocomelon - Nursery Rhymes	373	19238251
	3	4	5	SET India	27323	31180559
	4	5	5	WWE	36756	32852346

	4995	4996	1	Uras Benlioğlu	706	2072942
	4996	4997	1	HI-TECH MUSIC LTD	797	1055091
	4997	4998	1	Mastersaint	110	3265735
	4998	4999	1	Bruce McIntosh	3475	32990
	4999	5000	1	SehatAQUA	254	21172

4610 rows × 6 columns

In [45]: `ytube.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 4610 entries, 0 to 4999
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Rank         4610 non-null   int32
1   Grade        4610 non-null   int64
2   Channel name 4610 non-null   object
3   Video Uploads 4610 non-null   int32
4   Subscribers   4610 non-null   int32
5   Video views   4610 non-null   int64
dtypes: int32(3), int64(2), object(1)
memory usage: 198.1+ KB
```