

```
In [1]: 1 import numpy as np #linear algebra
2 import pandas as pd #data processing
3 pd.set_option('display.max_rows', None)
4 pd.set_option('display.max_columns', None)
5 pd.set_option('display.expand_frame_repr', False)
6 pd.set_option('max_colwidth', None)
```

```
In [2]: 1 df_IMDB_Akas = pd.read_csv("./Prj_Data/DownloadedData_Imdb/title.akas (1).ts
2 IM_title_Basics = pd.read_csv("./Prj_Data/DownloadedData_Imdb/title.basics (
3 df_MetaDataFromFinanceTables = pd.read_excel("./Prj_Data/ImdbScrapingData/df
```

```
In [3]: 1 #This table is used as a filter, filter on english and production type
2 df_IMDB_Akas = df_IMDB_Akas.loc[(df_IMDB_Akas['language']=='en')]
3 df_IMDB_Akas = df_IMDB_Akas.loc[(df_IMDB_Akas['types']=='imdbDisplay')]
```

```
In [4]: 1 df_IMDB_Akas = df_IMDB_Akas.drop_duplicates(subset='titleId', keep='first').
2
3 IM_title_Basics = IM_title_Basics.loc[(IM_title_Basics['isAdult']==0)]
4 IM_title_Basics.drop(['isAdult', 'endYear', 'titleType', 'originalTitle'], a
5 IM_title_Basics['runtimeMinutes'] = IM_title_Basics.runtimeMinutes.replace(r
6 IM_title_Basics['runtimeMinutes'] = IM_title_Basics['runtimeMinutes'].astype
7 IM_title_Basics = IM_title_Basics.loc[(IM_title_Basics['runtimeMinutes']>60)
8 IM_title_Basics['startYear'] = IM_title_Basics.startYear.replace(r'\N',0, re
9 IM_title_Basics['startYear'] = IM_title_Basics['startYear'].astype(int)
10 IM_title_Basics = IM_title_Basics.loc[(IM_title_Basics['startYear']>=2005) &
11 IM_title_Basics['startYear_str'] = IM_title_Basics['startYear'].astype(str)
12 IM_title_Basics['primaryTitle'] = IM_title_Basics['primaryTitle'].str.title(
13 IM_title_Basics['titleyear'] = IM_title_Basics['primaryTitle'] + IM_title_Ba
```

```
In [5]: 1 df_IMDB_Akas_english = IM_title_Basics.merge(df_IMDB_Akas, left_on="tconst",
```

```
In [6]: 1 df_IMDB_Akas_english.info()
```

```
In [7]: 1 df_IMDB_Akas_english.drop(["ind_Link"], axis=1, inplace=True)
```

```
In [8]: 1 df_IMDB_Akas_english.head()
```

```
In [9]: 1 df_MetaDataFromFinanceTables.info()
```

```
In [10]: 1 df_MetaDataFromFinanceTables.drop(["year", "RunningTime", "genres", "title", "
```

```
In [11]: 1 df_IMDB_Eng_with_metadata = df_IMDB_Akas_english.merge(df_MetaDataFromFinanc
```

```
In [12]: 1 df_IMDB_Eng_with_metadata.info()
```

```
...
```

```
In [13]: 1 df_IMDB_Eng_with_metadata.drop(["titleId", "Unnamed: 0", "Merg_MetaData" ], ax
```

```
In [14]: 1 df_IMDB_Eng_with_metadata.info()
```

```
...
```

```
In [15]: 1 df_IMDB_Eng_with_metadata.to_excel("df_IMDB_MovieCatalog.xlsx")
```

Data belwo is for quick cutting and copying for testing :)

```
In [ ]: 1 # df_IMDB_Akas_english.drop(["ind_Link", "Unnamed: 0", "titleyear_fin", "title
2 # "title_fin", "Merg_MetaData", "year_fin", "Runnin
3
4 # df_IMDB_Akas_english.drop(["ind_Link", "Unnamed: 0", "titleyear_fin", "title
5 # "title_fin", "Merg_MetaData", "year_fin", "Runnin
```

```
In [ ]: 1 testingLink.loc[testingLink['tconst'] == "tt0499549"][["primaryTitle_x", "tc
2 # .loc[df_IMDB_Akas['tconst'] == "tt0499549"][["primaryTitle", "tconst"]]
3 # df_IMDB_Akas_english[df_IMDB_Akas_english.primaryTitle.str.contains('Happy
```