

This file was used to blend our financial data sources

```
In [54]: 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)
7
8 # Importing into dataframe
9 df_Bom_MovieGross = pd.read_csv("./Prj_Data/DownloadedData_FlatIron/bom.movi
10 df_TN_Movie_Budgets = pd.read_csv("./Prj_Data/DownloadedData_FlatIron/tn.mov
11 df_TN_Movie_Budgets = pd.read_csv("./Prj_Data/DownloadedData_FlatIron/tn.mov
12 df_Scrp_Financials = pd.read_excel("./Prj_Data/ImdbScrapingData/df_Financial
13 df_IMDB_Akas_english = pd.read_excel("./Prj_Data/DownloadedData_Imdb/df_IMDB
14 df_IMDB_InflationAdjuster = pd.read_excel("./Prj_Data/ImdbScrapingData/Infla
15
```

In [55]:

```

1  #PREPARE THE FINACIAL DATA: WE USED 2 SOURCES,1 FROM FLATIRON DOWNLOAD, 2 F
2  #*****
3  #preparing to work with Movie Gross Df, cleaning up data and create right d
4  df_Bom_MovieGross.name = "df_Bom_MovieGross"
5  df_Bom_MovieGross['foreign_gross'] = pd.to_numeric(df_Bom_MovieGross['foreign
6  df_Bom_MovieGross['year_str_BOM'] = df_Bom_MovieGross['year'].astype(str)
7  df_Bom_MovieGross['year'] = df_Bom_MovieGross['year'].astype(int)
8  df_Bom_MovieGross_values = {'domestic_gross':0, 'foreign_gross': 0}
9  df_Bom_MovieGross.fillna(value=df_Bom_MovieGross_values, inplace=True)
10 df_Bom_MovieGross["wwg_calc_BOM"] = 0
11 df_Bom_MovieGross['title'] = df_Bom_MovieGross['title'].str.title() #*****K
12 df_Bom_MovieGross["titleyear"] = df_Bom_MovieGross['title'] + df_Bom_MovieGr
13 df_Bom_MovieGross["wwg_calc_BOM"] = (df_Bom_MovieGross['domestic_gross'] + d
14
15 #Renaming to aid in consolidating between the three sources
16 df_Bom_MovieGross.rename(columns={"foreign_gross": "fg_BOM", "domestic_gross
17                               "title": "title_BOM", "studio": "s
18 #-----
19
20 #preparing to work with Movie budgets, cleaning up data and create right da
21 df_TN_Movie_Budgets.name = "df_TN_Movie_Budgets"
22 df_TN_Movie_Budgets[df_TN_Movie_Budgets.columns[1:]] = df_TN_Movie_Budgets[d
23 df_TN_Movie_Budgets['production_budget'] = pd.to_numeric(df_TN_Movie_Budgets
24 df_TN_Movie_Budgets['domestic_gross'] = pd.to_numeric(df_TN_Movie_Budgets['d
25 df_TN_Movie_Budgets['worldwide_gross'] = pd.to_numeric(df_TN_Movie_Budgets['
26
27 df_TN_Movie_Budgets_values = {'domestic_gross': 0, 'year':0 , 'worldwide_gro
28 df_TN_Movie_Budgets.fillna(value=df_TN_Movie_Budgets_values, inplace=True)
29 df_TN_Movie_Budgets["fg_calc"] = 0
30
31 df_TN_Movie_Budgets['movie'] = df_TN_Movie_Budgets['movie'].str.title() #***
32 df_TN_Movie_Budgets["year"] = df_TN_Movie_Budgets['release_date'].str[-4:].a
33 df_TN_Movie_Budgets["year_str_TN"] = df_TN_Movie_Budgets['release_date'].str
34 df_TN_Movie_Budgets["titleyear"] = df_TN_Movie_Budgets['movie'] + df_TN_Movi
35 df_TN_Movie_Budgets["fg_calc_TN"] = (df_TN_Movie_Budgets['worldwide_gross']
36 df_TN_Movie_Budgets = df_TN_Movie_Budgets.drop('id', axis = 1)
37 df_TN_Movie_Budgets.rename(columns={"domestic_gross": "dg_TN", "worldwide_gr
38                               "production_budget": "pb_TN", "movie": "t
39
40 #needed to do this steep to collapse data given there are dups in a few movi
41 from pandasql import sqldf
42 pysqldf = lambda q: sqldf(q, globals())
43 q3 = """SELECT titleyear as titleyear, max(title_TN) as title_TN, max(year_T
44 min(rd_TN) as rd_TN, Sum(pb_TN) as pb_TN,sum(dg_TN) as dg_TN, sum(wwg_TN) as
45         FROM df_TN_Movie_Budgets
46         GROUP BY titleyear, title_TN
47         """
48 df_TN_Movie_Budgets = pysqldf(q3)
49 #-----
50
51 # Finanally merger data from flatiron
52 df_fI_Financials = df_TN_Movie_Budgets.merge(df_Bom_MovieGross,
53         on='titleyear', how='outer', indicator='Combing_FL_Financials', suffixe
54
55 # Add tconst Key to FFlatIronTables then drop
56 df_FI_financials_With_tconst = df_fI_Financials.merge(df_IMDB_Akas_english[[

```

57

In [56]: 1 `# df_FI_financials_With_tconst.info()`

In [57]: 1 `# df_FI_financials_With_tconst.nunique()`

In [58]: 1 `#USED TO CHECK FOR MERGING ERRORS`
 2 `# # df_FI_financials_With_tconst[["title_TN", "title_BOM", "title_FI",]].sort_`
 3 `# df_FI_financials_With_tconst[["year_TN", "year_BOM", "year_str_FI",]].sort_v`
 4 `# df_FI_financials_With_tconst[["year_TN", "year_BOM", "year_FI",]].sort_value`
 5 `# # df_Bom_MovieGross[["studio", "studio_short", "Studio_Desc"]].sort_values(b`

In [59]: 1 `#Use all data scraped from IMDB as the default for domestic, foreign ww sale`
 2 `df_FI_financials_With_tconst.title_TN.fillna("0", inplace=True)`
 3 `df_FI_financials_With_tconst.year_str_TN.fillna("0", inplace=True)`
 4 `df_FI_financials_With_tconst.year_TN.fillna(0, inplace=True)`
 5
 6 `df_FI_financials_With_tconst.title_BOM.fillna("0", inplace=True)`
 7 `df_FI_financials_With_tconst.year_str_BOM.fillna("0", inplace=True)`
 8 `df_FI_financials_With_tconst.year_BOM.fillna(0, inplace=True)`
 9
 10 `df_FI_financials_With_tconst['title_FI'] = df_FI_financials_With_tconst['tit`
 11 `df_FI_financials_With_tconst['year_str_FI'] = df_FI_financials_With_tconst['`
 12 `df_FI_financials_With_tconst['year_FI'] = df_FI_financials_With_tconst['year`
 13
 14 `df_FI_financials_With_tconst['title_FI'] = df_FI_financials_With_tconst.titl`
 15 `df_FI_financials_With_tconst['year_str_FI'] = df_FI_financials_With_tconst.y`
 16 `df_FI_financials_With_tconst['year_FI'] = df_FI_financials_With_tconst.year_`
 17
 18 `df_FI_financials_With_tconst['title_FI'] = df_FI_financials_With_tconst.appl`
 19 `df_FI_financials_With_tconst['year_str_FI'] = df_FI_financials_With_tconst.a`
 20 `df_FI_financials_With_tconst['year_FI'] = df_FI_financials_With_tconst.apply`
 21

In [60]: 1 df_FI_financials_With_tconst.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2609 entries, 0 to 2608
Data columns (total 22 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   titleyear                             2609 non-null   object
1   title_TN                              2609 non-null   object
2   year_TN                               2609 non-null   float64
3   year_str_TN                           2609 non-null   object
4   rd_TN                                 1983 non-null   object
5   pb_TN                                 1983 non-null   float64
6   dg_TN                                 1983 non-null   float64
7   wwg_TN                                1983 non-null   float64
8   fg_calc_TN                            1983 non-null   float64
9   title_BOM                             2609 non-null   object
10  studio_BOM                            1646 non-null   object
11  dg_BOM                                 1647 non-null   float64
12  fg_BOM                                 1647 non-null   float64
13  year_BOM                              2609 non-null   float64
14  year_str_BOM                           2609 non-null   object
15  wwg_calc_BOM                           1647 non-null   float64
16  Combing_FL_Financials                 2609 non-null   category
17  tconst                                 2609 non-null   object
18  Adding tconst                         2609 non-null   category
19  title_FI                              2609 non-null   object
20  year_str_FI                           2609 non-null   object
21  year_FI                               2609 non-null   float64
dtypes: category(2), float64(10), object(10)
memory usage: 433.3+ KB
```

In [61]: 1 # df_FI_financials_With_tconst.sort_values(by="year_FI").head(100)

In [62]: 1 #DROP UNNEEDED FIELDS
2 df_FI_financials_With_tconst.drop(["year_TN", "year_BOM", "year_TN", "year_BOM"]

In [63]: 1 # df_Scrp_Financials.info()

In []: 1

In [64]: 1 #Merg flatiron financials with scraping financials IMDB Site :https://www.box
2 #*****
3 #preparing to work with Movie Gross Df, cleaning up data and create right d
4 df_MasterFinancials = df_Scrp_Financials.merge(df_FI_financials_With_tconst,
5 on='tconst', how='outer', suf
6 fieldsToConvert = {'dg_IMDB': 0, 'fg_IMDB':0, 'ww_IMDB': 0, 'Domestic Openi
7 'dg_TN': 0, 'wwg_TN':0, 'fg_calc_TN': 0, 'dg_BOM': 0,
8 'wg_calc_BOM': 0}
9 df_MasterFinancials.fillna(value=fieldsToConvert, inplace=True)

In [65]: 1 # df_MasterFinancials.info()

In [66]:

```

1  #Clean up dataframe - consolidate, remove unwanted columns, all remove metad
2  df_MasterFinancials = df_MasterFinancials.reindex(sorted(df_MasterFinancials
3  col_names = ['tconst', 'Adding tconst', 'Combing_FL_Financials',
4              'Genres_IMDB', 'MPAA', 'Running Time_IMDB', 'Unnamed: 0', 'Domestic O
5              'dg_IMDB', 'dg_TN', 'fg_BOM', 'fg_IMDB', 'fg_calc_TN', 'genres',
6              'isAdult', 'mergingfinancials', 'originalTitle', 'pb_IMDB', 'pb_TN',
7              'rd_IMDB', 'rd_TN', 'runtimeMinutes', 'studio_BOM', 'studio_IMDB',
8              'title_BOM', 'title_IMDB', 'title_TN', 'titleyear_IM',
9              'titleyear_fl', 'wwg_calc_BOM', 'ww_IMDB', 'wwg_TN',
10             'year', 'year_BOM', 'year_TN']
11  df_MasterFinancials.reindex(columns=col_names)
12
13
14  # df_MasterFinancials.drop(["Adding tconst", 'Unnamed: 0', 'isAdult', 'mergi
15
16  df_MasterFinancials_Only = df_MasterFinancials[['tconst', 'Domestic Opening',
17                                                  'fg_BOM', 'fg_calc_TN', 'fg_I
18                                                  'title_BOM', 'title_IMDB', '
19                                                  'wwg_calc_BOM', 'wwg_TN', 'ww
20

```

In [67]:

```

1  #Use all data scraped from IMDB as the default for domestic, foreign ww sale
2  df_MasterFinancials_Only['WW_Gross'] = df_MasterFinancials_Only['ww_IMDB'].a
3  df_MasterFinancials_Only['Dom_Gross'] = df_MasterFinancials_Only['dg_IMDB'].
4  df_MasterFinancials_Only['Frqn_Gross'] = df_MasterFinancials_Only['fg_IMDB']
5  df_MasterFinancials_Only['P_Cost'] = df_MasterFinancials_Only['pb_IMDB'].app
6  df_MasterFinancials_Only['year'] = df_MasterFinancials_Only['year'].apply(la
7
8
9
10 #fill in the blanks from max valus from the flatIron finance tables
11 df_MasterFinancials_Only['dg_Fl_Max'] = df_MasterFinancials_Only[["dg_BOM",
12 df_MasterFinancials_Only['fg_Fl_Max'] = df_MasterFinancials_Only[["fg_BOM",
13 df_MasterFinancials_Only['ww_Fl_Max'] = df_MasterFinancials_Only['dg_Fl_Max'
14
15 df_MasterFinancials_Only['WW_Gross'] = df_MasterFinancials_Only.apply(lambda
16 df_MasterFinancials_Only['Dom_Gross'] = df_MasterFinancials_Only.apply(lambd
17 df_MasterFinancials_Only['Frqn_Gross'] = df_MasterFinancials_Only.apply(lamb
18 df_MasterFinancials_Only['P_Cost'] = df_MasterFinancials_Only.apply(lambda x
19 df_MasterFinancials_Only['P_Cost']=df_MasterFinancials_Only['P_Cost'].replac
20
21 df_MasterFinancials_Only['year'] = df_MasterFinancials_Only['year'].fillna(0
22 df_MasterFinancials_Only['year'] = df_MasterFinancials_Only.apply(lambda x:
23
24
25
26 df_MasterFinancials_Only['titleyear'] = df_MasterFinancials_Only["titleyear_
27 df_MasterFinancials_Only.titleyear.fillna(df_MasterFinancials_Only["titleyea
28
29 # df_MasterFinancials_MetaData['titleyear'] = df_MasterFinancials_Only["titl
30 # df_MasterFinancials_MetaData.titleyear.fillna(df_MasterFinancials_Only["ti
31
32

```

<ipython-input-67-e68ce41ead83>:2: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df_MasterFinancials_Only['WW_Gross'] = df_MasterFinancials_Only['ww_IMDB'].
apply(lambda x: x if x>0 else 0)
```

<ipython-input-67-e68ce41ead83>:3: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df_MasterFinancials_Only['Dom_Gross'] = df_MasterFinancials_Only['dg_IMD
```

In [68]: 1 df_MasterFinancials_Only.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 5540 entries, 0 to 5539
Data columns (total 28 columns):
#   Column                Non-Null Count  Dtype
---  -
0   tconst                 5540 non-null   object
1   Domestic Opening      5540 non-null   float64
2   dg_BOM                5540 non-null   float64
3   dg_TN                 5540 non-null   float64
4   dg_IMDB               5540 non-null   float64
5   fg_BOM                5540 non-null   float64
6   fg_calc_TN            5540 non-null   float64
7   fg_IMDB               5540 non-null   float64
8   pb_IMDB               5540 non-null   float64
9   pb_TN                 5540 non-null   float64
10  title_BOM              2609 non-null   object
11  title_IMDB             4920 non-null   object
12  title_TN               2609 non-null   object
13  titleyear_IM           4920 non-null   object
14  titleyear_fl           2609 non-null   object
15  wwg_calc_BOM           5540 non-null   float64
16  wwg_TN                 5540 non-null   float64
17  ww_IMDB               5540 non-null   float64
18  year                   5540 non-null   float64
19  year_FI                2609 non-null   float64
20  WW_Gross               5540 non-null   float64
21  Dom_Gross              5540 non-null   float64
22  Frgn_Gross             5540 non-null   float64
23  P_Cost                 5540 non-null   float64
24  dg_Fl_Max              5540 non-null   float64
25  fg_Fl_Max              5540 non-null   float64
26  ww_Fl_Max              5540 non-null   float64
27  titleyear              5540 non-null   object
dtypes: float64(21), object(7)
memory usage: 1.2+ MB
```

In [69]: 1 # df_MasterFinancials_Only.sort_values(by="year").head(100)

In [70]: 1 df_MasterFinancials_Only = df_MasterFinancials_Only[['tconst', 'P_Cost', 'Dom_

In [71]: 1 # Add Justments for Inflation and Calculate ROI

In []: 1

In [72]: 1 df_MasterFinancials_Only = df_MasterFinancials_Only.merge(df_IMDB_InflationA

In [73]: 1 df_MasterFinancials_Only.head()

Out[73]:

	tconst	P_Cost	Dom_Gross	Frgn_Gross	WW_Gross	Domestic Opening	year	TicketPrice
0	tt0120667	100000000.0	154696080.0	178839854.0	333535934.0	56061504.0	2005.0	6.41
1	tt0121164	300000000.0	53359111.0	64731725.0	118090836.0	388166.0	2005.0	6.41
2	tt0121766	113000000.0	380270577.0	488119983.0	868390560.0	108435841.0	2005.0	6.41
3	tt0200465	20000000.0	30060660.0	34767761.0	64828421.0	5935256.0	2008.0	7.18
4	tt0206634	76000000.0	35552383.0	35043081.0	70595464.0	501003.0	2006.0	6.55

In [74]: 1 df_MasterFinancials_Only.columns

Out[74]: Index(['tconst', 'P_Cost', 'Dom_Gross', 'Frgn_Gross', 'WW_Gross', 'Domestic Opening', 'year', 'TicketPrice', 'EstInflation', 'Multiplier'], dtype='object')

```
In [75]: 1 df_MasterFinancials_Only["adj_P_Cost"] = df_MasterFinancials_Only["P_Cost"]
2 df_MasterFinancials_Only["adj_Frgn_Gross"] = df_MasterFinancials_Only["Frgn_Gross"]
3 df_MasterFinancials_Only["adj_WW_Gross"] = df_MasterFinancials_Only["WW_Gross"]
4 df_MasterFinancials_Only["Profits"] = df_MasterFinancials_Only["WW_Gross"] -
5 df_MasterFinancials_Only["adj_P_Cost"]
6
7 df_MasterFinancials_Only["ROI"] = df_MasterFinancials_Only.apply(lambda row
8
```

In [76]: 1 df_MasterFinancials_Only.head()

Out[76]:

	tconst	P_Cost	Dom_Gross	Frgn_Gross	WW_Gross	Domestic Opening	year	TicketPrice
0	tt0120667	100000000.0	154696080.0	178839854.0	333535934.0	56061504.0	2005.0	6.41
1	tt0121164	300000000.0	53359111.0	64731725.0	118090836.0	388166.0	2005.0	6.41
2	tt0121766	113000000.0	380270577.0	488119983.0	868390560.0	108435841.0	2005.0	6.41
3	tt0200465	20000000.0	30060660.0	34767761.0	64828421.0	5935256.0	2008.0	7.18
4	tt0206634	76000000.0	35552383.0	35043081.0	70595464.0	501003.0	2006.0	6.55

In [77]: 1 df_MasterFinancials_Only.to_excel("df_Movie_Financials.xlsx")

In []: 1

In []:

```
1  #THIS CODE IS FOR CHECKING THE SYNCHING STEP
2  # df_MasterFinancials_MetaData[["title_TN","title_BOM","title_IMDB","title_f
3  # df_MasterFinancials_MetaData[["genres_fin","genres",]].sort_values(by="gen
4  # df_MasterFinancials_MetaData[["year_BOM","year_TN","year"]].sort_values(by
5  # df_MasterFinancials_MetaData[["titleyear_fl","titleyear_IM", "titleyear_fi
6  # df_MasterFinancials_MetaData[["Running Time_IMDB","runtimeMinutes", "Runni
7  #df_MasterFinancials_MetaData[["rd_IMDB","rd_TN", "rd_fin", "rd_fin1", "rd_f
8
```

<

>