This File was used to get our baseline financial data from imdb, we blended this data with the flatiron financial data

```python
In [ ]:  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
```

```python
In [ ]:  1  df_IMDB_Title_Akas = pd.read_csv("./Prj_Data/DownLoadedData_Imdb/imdb_title_
         2  df_IMDB_Title_Akas.name = "df_IMDB_Title_Akas"
         3
         4  df_IMDB_Title_Akas = df_IMDB_Title_Akas.loc[(df_IMDB_Title_Akas['titleType']
         5  df_IMDB_Title_Akas = df_IMDB_Title_Akas.loc[(df_IMDB_Title_Akas['isAdult']==
         6  df_IMDB_Title_Akas = df_IMDB_Title_Akas.loc[~(df_IMDB_Title_Akas['genres']==
         7  df_IMDB_Title_Akas.drop(['endYear', 'isAdult'], axis=1, inplace = True)
         8
         9  df_IMDB_Title_Akas["startYear"] = df_IMDB_Title_Akas.startYear.replace(r'\N'
        10  df_IMDB_Title_Akas['startYear'] = df_IMDB_Title_Akas['startYear'].astype('in
        11
        12
        13  df_IMDB_Title_Akas["runtimeMinutes"] = df_IMDB_Title_Akas.runtimeMinutes.rep
        14  df_IMDB_Title_Akas['runtimeMinutes'] = df_IMDB_Title_Akas['runtimeMinutes'].
        15
        16  # df_IMDB_Title_Akas = df_IMDB_Title_Akas[df_IMDB_Title_Akas.primaryTitle.st
        17  df_IMDB_Title_Akas.dropna(subset=['primaryTitle'], how='all', inplace=True)
        18  df_IMDB_Title_Akas.fillna({"startYear":0,"runtimeMinutes":0}, inplace=True)
        19  df_IMDB_Title_Akas = df_IMDB_Title_Akas.loc[(df_IMDB_Title_Akas['startYear']
        20
        21
        22  df_IMDB_Title_Akas = df_IMDB_Title_Akas[df_IMDB_Title_Akas['primaryTitle'].s
        23  df_IMDB_Title_Akas = df_IMDB_Title_Akas.loc[~(df_IMDB_Title_Akas['primaryTit
        24
        25
        26  df_IMDB_Title_Ratings = pd.read_csv("./Prj_Data/DownLoadedData_Imdb/Ratings.
        27  df_Imdb_MoviesWithRatings = df_IMDB_Title_Akas.merge(df_IMDB_Title_Ratings,
        28
```

```
In [ ]:    1  #_____Indexes To Feed Into Scraping
           2
           3  #df_ttTolookup = pd.read_excel("df_ttTolookup.xlsx")
           4
           5  #df_0_7776 = df_ttTolookup.iloc[0:7776,]
           6        #df_0_7776.index = df_ttTolookup.iloc[0:7776,].index last row to com
           7  #df_11236_14000 = df_ttTolookup.iloc[11236:14000,]
           8  #df_11236_14000["index_O"] = df_ttTolookup.iloc[11236:14000,].index
           9  # df_14000_19999 = df_ttTolookup.iloc[14000:20000,]
          10  # df_14000_19999["index_O"] = df_ttTolookup.iloc[14000:20000,].index
          11
          12  #df_20000_24999v1 = df_ttTolookup_Backup.iloc[20000:25000,]
          13  # df_20000_24999v1["index_O"] = df_ttTolookup_Backup.iloc[20000:25000,].inde
          14
          15
          16  # df_20000_21910 = df_ttTolookup.iloc[20000:21910,]
          17  # df_21910_24999["index_O"] = df_ttTolookup.iloc[20000:21910,].index
          18
          19  # df_21909_24999 = df_ttTolookup.iloc[21910:25000,]
          20  # df_21909_24999["index_O"] = df_ttTolookup.iloc[21910:25000,].index
          21
          22  # df_25000_29999 = df_ttTolookup.iloc[21910:30000,]
          23  # df_25000_29999["index_O"] = df_ttTolookup.iloc[21910:30000,].index
          24
          25  # df_29290_36177 = df_ttTolookup.iloc[29290:36177,]
          26  # df_29290_36177["index_O"] = df_ttTolookup.iloc[29290:36177,].index
          27
          28
          29  # df_36177_39999 = df_ttTolookup.iloc[36177:40000,]
          30  # df_36177_39999["index_O"] = df_ttTolookup.iloc[36177:40000,].index
          31
          32  # df_40000_40629 = df_ttTolookup.iloc[40000:40629,]
          33  # df_40000_40629["index_O"] = df_ttTolookup.iloc[40000:40629,].index
          34
          35  # df_40628_45000 = df_ttTolookup.iloc[40629:45000,]
          36  # df_40628_45000["index_O"] = df_ttTolookup.iloc[40629:45000,].index
          37
          38  # df_45000_47252 = df_ttTolookup.iloc[45000:47252,]
          39  # df_45000_47252["index_O"] = df_ttTolookup.iloc[45000:47252,].index
          40
          41  # df_47252_47778 = df_ttTolookup.iloc[47252:47778,]
          42  # df_47252_47778["index_O"] = df_ttTolookup.iloc[47252:47778,].index
          43
          44  # df_47778_48108 = df_ttTolookup.iloc[47778:48108,]
          45  # df_47778_48108["index_O"] = df_ttTolookup.iloc[47778:48108,].index
          46
          47  # df_48108_49999 = df_ttTolookup.iloc[48108:50000,]
          48  # df_48108_49999["index_O"] = df_ttTolookup.iloc[48108:50000,].index
          49
          50  # df_50000_51418 = df_ttTolookup.iloc[50000:51418,]
          51  # df_50000_51418["index_O"] = df_ttTolookup.iloc[50000:51418,].index
          52
          53
          54  # df_51418_52224 = df_ttTolookup.iloc[51418:52224,]
          55  # df_51418_52224["index_O"] = df_ttTolookup.iloc[51418:52224,].index
          56
```

```
57
58   # df_52224_52561 = df_ttTolookup.iloc[52224:52561,]
59   # df_52224_52561["index_0"] = df_ttTolookup.iloc[52224:52561,].index
60
61   # df_52561_55000 = df_ttTolookup.iloc[52561:55000,]
62   # df_52561_55000["index_0"] = df_ttTolookup.iloc[52561:55000,].index
63
64   # df_55000_60000 = df_ttTolookup.iloc[55000:60000,]
65   # df_55000_60000["index_0"] = df_ttTolookup.iloc[55000:60000,].index
66
67
68   # df_60000_63054 = df_ttTolookup.iloc[60000:,]
69   # df_60000_63054["index_0"] = df_ttTolookup.iloc[60000:,].index
70
71
72
73
```

In [ ]:

```python
# step 1 get top line
from bs4 import BeautifulSoup as bs
import requests as rq # get url
import re

df_Summary_financials = pd.DataFrame(columns = ["tconst", "Index_O", "Domest
df_Summary_details = pd.DataFrame(columns = ["tconst", "Domestic Distributor
                                            "Earliest Release Date", "MPAA"
#if True:
#df_ttTolookup

for index, row in df_60000_63054.iterrows():
    #print(index)
    #print (row)

#for index, row in df_IMDB_Title_Akas.iterrows():
    #print(row["tconst"])
    #print(row["index_O"])
    #break
#      #current_ttconst = 'tt2527338'

    current_ttconst = row["tconst"]
    ddToLookupIndex = row["index_O"]

    summaryheaders = []
    summaryfinancials = []
    summarydata = {}

    detailheaders = []
    detaildata = []
    detailsummary = {}

    the_getString = 'https://www.boxofficemojo.com/title/'+current_ttconst
    r=rq.get(the_getString)
    p=bs(r.text,'html.parser')

    summaryheaders.append("tconst")
    summaryheaders.append("index")

    summaryfinancials.append(current_ttconst)
    summaryfinancials.append(ddToLookupIndex)

    # get Summary data
    b=p.find('div', class_="a-section a-spacing-none mojo-performance-summar
    if b:
        b=p.find('div', class_="a-section a-spacing-none mojo-performance-su
        divs=b.find_all('div', class_="a-section a-spacing-none")
        if divs:
        #append keys
            for div in divs:
                spans = div.find_all('span', class_=["a-size-small","money"]
                if spans:
                    for span in spans:
                        #print(span.text.strip())

                        #remove () and from values so you can have consisten
```

```
57                            thedata = re.sub(r'\([^)]*\)', '', span.text.strip()
58
59                            # must be a # remove $ signes so we can sum later
60                            if re.sub('[^0-9]',"", thedata):
61                                finanicals=re.sub('[^0-9-]',"", thedata)
62                                finanicals = int(finanicals)
63                                summaryfinancials.append(finanicals)
64                            else :
65                                #if no # in detected above must be a header so i
66                                summaryheaders.append(thedata)
67
68        summarydata = dict(zip(summaryheaders, summaryfinancials))
69        df_Summary_financials = df_Summary_financials.append(summarydata, ignore
70                #print(type(summarydata))
```

In [ ]:
```
df_ttTolookup = pd.read_excel("df_ttTolookup_DomensticWithSales.xlsx")
```

```
In [ ]:    1  #_____
           2
           3  # step 2 get detailed data
           4  from bs4 import BeautifulSoup as bs
           5  import requests as rq # get url
           6  import re
           7
           8  df_Summary_financials = pd.DataFrame(columns = ["tconst", "Domestic", "Inter
           9  df_Summary_details = pd.DataFrame(columns = ["tconst", "Index_O" "Domestic D
          10                                                "Earliest Release Date", "MPAA"
          11  #if True:
          12  #df_ttTolookup
          13
          14  for index, row in df_ttTolookup.iterrows():
          15      #print(index)
          16      #print (row)
          17
          18  #for index, row in df_IMDB_Title_Akas.iterrows():
          19      #print(row["tconst"])
          20      #print(row["index_O"])
          21      #break
          22  #     #current_ttconst = 'tt2527338'
          23
          24      current_ttconst = row["tconst"]
          25  #     ddToLookupIndex = row["Index_O"]
          26
          27      detailheaders = []
          28      detaildata = []
          29      detailsummary = {}
          30
          31      the_getString = 'https://www.boxofficemojo.com/title/'+current_ttconst
          32      r=rq.get(the_getString)
          33      p=bs(r.text,'html.parser')
          34
          35
          36      detailheaders.append("tconst")
          37  #     detailheaders.append("Index_O)
          38      detaildata.append(current_ttconst)
          39      #detaildata.append(ddToLookupIndex)
          40
          41
          42      b=p.find('div', class_="a-section a-spacing-none mojo-summary-values moj
          43      if b:
          44          if b.find_all('span'):
          45              spans=b.find_all('span')
          46              #print(len(spans))
          47              iteration = 1
          48              for span in spans:
          49                  # get rid of "a" tags
          50                  if span.a:
          51                      next
          52                  # get rid of sub span tags
          53                  if span.span:
          54                      next
          55                  else:
          56                      if (iteration % 2) == 0:
```

```
57                      iteration += 1
58                      # this is even - meaning a detail row
59                      detaildata.append(span.text.strip())
60                  else:
61                      # this is odd - meaning a headerrow
62                      #print(span.text)
63                      iteration += 1
64                      detailheaders.append(span.text.strip())
65                      #data=[item for item in data]
66                      #data=[re.sub('[^0-9]',"", str(item)) for item in da

68          detailsummary = dict(zip(detailheaders, detaildata))
69          #alldata.append(ttconst)
70          #alldata.append(summarydata)
71          df_Summary_details = df_Summary_details.append(detailsummary, ignore_ind
72
```

In [ ]:
```
1  df_Summary_details.to_excel("df_Summary_details.xlsx", header=True, index=Tr
```