```
In [20]:   1  import json
           2  import pandas as pd
           3  import numpy as np
           4
           5  import re
           6
           7  from sqlalchemy import create_engine
           8  import psycopg2
           9
          10  from config import db_password
          11
          12  import time
```

```
In [21]:   1  # 1. Add the clean movie function that takes in the argument, "movie".
           2  def clean_movie(movie):
           3      movie = dict(movie) #create a non-destructive copy
           4      alt_titles = {}
           5      # combine alternate titles into one list
           6      for key in ['Also known as','Arabic','Cantonese','Chinese','French'
           7                  'Hangul','Hebrew','Hepburn','Japanese','Literally',
           8                  'Mandarin','McCune-Reischauer','Original title','Polish
           9                  'Revised Romanization','Romanized','Russian',
          10                  'Simplified','Traditional','Yiddish']:
          11          if key in movie:
          12              alt_titles[key] = movie[key]
          13              movie.pop(key)
          14      if len(alt_titles) > 0:
          15          movie['alt_titles'] = alt_titles
          16
          17      # merge column names
          18      def change_column_name(old_name, new_name):
          19          if old_name in movie:
          20              movie[new_name] = movie.pop(old_name)
          21      change_column_name('Adaptation by', 'Writer(s)')
          22      change_column_name('Country of origin', 'Country')
          23      change_column_name('Directed by', 'Director')
          24      change_column_name('Distributed by', 'Distributor')
          25      change_column_name('Edited by', 'Editor(s)')
          26      change_column_name('Length', 'Running time')
          27      change_column_name('Original release', 'Release date')
          28      change_column_name('Music by', 'Composer(s)')
          29      change_column_name('Produced by', 'Producer(s)')
          30      change_column_name('Producer', 'Producer(s)')
          31      change_column_name('Productioncompanies ', 'Production company(s)')
          32      change_column_name('Productioncompany ', 'Production company(s)')
          33      change_column_name('Released', 'Release Date')
          34      change_column_name('Release Date', 'Release date')
          35      change_column_name('Screen story by', 'Writer(s)')
          36      change_column_name('Screenplay by', 'Writer(s)')
          37      change_column_name('Story by', 'Writer(s)')
          38      change_column_name('Theme music composer', 'Composer(s)')
          39      change_column_name('Written by', 'Writer(s)')
          40
          41      return movie
```

In [22]:

```python
# 2 Add the function that takes in three arguments;
# Wikipedia data, Kaggle metadata, and MovieLens rating data (from Kag

def movies_function():
    # Read in the kaggle metadata and MovieLens ratings CSV files as P
    kaggle_metadata = pd.read_csv ('movies_metadata.csv', low_memory =
    ratings = pd.read_csv('ratings.csv')

    kaggle_metadata_df = pd.DataFrame(kaggle_metadata)
    ratings_df = pd.DataFrame(ratings)

    # Open and read the Wikipedia data JSON file.
    file_dir = "/Users/caroline/Documents/Data Boot Camp/Module 8/Movi
    with open(f'{file_dir}/wikipedia-movies.json', mode='r') as file:
        wiki_movies_raw = json.load(file)

    # 3. Write a list comprehension to filter out TV shows.
    wiki_movies = [movie for movie in wiki_movies_raw if 'No. of episo

    # 4. Write a list comprehension to iterate through the cleaned wik
    # and call the clean_movie function on each movie.
    clean_wiki_movies = [clean_movie(movie) for movie in wiki_movies]

    # 5. Read in the cleaned movies list from Step 4 as a DataFrame.
    wiki_movies_df = pd.DataFrame(clean_wiki_movies)

    # 6. Write a try-except block to catch errors while extracting the
    #  dropping any imdb_id duplicates. If there is an error, capture
    try:
        wiki_movies_df['imdb_id'] = wiki_movies_df['imdb_link'].str.ex

        wiki_movies_df.drop_duplicates(subset = 'imdb_id', inplace = T

    except:
        print("This is an error from step 6")

    #  7. Write a list comprehension to keep the columns that don't ha
    wiki_columns_to_keep = [column for column in wiki_movies_df.column
                            < len(wiki_movies_df) * 0.9]
    wiki_movies_df = wiki_movies_df[wiki_columns_to_keep]

    # 8. Create a variable that will hold the non-null values from the
    box_office = wiki_movies_df['Box office'].dropna()

    # 9. Convert the box office data created in Step 8 to string value
    box_office[box_office.map(lambda x: type(x) != str)]

    # 10. Write a regular expression to match the six elements of "for
    form_one = r'\$\d+\.?\d*\s*[mb]illion'
    matches_form_one = box_office.str.contains(form_one, flags=re.IGNO

    # 11. Write a regular expression to match the three elements of "f
    form_two = r'\$\d{1,3}(?:,\d{3})+'
    matches_form_two = box_office.str.contains(form_two, flags=re.IGNO

    # 12. Add the parse_dollars function.
```

```python
57      def parse_dollars(s):
58          if type(s) != str:
59              return np.nan
60
61          if re.match(r'\$\s*\d+\.?\d*\s*milli?on', s, flags=re.IGNORECA
62              s = re.sub('\$|\s|[a-zA-Z]','', s)
63              value = float(s) * 10**6
64              return value
65
66          elif re.match(r'\$\s*\d+\.?\d*\s*billi?on', s, flags=re.IGNORE(
67              s = re.sub('\$|\s|[a-zA-Z]','', s)
68              value = float(s) * 10**9
69              return value
70
71          elif re.match(r'\$\s*\d{1,3}(?:[,\.]\d{3})+(?!\s[mb]illion)',
72              s = re.sub('\$|,','', s)
73              value = float(s)
74              return value
75
76          else:
77              return np.nan
78
79      # 13. Clean the box office column in the wiki_movies_df DataFrame.
80      wiki_movies_df['box_office'] = box_office.str.extract(f'({form_one
81      wiki_movies_df.drop('Box office', axis=1, inplace=True)
82
83      # 14. Clean the budget column in the wiki_movies_df DataFrame.
84      budget = wiki_movies_df['Budget'].dropna().apply(lambda x: ' '.joi
85      budget = budget.str.replace(r'\$.*[-—](?![a-z])', '$', regex=True
86      budget = budget.str.replace(r'\[\d+\]\s*', '')
87      wiki_movies_df['budget'] = budget.str.extract(f'({form_one}|{form_
88
89      # 15. Clean the release date column in the wiki_movies_df DataFram
90      release_date = wiki_movies_df['Release date'].dropna().apply(lambd
91      date_form_one = r'(?:January|February|March|April|May|June|July|Au
92      date_form_two = r'\d{4}.[01]\d.[123]\d'
93      date_form_three = r'(?:January|February|March|April|May|June|July|
94      date_form_four = r'\d{4}'
95      wiki_movies_df['release_date'] = pd.to_datetime(release_date.str.e
96
97      # 16. Clean the running time column in the wiki_movies_df DataFram
98      running_time = wiki_movies_df['Running time'].dropna().apply(lambd
99      running_time_extract = running_time.str.extract(r'(\d+)\s*ho?u?r?s
100     running_time_extract = running_time_extract.apply(lambda col: pd.t(
101     wiki_movies_df['running_time'] = running_time_extract.apply(lambda
102     wiki_movies_df.drop('Running time', axis=1, inplace=True)
103
104     # Return three variables. The first is the wiki_movies_df DataFram
105     return wiki_movies_df, kaggle_metadata, ratings
```

```
In [23]:   1  # 17. Create the path to your file directory and variables for the thre
           2  file_dir = "/Users/caroline/Documents/Data Boot Camp/Module 8/Movies-ET
           3  # The Wikipedia data
           4  wiki_file = f'{file_dir}/wikipedia.movies.json'
           5  # The Kaggle metadata
           6  kaggle_file = f'{file_dir}/movies_metadata.csv'
           7  # The MovieLens rating data.
           8  ratings_file = f'{file_dir}/ratings.csv'
```

```
In [24]:   1  # 18. Set the three variables equal to the function created in D1.
           2  wiki_file, kaggle_file, ratings_file = movies_function()
```

```
In [25]:   1  # 19. Set the wiki_movies_df equal to the wiki_file variable.
           2  wiki_movies_df = wiki_file
```

```
In [26]:   1  # 20. Check that the wiki_movies_df DataFrame looks like this.
           2  wiki_movies_df.head()
```

Out[26]:

| | url | year | imdb_link | tit |
|---|---|---|---|---|
| 0 | https://en.wikipedia.org/wiki/The_Adventures_o... | 1990.0 | https://www.imdb.com/title/tt0098987/ | Th Adventure of Fo Fairlar |
| 1 | https://en.wikipedia.org/wiki/After_Dark,_My_S... | 1990.0 | https://www.imdb.com/title/tt0098994/ | After Dar My Swe |
| 2 | https://en.wikipedia.org/wiki/Air_America_(film) | 1990.0 | https://www.imdb.com/title/tt0099005/ | A Americ |
| 3 | https://en.wikipedia.org/wiki/Alice_(1990_film) | 1990.0 | https://www.imdb.com/title/tt0099012/ | Ali |
| 4 | https://en.wikipedia.org/wiki/Almost_an_Angel | 1990.0 | https://www.imdb.com/title/tt0099018/ | Almost a Ang |

5 rows × 23 columns

```
In [27]:  1  # 21. Check that wiki_movies_df DataFrame columns are correct.
          2  wiki_movies_df.columns.to_list()
```

```
Out[27]: ['url',
          'year',
          'imdb_link',
          'title',
          'Based on',
          'Starring',
          'Cinematography',
          'Release date',
          'Country',
          'Language',
          'Budget',
          'Director',
          'Distributor',
          'Editor(s)',
          'Composer(s)',
          'Producer(s)',
          'Production company(s)',
          'Writer(s)',
          'imdb_id',
          'box_office',
          'budget',
          'release_date',
          'running_time']
```

```
In [ ]:  1
```