

```
In [1]: 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 [2]: 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 [10]: 1 # 1 Add the function that takes in three arguments;
2 # Wikipedia data, Kaggle metadata, and MovieLens rating data (from Kag
3
4 def movies_function():
5     # Read in the kaggle metadata and MovieLens ratings CSV files as P
6     kaggle_metadata = pd.read_csv('movies_metadata.csv', low_memory =
7     ratings = pd.read_csv('ratings.csv')
8
9     kaggle_metadata_df = pd.DataFrame(kaggle_metadata)
10    ratings_df = pd.DataFrame(ratings)
11
12    # Open and read the Wikipedia data JSON file.
13    file_dir = "/Users/caroline/Documents/Data Boot Camp/Module 8/Movi
14    with open(f'{file_dir}/wikipedia-movies.json', mode='r') as file:
15        wiki_movies_raw = json.load(file)
16
17    # 3. Write a list comprehension to filter out TV shows.
18    wiki_movies = [movie for movie in wiki_movies_raw if 'No. of episo
19
20    # 4. Write a list comprehension to iterate through the cleaned wik
21    # and call the clean_movie function on each movie.
22    clean_wiki_movies = [clean_movie(movie) for movie in wiki_movies]
23
24    # 5. Read in the cleaned movies list from Step 4 as a DataFrame.
25    wiki_movies_df = pd.DataFrame(clean_wiki_movies)
26
27    # 6. Write a try-except block to catch errors while extracting the
28    # dropping any imdb_id duplicates. If there is an error, capture
29    try:
30        wiki_movies_df['imdb_id'] = wiki_movies_df['imdb_link'].str.ex
31
32        wiki_movies_df.drop_duplicates(subset = 'imdb_id', inplace = T
33
34    except:
35        print("This is an error from step 6")
36
37    # 7. Write a list comprehension to keep the columns that don't ha
38    wiki_columns_to_keep = [column for column in wiki_movies_df.column
39        < len(wiki_movies_df) * 0.9]
40    wiki_movies_df = wiki_movies_df[wiki_columns_to_keep]
41
42    # 8. Create a variable that will hold the non-null values from the
43    box_office = wiki_movies_df['Box office'].dropna()
44
45    # 9. Convert the box office data created in Step 8 to string value
46    box_office[box_office.map(lambda x: type(x) != str)]
47
48    # 10. Write a regular expression to match the six elements of "for
49    form_one = r'\$\d+\.?\d*\s*[mb]illion'
50    matches_form_one = box_office.str.contains(form_one, flags=re.IGNO
51
52    # 11. Write a regular expression to match the three elements of "f
53    form_two = r'\$\d{1,3}(?:,\d{3})+'
54    matches_form_two = box_office.str.contains(form_two, flags=re.IGNO
55
56    # 12. Add the parse_dollars function.

```

```

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.IGNORECASE):
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.IGNORECASE):
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)', s):
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: ' '.join
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 DataFrame.
90 release_date = wiki_movies_df['Release date'].dropna().apply(lambda x:
91 date_form_one = r'(\?:January|February|March|April|May|June|July|August|
92 date_form_two = r'\d{4}\.[01]\d\.[123]\d'
93 date_form_three = r'(\?:January|February|March|April|May|June|July|August|
94 date_form_four = r'\d{4}'
95 wiki_movies_df['release_date'] = pd.to_datetime(release_date.str.extract
96
97 # 16. Clean the running time column in the wiki_movies_df DataFrame.
98 running_time = wiki_movies_df['Running time'].dropna().apply(lambda x:
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.to
101 wiki_movies_df['running_time'] = running_time_extract.apply(lambda x:
102 wiki_movies_df.drop('Running time', axis=1, inplace=True)
103
104 # 2. Clean the Kaggle metadata.
105 kaggle_metadata = kaggle_metadata[kaggle_metadata['adult'] == 'False']
106 kaggle_metadata['video'] = kaggle_metadata['video'] == 'True'
107 kaggle_metadata['budget'] = kaggle_metadata['budget'].astype(int)
108 kaggle_metadata['id'] = pd.to_numeric(kaggle_metadata['id'], error
109 kaggle_metadata['popularity'] = pd.to_numeric(kaggle_metadata['popularity'], error
110 kaggle_metadata['release_date'] = pd.to_datetime(kaggle_metadata['release_date'])
111
112 # 3. Merged the two DataFrames into the movies DataFrame.
113 movies_df = pd.merge(wiki_movies_df, kaggle_metadata, on='imdb_id')

```

```

114
115 # 4. Drop unnecessary columns from the merged DataFrame.
116 movies_df.drop(columns=['title_wiki', 'release_date_wiki', 'Language
117
118 # 5. Add in the function to fill in the missing Kaggle data.
119 def fill_missing_kaggle_data(df, kaggle_column, wiki_column):
120     df[kaggle_column] = df.apply(lambda row: row[wiki_column] if r
121     df.drop(columns=wiki_column, inplace=True)
122
123 # 6. Call the function in Step 5 with the DataFrame and columns as
124 fill_missing_kaggle_data(movies_df, 'runtime', 'running_time')
125 fill_missing_kaggle_data(movies_df, 'budget_kaggle', 'budget_wiki'
126 fill_missing_kaggle_data(movies_df, 'revenue', 'box_office')
127
128 # 7. Filter the movies DataFrame for specific columns.
129 movies_df = movies_df.loc[:, ['imdb_id', 'id', 'title_kaggle', 'origi
130                             'runtime', 'budget_kaggle', 'revenue', 'release_da
131                             'genres', 'original_language', 'overview', 'spoken
132                             'production_companies', 'production_countries', '
133                             'Producer(s)', 'Director', 'Starring', 'Cinematogr
134                             ]]
135
136 # 8. Rename the columns in the movies DataFrame.
137 movies_df.rename({'id': 'kaggle_id',
138                  'title_kaggle': 'title',
139                  'url': 'wikipedia_url',
140                  'budget_kaggle': 'budget',
141                  'release_date_kaggle': 'release_date',
142                  'Country': 'country',
143                  'Distributor': 'distributor',
144                  'Producer(s)': 'producers',
145                  'Director': 'director',
146                  'Starring': 'starring',
147                  'Cinematography': 'cinematography',
148                  'Editor(s)': 'editors',
149                  'Writer(s)': 'writers',
150                  'Composer(s)': 'composers',
151                  'Based on': 'based_on'
152                  }, axis='columns', inplace=True)
153
154 # 9. Transform and merge the ratings DataFrame.
155 rating_counts = ratings.groupby(['movieId', 'rating'], as_index=False
156                               .rename({'userId': 'count'}, axis=1) \
157                               .pivot(index='movieId', columns='rating', values='count'
158 rating_counts.columns = ['rating_' + str(col) for col in rating_co
159 movies_with_ratings_df = pd.merge(movies_df, rating_counts, left_o
160 movies_with_ratings_df[rating_counts.columns] = movies_with_rating

```

```
In [13]: 1 #creating a SQL database
2 db_string = f"postgres://postgres:{db_password}@127.0.0.1:5432/movies_d
3 engine = create_engine(db_string)
4 movies_df.to_sql(name='movies', con=engine)
5
6 rows_imported = 0
7 start_time = time.time()
8 for data in pd.read_csv('ratings.csv', chunksize=1000000):
9     print(f'importing rows {rows_imported} to {rows_imported + len(data)
10         data.to_sql(name='ratings', con=engine, if_exists='append')
11         rows_imported += len(data)
12
13     print(f'Done. {time.time() - start_time} total seconds elapsed')
```

```
importing rows 0 to 1000000...Done. 252.13348007202148 total seconds elapsed
importing rows 1000000 to 2000000...Done. 513.536689043045 total seconds elapsed
importing rows 2000000 to 3000000...Done. 922.1973788738251 total seconds elapsed
importing rows 3000000 to 4000000...Done. 1332.1677091121674 total seconds elapsed
importing rows 4000000 to 5000000...Done. 1619.4044170379639 total seconds elapsed
importing rows 5000000 to 6000000...Done. 1927.9038622379303 total seconds elapsed
importing rows 6000000 to 7000000...Done. 2185.834051847458 total seconds elapsed
importing rows 7000000 to 8000000...Done. 2467.8625650405884 total seconds elapsed
importing rows 8000000 to 9000000...Done. 2729.9847791194916 total seconds elapsed
importing rows 9000000 to 10000000...Done. 2983.683515071869 total seconds elapsed
importing rows 10000000 to 11000000...Done. 3214.713099002838 total seconds elapsed
importing rows 11000000 to 12000000...Done. 3452.246883869171 total seconds elapsed
importing rows 12000000 to 13000000...Done. 3677.7869729995728 total seconds elapsed
importing rows 13000000 to 14000000...Done. 3879.994642972946 total seconds elapsed
importing rows 14000000 to 15000000...Done. 4225.885064125061 total seconds elapsed
importing rows 15000000 to 16000000...Done. 4503.6516308784485 total seconds elapsed
importing rows 16000000 to 17000000...Done. 4778.14617395401 total seconds elapsed
importing rows 17000000 to 18000000...Done. 5043.797405004501 total seconds elapsed
importing rows 18000000 to 19000000...Done. 5295.373776912689 total seconds elapsed
importing rows 19000000 to 20000000...Done. 5541.271307229996 total seconds elapsed
importing rows 20000000 to 21000000...Done. 5775.770133018494 total seconds elapsed
```

```
importing rows 21000000 to 22000000...Done. 5950.776211023331 total seconds elapsed
importing rows 22000000 to 23000000...Done. 6132.967226028442 total seconds elapsed
importing rows 23000000 to 24000000...Done. 6302.955336093903 total seconds elapsed
importing rows 24000000 to 25000000...Done. 6479.958534002304 total seconds elapsed
importing rows 25000000 to 26000000...Done. 6669.753993034363 total seconds elapsed
importing rows 26000000 to 26024289...Done. 6676.92467212677 total seconds elapsed
```

```
In [14]: 1 # 10. Create the path to your file directory and variables for the three
2 file_dir = "/Users/caroline/Documents/Data Boot Camp/Module 8/Movies-ETL"
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 [ ]: 1 # 11. Set the three variables equal to the function created in D1.
2 wiki_file, kaggle_file, ratings_file = movies_function()
```