SI 206 Final Project:

Common Words and Sentiment Analysis in Various Forms of Media Titles

Carolyn Cullen, Eugene Cousino

University of Michigan, School of Information

SI 206: Data-Oriented Programming

Dr. Barbara Ericson

December 12, 2023

GitHub link: https://github.com/genocous/cousino_cullen_final.git

I.    **Project Goal**

Originally, we planned to use the Spotify API, Apple Music API, and Google Books API.

We planned on gathering title and genre information in order to find the most common

words in titles grouped by genre.

II.    **Achieved Goals**

During our initial attempts at gathering data, we realized that the APIs we chose would

not be sufficient for all calculations. So, we instead used the Spotify API, IMDb API, and

Open Library API. We successfully found the most common words in titles by genre and

decided to take our project a step further by performing sentiment analysis on all titles by

genre.

III.    **The problems that you faced**

We faced a number of problems, including in our initial stages when we discovered a few

things. One, the Apple Music API required $99.99 developer "membership" in order to

pull data from the API. We had to find a different API and ended up using the IMDb API

instead. We also found that the Google Books API was hard to work with, so we used the

Open Library API instead, and it was much easier to use. We also dealt with significant

runtime issues of data collection from our API's, specifically the IMDb API. We

considered switching to another database, but ultimately decided against it, because we

only had to run the code a few times to get the data that we needed. Once we had the

data, we had to deal with cleaning it. We created a list of excluded words for each source

of media, because they skewed the results and made it harder to analyze the content. We

also had to ensure that we were using English titles, because particularly in the Spotify

database, there were many foreign tracks, which were not helpful to us.

### IV.    Calculations

Music Results:

```
Most common words in Pop song titles:    Most common words in Rap song titles:    Most common words in Hip hop song titles:
love: 23                                 love: 9                                   baby: 17
vault: 15                                savage: 8                                 young: 7
girl: 11                                 remix: 7                                  broke: 6
good: 9                                  bad: 6                                    never: 6
remix: 7                                 go: 5                                     rich: 5
night: 6                                 around: 4                                 time: 5
last: 5                                  party: 4                                  remix: 5
know: 5                                  money: 4                                  love: 5
christmas: 5                             polo: 4                                   best: 4
bonus: 5                                 girl: 4                                   future: 4
Most common words in Alt song titles:    Most common words in Country song titles: Most common words in Indie song titles:
love: 7                                  good: 7                                   love: 4
boy: 3                                   one: 6                                    shut: 3
interlude: 3                             country: 5                                close: 3
little: 3                                whiskey: 5                                moms: 2
hate: 3                                  girl: 5                                   calling: 2
forever: 3                               time: 5                                   end: 2
life: 3                                  hell: 4                                   song: 2
night: 2                                 life: 4                                   animal: 2
devil: 2                                 man: 4                                    life: 2
see: 2                                   cold: 4                                   mind: 2
```

Movie Results:

```
Most common words in Drama movie titles: Most common words in Comedy movie titles: Most common words in Action movie titles:
love: 10                                 love: 17                                  part: 4
part: 9                                  part: 15                                  death: 3
death: 8                                 bob: 8                                    force: 3
heart: 6                                 mike: 7                                   deep: 3
time: 6                                  time: 7                                   space: 3
life: 5                                  house: 6                                  king: 3
trail: 5                                 day: 5                                    love: 3
girl: 5                                  lady: 5                                   murder: 3
day: 5                                   christmas: 5                              day: 3
moment: 4                                lucy: 5                                   close: 2
Most common words in Adventure movie titles: Most common words in Family movie titles: Most common words in Documentary movie titles:
part: 4                                  part: 7                                   american: 7
boy: 2                                   magician: 3                               live: 7
runaway: 2                               house: 3                                  world: 5
challenge: 2                             hare: 3                                   story: 5
adventures: 2                            celebrities: 2                            life: 5
robin: 2                                 pressure: 2                               years: 4
hood: 2                                  friend: 2                                 john: 4
fire: 2                                  mr: 2                                     trip: 3
strange: 2                               leakey: 2                                 west: 3
magic: 2                                 daffy: 2                                  tv: 3
Most common words in Crime movie titles:
part: 12
big: 6
case: 5
time: 4
blood: 4
family: 3
run: 3
house: 2
murder: 2
blind: 2
```

Book Results:

```
Most common words in Film book titles: Most common words in Music book titles: Most common words in Fantasy book titles:
film: 52                                history: 15                              book: 10
cinema: 51                              musical: 12                              oz: 7
american: 14                            business: 7                              harry: 7
culture: 11                             dictionary: 6                            potter: 7
hollywood: 10                           musicians: 6                             magic: 7
films: 10                               western: 6                               dark: 7
history: 7                              art: 6                                   rising: 7
movies: 7                               english: 6                               doctor: 6
america: 7                              studies: 5                               dolittle: 6
politics: 6                             blues: 5                                 tales: 6
Most common words in Historical_fiction book titles: Most common words in Horror book titles: Most common words in Humor book titles:
sharpe: 12                              stories: 59                              history: 12
king: 7                                 tales: 30                                book: 8
white: 4                                short: 20                                big: 6
pimpernel: 4                            house: 19                                handbook: 6
scarlet: 3                              poems: 19                                letters: 5
red: 3                                  usher: 15                                dave: 5
black: 3                                fall: 14                                 barry: 5
lady: 3                                 night: 13                                humor: 5
brother: 3                              mystery: 12                              stories: 4
pearl: 3                                dark: 11                                 thurber: 4
Most common words in Science_fiction book titles: Most common words in Young_adult book titles:
time: 13                                book: 9
foundation: 11                          health: 9
book: 8                                 adolescents: 7
star: 8                                 death: 6
moon: 7                                 love: 5
dune: 7                                 story: 5
space: 7                                girls: 5
world: 7                                soul: 5
mars: 6                                 stories: 4
gods: 5                                 tree: 4
```
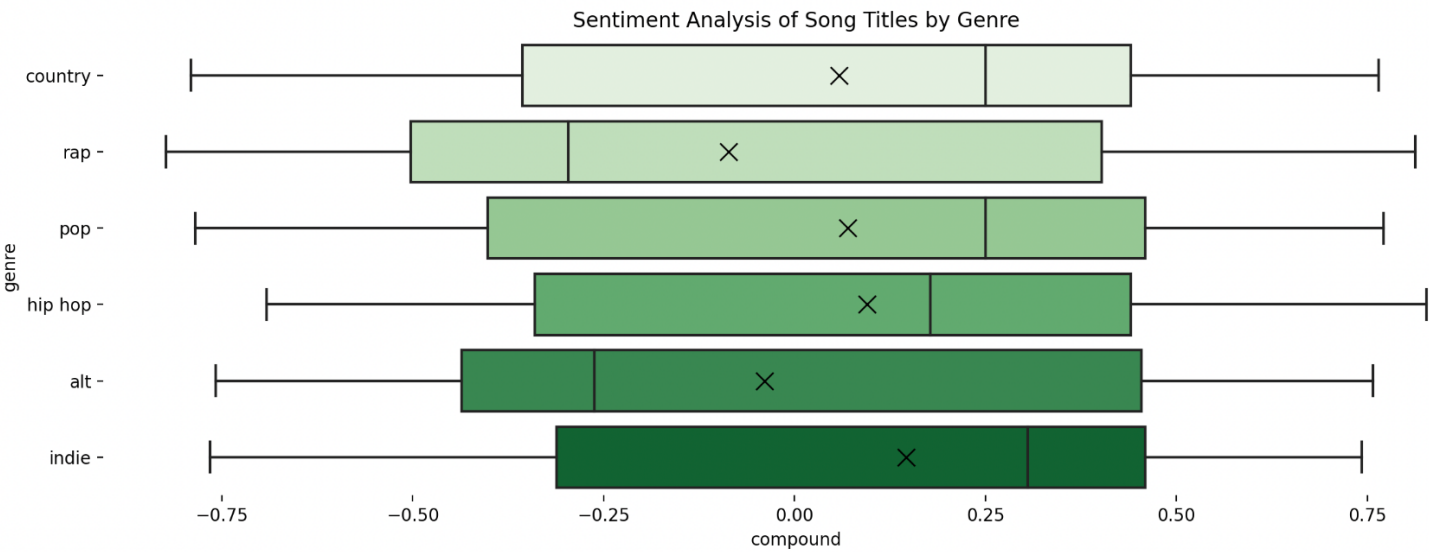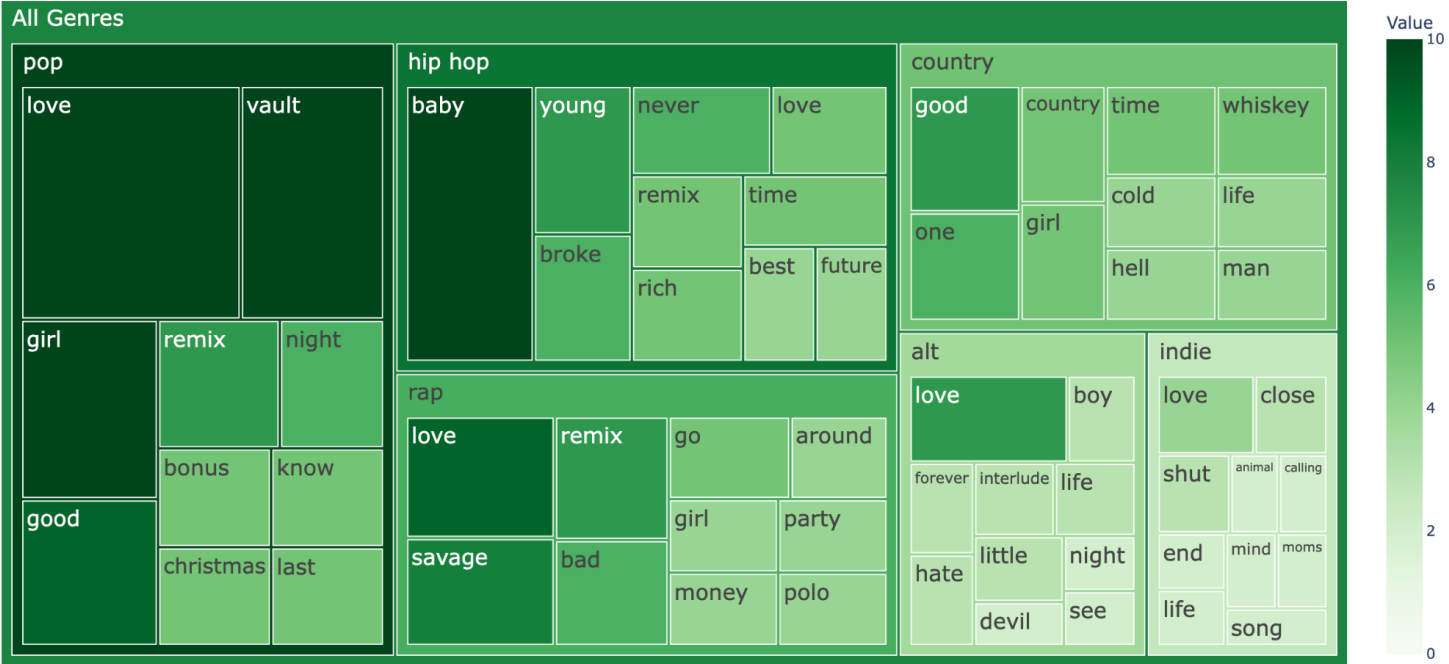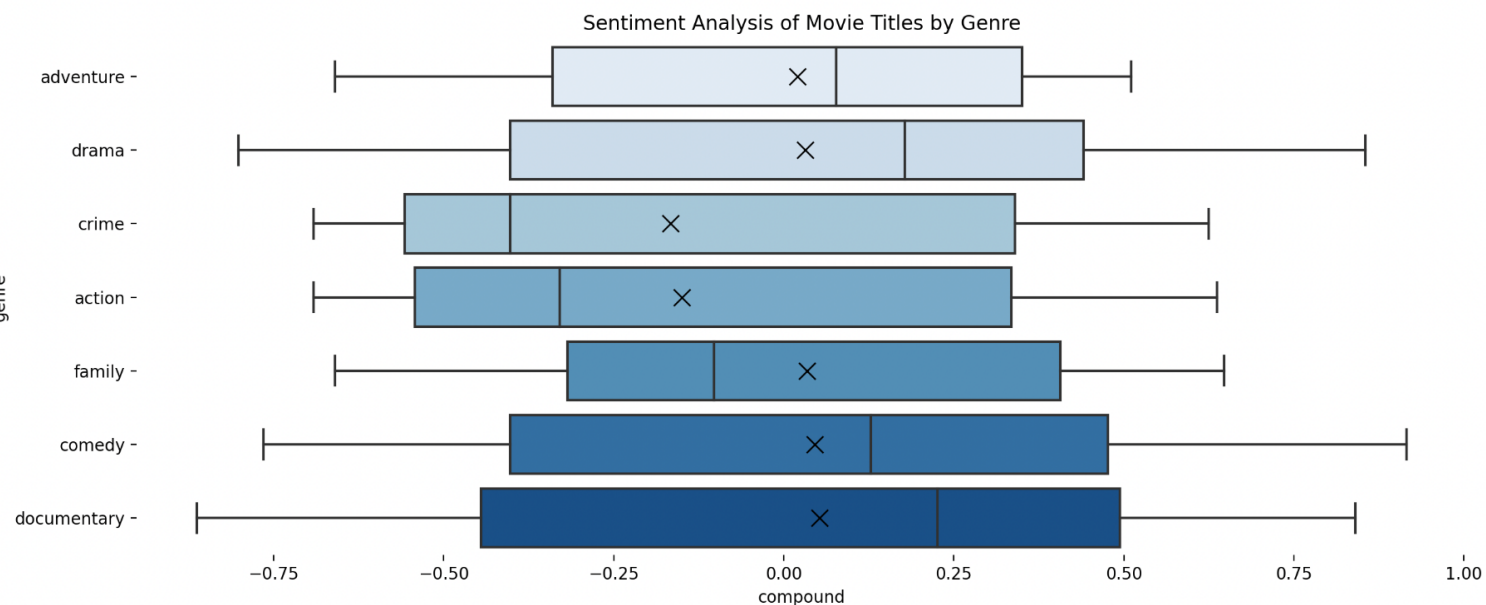
Book/Movie Horror Results:

```
Most common words in horror titles (books and movies):
stories: 1419
tales: 1290
house: 1167
death: 1087
dead: 1001
blood: 915
night: 909
one: 872
terror: 872
poe: 860
```
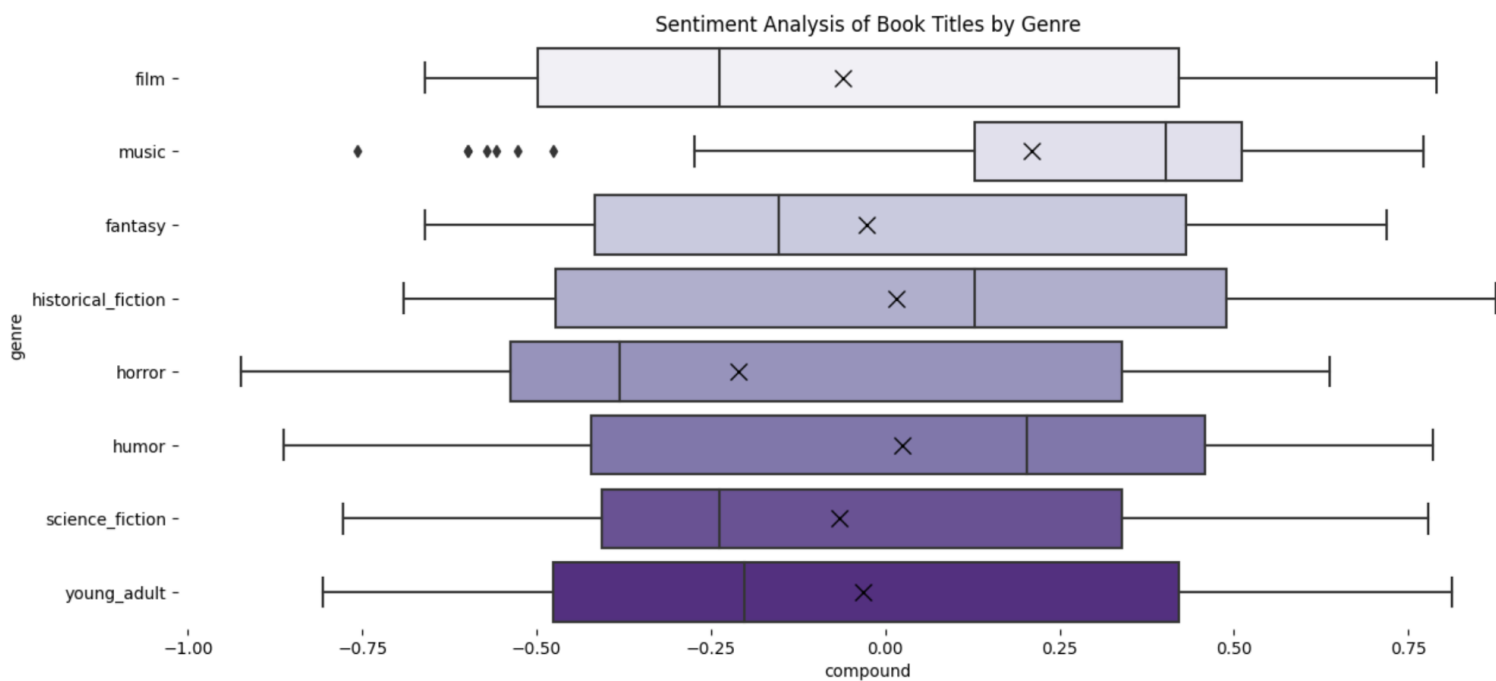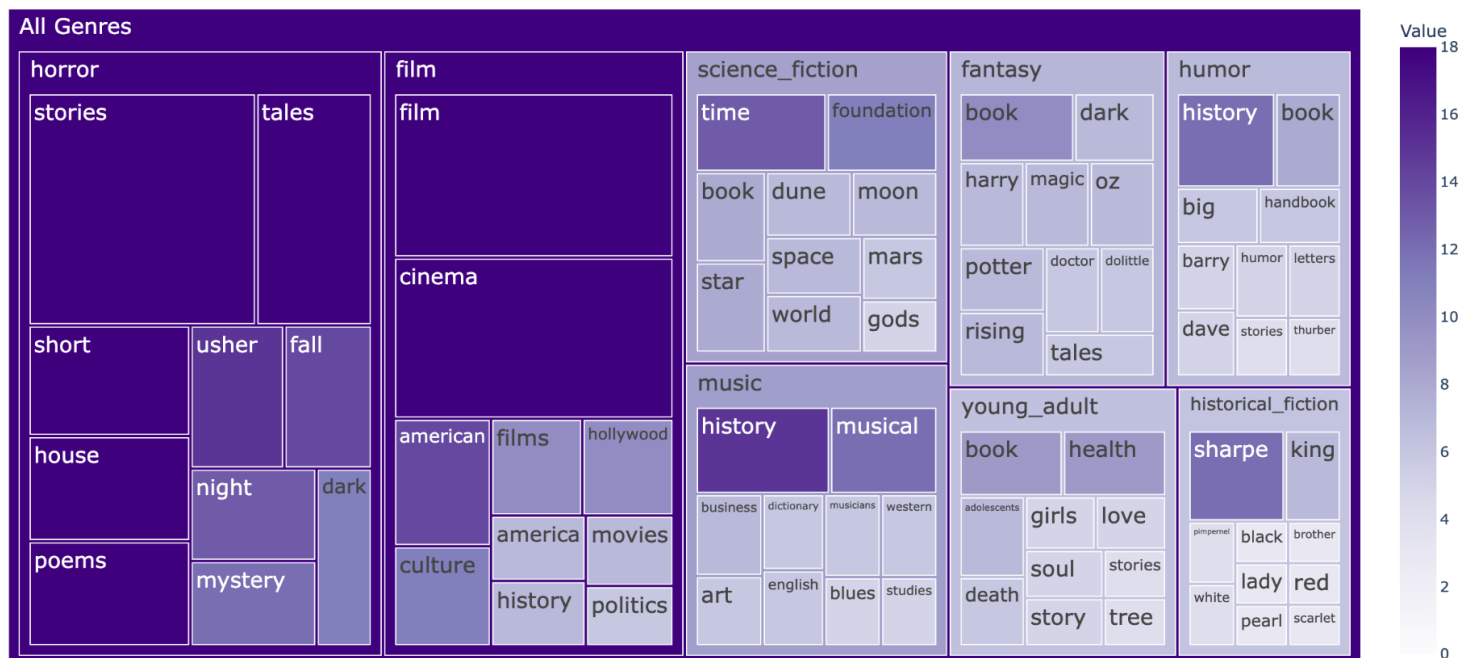
## V.    Visualizations

Common Words in Music Titles by Genre

Common Words in Movie Titles by Genre



Sentiment Analysis of Movie Titles by Genre

Common Words in Book Titles by Genre



Sentiment Analysis of Book Titles by Genre

**VII.    Instructions for running code**

The first step is collecting all of the data and placing it in a database. There are three

separate files for collecting data from each of the API's that we used. The 'spotify_api' is the file

that you should start with, and collect as much data as you want. We have it limited at 25 tracks

per runtime, but you can change that number to whatever you want it to be, and that is the case

for all of these files as well. The Spotify file puts its data in a table called tracks within the

database 'title_analysis'. You should then run the IMDb file (imdb_api), which will take around

6-7 minutes to collect 25 titles and genres. This happened every time we ran the code, so expect

to wait for the data to load. That file puts movie titles in a table called movies within the

title_analysis database. Then, run the open_library_api function, which puts 25 book titles in a

table called books within the title_analysis database. After you run these files, ensure that the

data is loaded correctly by viewing the tables in SQLite, or a related database server.

Your next step is to run the 'visualization_functions' file, which will generate treemap

visuals for each category of data (tracks, movies, books). Within this file, you can specify certain

genres that you want to see in the main() function, as well as edit words that you want excluded

from the analysis. This file also writes the top 10 most common words by genre into a text file.

Then, you can run the 'join_analysis' function, which looks at the most common words in horror

titles across both books and movies, and writes the result into a text file.

Finally, you should run the analysis_function file, which will perform the sentiment

analysis of the various titles. You can also change the genres that you want to see results for in

the main() function. The results of the analysis_function file are visualizations for each media

type. We also added another file called join_analysis, where we joined together titles from horror

movies and books to compare them.

**VIII.** **Documentation for each function that you wrote. This includes describing the input and output for each function**

## spotify_api:

create_table():

      Description: Creates a table called *tracks* to store track titles and genres

      Parameters:

        - None

      Returns:

        - None

      Example:

      create_table()

fetch_tracks():

      Description: Stores 25 tracks and their corresponding genres to the *tracks* table

      Parameters:

        - None

      Returns:

        - None

      Example:

      fetch_tracks()

## imdb_api:

create_table():

      Description: Creates a table called *movies* to store movie titles and genres

      Parameters:

        - None

      Returns:

        - None

      Example:

      create_table()

fetch_movies():

      Description: Stores 25 movies and their corresponding genres to the *movies* table

      Parameters:

        - None

      Returns:

        - None

      Example:

      fetch_movies()

**open_library_api:**
create_table():
      Description: Creates a table called books to store book titles and genres
      Parameters:
         - None
      Returns:
         - None
      Example:
      create_table()

fetch_books():
      Description: Stores 25 books and their corresponding genres to the books table
      Parameters:
         - None
      Returns:
         - None
      Example:
      fetch_books()

**visualization_functions:**
common_words(db_name, file_name, genre_keywords, excluded_words, data_row):
      Description: Create a list of the 10 most common words by genre
      Parameters:
         - db_name (table): name of the table that is providing data
         - file_name (string) : file name for the results file
         - genre_keywords (list): keywords to use as genres stored within a list
         - excluded_words (list): words that should be excluded from analysis, stored in a list
         - data_row (list): list of columns in the given table
      Returns:
         Dictionary containing the 10 most common words by genre
      Example:
      common_words(tracks, 'spotify_output.txt', gernre_keywords, excluded_words, ['id', 'name', 'genre'])

visualizations(common_words_dict, color_scheme, color_range, title)
      Description: Create a visualization of common words per genre
      Parameters:
         - common_words_dict (dict): dictionary containing the 10 most common words by genre
         - color_scheme (string): indicates which color scheme the visualization should be
         - color_range (tuple): indicates the range of the color scheme for the visualization
         - title (string): title of the visualization

Returns:

Treemap visualization of the most common words per genre

Example:

visualizations(common_words_result, 'Purples', [0,18], 'Common Words in Book Titles by Genre')

## analysis_function:

visualize_sentiment_analysis(db_name, genre_keywords, data_row, color_scheme, g_title):

Description: Create a visualization of common words per genre

Parameters:

- common_words_dict (dict): dictionary containing the 10 most common words by genre
- color_scheme (string): indicates which color scheme the visualization should be
- color_range (tuple): indicates the range of the color scheme for the visualization
- title (string): title of the visualization

Returns:

Treemap visualization of the most common words per genre

Example:

visualizations(common_words_result, 'Purples', [0,18], 'Common Words in Book Titles by Genre')

## IX.    Resources used

| Date | Issue Description | Location of Resource | Result (did it solve the issue?) |
|------|-------------------|----------------------|----------------------------------|
| 12/2/2023 | IMDb API running slow | chat.openai.com | Issue solved - prepped statement to insert data into the tables. Runtime improved by 2-3 minutes per run. |
| 12/3/2023 | General issues with sentiment analysis | https://www.kaggle.com/code/robikscube/sentiment-analysis-python-youtube-tutorial | Issue solved - learned how to use nltk SentimentIntensityAnalysis |