

CSC111 Project Report: ScreenSelect

Personalized Movie Recommendation System

Dogyu Lee, Narges Movahedian Nezhad, Sidharth Sawhney, Aastha Sharma

April 3rd, 2023

Part 1: Goal

Due to the advent of streaming services, which provide customers immediate access to a vast collection of information, the entertainment industry has undergone a vital change. Finding the right things to watch, for example, might just be more difficult than ever for viewers. This is due to the overflow of available content. As a solution to this issue, recommendation systems have been created; personalizing and recommending material based on users' viewing patterns and interests. With this in mind, **our project intends to develop a movie customization and recommendation app that takes inspiration from Netflix's recommendation system, leveraging watching patterns, user input, and clever algorithms to provide tailored movie suggestions.**

Our project intends to build movie customization that *varies* from that of others in regard to algorithm, despite the fact that the issue of content overload and the demand for tailored movie recommendations have been thoroughly examined. For instance, we will employ contextual data in addition to the users' viewing history to provide recommendations that are more customized. This goes beyond Netflix's simple recommendations of related content based only on user watch history. Additionally, our project adds a sense of individualism to every user. Unlike other recommendation platforms, we have deliberately chosen not to take into consideration the location of our consumers and users nearby because we understand that everyone has unique taste. Having geographically been near someone in no way guarantees similar preferences when it comes to movies. For this reason, we would go as far as to argue that taking location into account ruins the algorithm and displays less desirable content.

Moreover, ScreenSelect takes into account only the past 10 movies for every user, as opposed to all viewing history—which is the case with other projects in the market. This is largely due to the fact that people's tastes are likely to change over time. By doing so, we are ensuring that our recommendations are updated to appeal to consumers' most recent interests in contrast to the movies they enjoyed in the distant past. This strategy also ensures user happiness and increased engagement with the app. Also, by focusing on the viewer's most recent 10 films rather than their complete viewing history, we may prevent constantly recommending identical films, which can cause user boredom and disengagement. As the user's preferences change over time, our strategy utilizes both the user's viewing history and contextual information to provide recommendations that are more individualized. Overall, our project provides a fresh and original perspective on the issue of information overload. We strive to deliver a highly customized, enjoyable, and useful movie recommendation app to increase user engagement and commitment.

Our project's objective is to boost user engagement by means of providing a highly-customized watching experience that changes depending on each user. We will offer a platform that generates the best recommendations for each user every time by taking into account contextual elements—such as the last 10 watched movies, keywords, genre, and language—going beyond Netflix's straightforward suggestions of related material based on user watch history. We want to provide a highly engaging experience that keeps people riveted to the screen and coming back for more by putting a strong emphasis on tailored suggestions.

Our group has also made several considerations in regard to the ethics of ScreenSelect. We will put in place a number of privacy precautions to make sure that our project doesn't divulge any sensitive information. We

never request personally identifiable information when gathering user data. The username we take from our users is the only sort of identification verification input we take from them. This username can never be connected, as it can be anything users choose. This will assist in preserving user privacy and preventing the disclosure of private information. We have also put in place a number of safeguards to make sure that our program doesn't make any incorrect recommendations to vulnerable populations. Our group now has a dedicated email address on the site that we utilize to ask users for feedback, especially from members of vulnerable groups. This will make it easier to spot any offensive ideas or information. The input will help our team improve the app's suggestions and make sure they are appropriate for all users. These precautions will be regularly audited to identify potential privacy and security breaches and ensure that our app complies with privacy regulations.

We have concentrated on developing a recommendation system that prioritizes user satisfaction and engagement over other metrics like watch time or income in order to make sure that our project's ideas are valuable to the user as opposed to the Netflix engine. When creating our project, we concentrated on giving user-centered design concepts top priority. Aspects such as reading accessibility (through large fonts) can be seen in all parts of our implementation. This required comprehending the user's requirements and preferences and tailoring the platform to satisfy them. We boosted user engagement and improved the app's suggestions as a result of user testing and feedback. In addition to other helpful measures, we provide users with a wide variety of content, including works that might not be as well-known but are nevertheless well-regarded and pertinent to their interests. This is one of the reasons for our decision to work with a large data set. Users will find fresh material and engage more as a result of this.

In conclusion, the goal of our project is to develop a highly individualized, entertaining, and impactful movie recommendation app. We seek to offer consumers a watching experience that is customized to their interests and keeps them hooked with our platform over time. This is done through machine learning algorithms and contextual elements. Our ultimate purpose is to enhance user engagement and commitment. We will gauge our progress using a variety of test cases that reflect this goal.

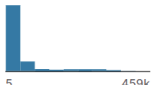
Part 2: Datasets

ScreenSelect used two datasets in the form of CSV files: `tmdb_5000_movies.csv` (Movies Dataset) and `tmdb_5000_credits.csv` (Credits Dataset). Both files were downloaded from Kaggle. The Movies Dataset consists of 20 columns with the information and statistics of 4803 movies. Out of the 20 columns in the Movies Dataset, ScreenSelect utilized a subset of nine columns: `id`, `title`, `genres`, `keywords`, `runtime`, `original_language`, `overview`, `vote_average`, and `release_date`. The Credits Dataset consists of four columns with the `id`, `title`, `cast`, and `crew` for each of 4803 movies. Out of the four columns in the Credits Dataset, ScreenSelect utilized a subset of two columns: `id` and `crew`.


We have extracted approximately 100 movies in a smaller version of the original Movies Dataset CSV file called `tmdb_5000_movies1` for regular, consistent testing purposes in a separate CSV file. We have done this because, upon testing, we found that it takes a long time to load the entire data for 4803 movies into the system. For the same reason, it also takes a long time to generate a recommendation as we have to filter through 4803 possibilities and analyze all of their content through our algorithm for the user. Neither one of these issues arise with the testing Dataset. Although we have tested our code on the larger dataset and confirmed that it performs as planned, we found that confirming everything on the smaller Dataset was much more time-effective as it provided desired results faster. For comparison, loading the entire graph with the original Movies Dataset takes around 15 minutes while loading the smaller version takes approximately 2 minutes.


We will provide both the smaller file along with the original Movies Dataset CSV file in our submissions. We suggest that tests be run on the smaller dataset, as it will take a maximum of 2 minutes to load the recommendations after the graph has been loaded. This is a significantly smaller waiting time as opposed to running the Movies Dataset.


CSV File 1

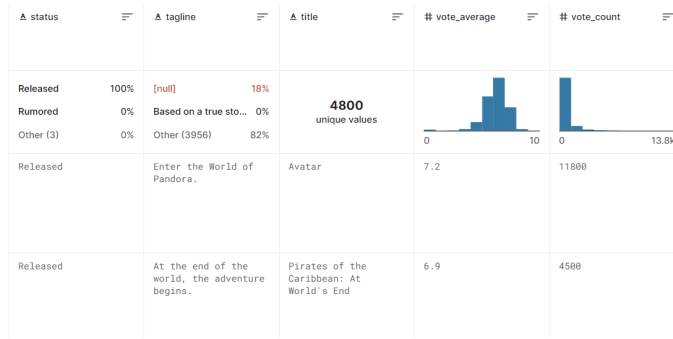
movie_id	title	cast	crew
ID of the movie	Name of the movie	Cast	Crew
	4800 unique values	4761 unique values	4776 unique values
19995	Avatar	[{"cast_id": 242, "character": "Jake Sully", "credit_id": "5602a8a7c3a3685532001c9a", "gender": 2, "...	[{"credit_id": "52fe48009251416c750aca23", "department": "Editing", "gender": 0, "id": 1721, "job": "...
285	Pirates of the Caribbean: At World's End	[{"cast_id": 4, "character": "Captain Jack Sparrow", "credit_id": "52fe4232c3a36847f800b58d", "gende...	[{"credit_id": "52fe4232c3a36847f800b579", "department": "Camera", "gender": 2, "id": 120, "job": "D...

CSV File 2

# budget	genres	homepage	id	keywords
Budget of the movie	Genres of the movie	Homepage link of the movie	ID of the movie	Keywords of the movie
	[{"id": 18, "name": "... 8% [{"id": 35, "name": "... 0% Other (4151) 86%	[null] 64% http://www.thehun... 0% Other (1708) 36%		[] 9% [{"id": 10183, "name... 1% Other (4336) 90%
237000000	[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": ...	http://www.avatarmovie.com/	19995	[{"id": 1463, "name": "culture clash"}, {"id": 2964, "name": "future"}, {"id": 3386, "name": "space ...
300000000	[{"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 28, "name": "Action"}]	http://disney.go.com/disneypictures/pirates/	285	[{"id": 270, "name": "ocean"}, {"id": 726, "name": "drug abuse"}, {"id": 911, "name": "exotic island...]

original_language	original_title	overview	popularity	production_companies
Original language of the movie	Original title of the movie	Overview of the movie	Popularity of the movie	Production companies of the movie
en 94% fr 1% Other (228) 5%	4801 unique values	4801 unique values		[] 7% [{"name": "Paramou... 1% Other (4394) 91%
en	Avatar	In the 22nd century, a paraplegic Marine is dispatched to the moon Pandora on a unique mission, but ...	150.437577	[{"name": "Ingenious Film Partners", "id": 289}, {"name": "Twentieth Century Fox Film Corporation", ...
en	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, has come back to life and is headed to the edge of the E...	139.082615	[{"name": "Walt Disney Pictures", "id": 2}, {"name": "Jerry Bruckheimer Films", "id": 130}, {"name": "...

production_count...	release_date	revenue	runtime	spoken_languages
[{"iso_3166_1": "US", "name": "United States of America"}, {"iso_3166_1": "GB", "name": "United Kingdom"}]				[{"iso_639_1": "en", "name": "English"}, {"iso_639_1": "es", "name": "Espa\u00f1ol"}]
[{"iso_3166_1": "US", "name": "United States of America"}, {"iso_3166_1": "GB", "name": "United Kingdom"}]	2009-12-10	2787065087	162	[{"iso_639_1": "en", "name": "English"}, {"iso_639_1": "es", "name": "Espa\u00f1ol"}]
[{"iso_3166_1": "US", "name": "United States of America"}]	2007-05-19	961000000	169	[{"iso_639_1": "en", "name": "English"}]



Part 3: Computational Overview

1. ScreenSelect's data is represented using a Graph class which consists of a dictionary containing all the vertices. The vertices of the graph are objects of the Vertex abstract superclass. The subclasses are different for the **movie vertex** and the **user vertex**, so each vertex represents either a user or a movie. An edge in the graph is made between the user vertex and the movie vertex when the user selects one movie from the given five movie recommendations. The selected movie becomes a neighbour of the user and vice versa, how this process works will be detailed in the following paragraphs. A maximum of 10 edges can exist between a movie and a user. The formed edges are utilized to suggest a new set of movies for an existing user. The information about the past 10 selected movies and new user preferences are considered when the user returns to the program for another movie recommendation. The existing edge also prevents the same vertices, a user and a movie, from being connected twice, avoiding duplicate suggestions.
2. The **user_log_in** function in *main_functions.py* maintains and builds the given graph. It retrieves and stores its first user input in the variable username. A while loop with a stopping condition of username not being an empty str ensures the username is valid and can be used in the program. Then, it calls the graph method **verify_vertex** to check whether the username exists in the system. If it is not in the graph's `_vertices` dictionary, it calls another graph method **add_user_vertex** to add a new user vertex in the graph. Then, a private function within *main_functions.py* is called to get user movie preference inputs and reassign the given user's preference attributes. The private function introduces a while loop with a stopping condition for keywords to ensure the user does not type more than three keywords and prevents jeopardizing the preference scores. The weights of 5 and 10 points for more important attributes such as genre and language will lose their meaning because when there are more than three matching keywords, the keywords' score will be greater than the sum of *genre* and *lang*'s score. Finally, the **user_log_in** function returns the str username for the functionality of the front end which will be elaborated upon in the 3rd portion, visualization, of the computation overview.

read_csv_and_create_data in *main_functions.py* processes movie data from csv files and add each movie as a vertex in the given graph. **read_csv_and_create_data** transforms data types such as change *vote_average* and *runtime* from str into float and int. Data filtering is done using private helper functions in the same file by breaking down a list of dictionaries of strings and looping through each element to find the correct information for each movie which is returned immediately or added to a set before returning the collection. The value from the column to be used and the returned filtered values from the helper functions for each row are assigned to variables. These variables are passed into a graph method **add_movie_vertex** which creates a new `_Movie` object and maps it to the movie title in the graph's `_vertices` dictionary.

Given a graph and a username, the **compute** function in *main_functions.py* uses `_Movie` class's method **score** to calculate the user's preference scores for all movie vertices and returns a dictionary that maps the top five matching movies for the user. With the given user vertex, the **score** method aggregates the preference points on the variable *score* by adding 5 or 10 points based on a match by comparing the user and the movie's *genre*, *lang*, *keywords*, and *director*. If the user has a past selection of recommended movies in its *past_10_neighbours* list, it applies the same procedure for each of the movies in the list. The first 5 movie vertices are mapped to their title in the dict *_top_scores* of the user vertex. However, when there are already 5 movies in *_top_scores*, only movies with scores greater than the minimum of the vertices' scores are added to *_top_scores* after removing the vertex with the minimum score. Thus, by the end of the computation, 5 movies with the highest preference scores remain in *_top_scores* which is returned by the compute function.

3. This part contains both how the PyQt6 library was used to accomplish its task and how it is responsible for the visual report of the resulting computations in ScreenSelect.

Screenselect's user interface is created using the PyQt6 library. Its `QtWidgets`, `QtCore`, and `QtGui` modules are imported and used in our program. In `screen_select_gui.py`, three subclasses, `RecommendationScreen`, `PreferenceScreen`, and `LogInScreen`, of `QtWidgets` exist. Each subclass represents a window the program displays for the user and can use `show()` or `close()` to appear or disappear while the program is running. The `QtWidgets` module's `QWidget`, `QPushButton`, `QLabel`, and `QLineEdit` classes are used to create graphical user interface (GUI) components such as buttons, labels, and textboxes, for each screen. `QLineEdit`'s `displayText()` is used multiple times to show the received input to the user. `QtWidgets` module's `QGridLayout` class and `QtCore` module's `Qt` class are employed to set a layout for each screen using `setLayout()`, to divide the screens into sections using `setSpacing()` and `setContentMargins()`, and add each widget to a layout and set their location using `addWidget()`. The `QtGui` module's `QFont` and `QIcon` classes are used to design the screens and the style of the components to be more pleasing to the eye and fit the cinematic theme of our project. Particularly, functions such as `setFont()` and `setStyleSheet()` are used to adjust the fonts and colors of label figures.

A key capability of the PyQt6 library is its ability to emit and receive signals between classes. A major constituent of each screen is the buttons the user can click on to move to the previous or the next step of the recommendation process. The users' clicks activate the system to change windows and display the steps graphically because the buttons are made by the `QtWidgets` module's `QPushButton` class. `QPushButton`'s method `clicked()` enables the button to emit a signal to a function connected to it [the connection is established using `QPushButton`'s `connect()`]. This functionality allows the user to interact with the system for each screen. On the `RecommendationScreen`, each button is connected to one of the five **screenselect** methods or a **back** method. When clicked on, the buttons will signal the functions to activate which results in the user either choosing a recommended movie or going back to `PreferenceScreen`. Similarly, the user on the `PreferenceScreen` could choose buttons that activate either the **recommendation_button** to move on to `RecommendationScreen` for movie suggestions or **log_out** to go back to the `LogInScreen`. Finally, the user on `LogInScreen` can choose either to click on **sign_in** or **sign_up** and be taken to the `PreferenceScreen`.

Part 4: Instructions for Running ScreenSelect

Part 5: Changes to Project Plan

Classes

- The `Graph` class was modified with five new public methods **verify_vertex**, **add_user_vertex**, **retrieve_item_obj**, **retrieve_vertex_dict**, and **add_movie_vertex**. The preceding methods add user and/or movie vertices to the graph and give access to its vertex objects which are mapped in a private instance attribute dictionary if other functions need it.
- The `_Vertex` superclass has changed from a class to a data class for cleaner code and not having to add any other abstract methods to it. Its instance attributes have all been changed to have a default `None` value when initialized. One of its instance attributes **genre**'s type contract changed from `str` to `Optional[set[str] | str]` to allow a set of `str` for the `_Movie` class which has more than one genre in the `Movies` dataset. A public instance attribute **neighbours** was removed from the superclass and was added separately to the `_Movie` and `_User` subclasses. Representation invariants were added to ensure the instance attributes were non-empty if not the default value.
- `_Movie` class's private instance attribute **_total_score**'s type contract changed from an `int` to a dictionary to store users' names and preference-matching scores for itself. The dictionary allowed the system to efficiently access the mapped preference scores of multiple users by their usernames for a specific movie instead of having to compute the score for each user and movie all over again. Representation invariants were added to ensure the instance attributes were not empty strings or 0.
- A new public instance attribute **past_10_neighbours** was added to the `_User` class to keep track of the recently watched/recommended movies and use them (their scores) to recommend the next movie. Two new public methods **retrieve_top_scores** and **modify_preferences** were added to this class. These methods give access to its private instance attribute dictionary **_top_scores** if other functions need it and reassign some of its public

instance attributes' values. Representation invariants were added to ensure the instance attributes were within the correct range of numbers.

Overall Program

- The plan to use a username and a password for the user login procedure changed to only using the username. Passwords were not used in the program except for the login procedure and each user object has a unique username as its key in the graph's `_vertices` dictionary. The role of passwords was not much important, so it was excluded from our program.
- The plan to display the recommended movies' posters on the user interface (screen/window) was dropped. It took longer than expected to retrieve the image from the website link and resize it to fit the UI screen.
- The plan to display the overviews of the recommended movies was dropped due to difficulties fitting different length str for each movie on a screen with other existing UI widgets in fixed positions.

Part 6: Discussion

Our goal in this project was to create a movie recommendation app that provides users with tailored recommendations based on their viewing habits and preferences. We understood that such platforms exist in the market, however, we were determined to offer the best suggestions with the highest accuracy for desire. This was done by disregarding our user's geographical location. In order to consistently provide the best suggestions for each user, we also employed a content-based filtering algorithm that takes into account contextual factors including the most recent 10 movies seen, keywords, genre, and language.

We are very proud to say that the results of our computational exploration helped us achieve our goal. Our computational investigation's findings (through self-testing) demonstrated that our system was very effective at giving consumers personalized movie selections. Compared to traditional content-based filtering algorithms that simply take into consideration the qualities of each piece of content, our system was able to provide more customized suggestions. This is largely due to the fact that it takes contextual factors and most recently watched into account. We are confident in this conclusion thanks to the countless tests we have conducted both on the larger as well as smaller datasets. We have also separately tested all possible filters for our movie dataset on our GUI to ensure that the recommendations are all successful in achieving the end goal for users. Moreover, we have tested all of this for multiple users, and every test returns the anticipated results. Numerous users can now interact with one another in terms of the movies in the dataset while each having their own separate, independent object.

However, we did encounter some limitations during the development of our app. One limitation was the difficulty in loading the movie data into the graph due to the large amount of data and the different formats in which it was presented. For example, in order to access the genre for a movie, we had to access a dictionary value that was in a list which was a string. We used for loops to access each piece and get rid of unnecessary details, but we were not able to typecast the string into different dictionaries of lists. Eventually, we went with our method of normal stripping of brackets, etc. in order to overcome this obstacle.

Another limitation was the difficulty in adding images to the graphic interface of each movie when displaying its title. As the GUI was quite small, we were not able to see how to add images without messing up the layout. Therefore, we chose to take this out and implement a clean version in the future by changing the GUI from the layout import that we have currently used to a fully custom-made file that supports images as well. We believe that this is an achievable goal to set for the future; and it would enhance our current UI even more.

One additional constraint that we encountered pertained to the PyQt6 library. A significant proportion of the features that we initially expected to employ for receiving input proved to be inadequate as there was no direct link between them and the graphical user interface (GUI). As a result, we were prompted to revise our strategy and gather input exclusively from the GUI, while gathering it as an object. We then process out computations on this created object. The basis for collecting user information differed from our original expectations, necessitating a shift in its source from the console to the GUI, with consequential ramifications.

In terms of the next steps for further exploration, we plan on adding strong security to our app by implementing encryption using concepts we learned in CSC110. We will also ask users for their ages and sort movies by age group so that children would never receive 13+/18+ content. Unfortunately, our datasets do not have age ratings for any of the movies listed, and we are limited from doing this in that regard. Our long-term plan is to cross reference all of our movie information with that of another data set that lists ages as well all in order to implement the safest algorithm for our younger users. Additionally, we could display the user's past top preferences, those that the user has already chosen, to provide more personalized recommendations. This would come in handy when users wish to re-watch past favorites.

To further improve our recommendation system, we also plan to implement a wishlist feature inspired by the "watch later" concept on YouTube. This feature will allow users to save movies they are interested in but do not have time to watch at the moment, thus providing our AI with more data to make personalized recommendations. By incorporating this feature, we hope to enhance the accuracy of our system and increase user engagement with ScreenSelect.

In our current implementation, a major limitation we encountered was the loss of our graph object upon closing the tab. To address this issue and facilitate the storage of user data in the future, we plan to develop a central database to store all information collected during each session. By storing user preferences in the format of a graph, we can easily manipulate this data and provide more personalized recommendations. Another crucial next step is optimizing the speed of our project. To achieve this, we are exploring technical options to streamline the loading and analysis process. Furthermore, we plan to expand the functionality of ScreenSelect by scraping the web for streaming platform links, allowing users to access third-party sites and watch their selected movies. Only after this would we record movies as "watched". Our current implementation determines that users have watched a movie if they simply click on it, which is not as accurate as we hope for the final version of the platform to be. By doing all of the above, we can enhance the accuracy of our recommendations and provide a more seamless user experience.

To conclude this written report, our project was successful in developing a movie recommendation app that offers tailored suggestions to users based on their watching patterns and interests. While we encountered some limitations, we were able to overcome them and offer a useful tool for movie enthusiasts. Future improvements could be made to enhance security, age-based filtering, and display past user preferences.

Part 7: References—IEEE Format

Course Notes

- D. Liu and M. Badr. "Course Notes for CSC110 and CSC111." Foundations of Computer Science. <https://www.teach.cs.toronto.edu/csc110y/fall/notes/> (accessed Mar. 5, 2023)

Data Set

- A. Dattatray, tmdb_5000_credits, vol.1, Kaggle: Kaggle, 2022. [Dataset]. Available: <https://www.kaggle.com/datasets/akshaydattatraykhare/movies-dataset>. [Accessed: Mar. 5, 2022].
- A. Dattatray, tmdb_5000_movies, vol.1, Kaggle: Kaggle, 2022. [Dataset]. Available: <https://www.kaggle.com/datasets/akshaydattatraykhare/movies-dataset>. [Accessed: Mar. 5, 2022].

Module Documentation

- Riverbank Computing and The Qt Company. "Reference Guide — PyQt Documentation v6.4.1." Riverbank Computing. <https://www.riverbankcomputing.com/static/Docs/PyQt6/> (Accessed: Mar. 5, 2022).
- J. Bodnar. "Python PyQt6." ZetCode. <https://zetcode.com/pyqt6/introduction/> (Accessed Mar. 7, 2023)
- The Qt Company. "Qt for Python" Qt. <https://doc.qt.io/qtforpython/index.html> (Accessed Mar. 8, 2023)

Idea Research

- M. Gavira. “How Netflix uses AI and Data to conquer the world.” LinkedIn.
<https://www.linkedin.com/pulse/how-netflix-uses-ai-data-conquer-world-mario-gavira/>
(Accessed Mar. 7, 2023)
- J. Ciancutti. “How We Determine Product Success.” Netflix Technology Blog.
<https://netflixtechblog.com/how-we-determine-product-success-980f81f0047e> (Accessed Mar. 7, 2023)
- D. Chong. “Deep Dive into Netflix’s Recommender System.” Towards Data Science.
<https://towardsdatascience.com/deep-dive-into-netflixs-recommender-system-341806ae3b48>
(Accessed Mar. 7, 2023)

Inline Citations (Listed in order of use)

- Netflix Help Center. (n.d.). How Netflix’s Recommendations System Works. Privacy and Security Help Page. Retrieved March 26, 2023, from <https://help.netflix.com/en/node/100639>
- Help Center. (n.d.). How to See Viewing History and Device Activity. Privacy and Security Help Page. Retrieved March 28, 2023, from <https://help.netflix.com/en/node/101917>