

ASSIGNMENT 1C

Joshua Paterson, Kaiyun Yu
N10193197, n9889663

Assignment 1C

Contribution table

Author	Completed
Joshua Paterson – N10193197	Problem 1
Kaiyun Yu - N9889663	Problem 2

Problem 1

Dataset

The data for this problem comes from the MovieLens user rating dataset. The relevant files being used are rating.csv which contain user ratings from movies and movies.csv which relate is used to reference movie names. This dataset contains 100836 ratings of 9742 movies from 610 users.

Dataset Pre-processing

As the objective of the problem is to find movie recommendations for users the data in its current form will not be able to produce this. As users are the object that will be clustered with the information used to cluster them will be their ratings of varies movies. The data will be transformed to create a table where the row of the dataset will be each individual user and the columns will represent a specific movie rating. In figure 1 below a visual representation of the data can be seen.

It should be noted that a lot of the data is unrated or NaN as most users will not have rated every movie.

It was decided to use all movies and users in the model as it seemed unlikely that every user will have watched a high number of the 9742 movies available causing a sparse dataset.

```
# Merge the two tables then pivot so we have Users X Movies dataframe
ratings_title = pd.merge(ratings, movies[['movieId', 'title']], on='movieId')
user_movie_ratings = pd.pivot_table(ratings_title, index='userId', columns= 'title', values='rating')
print('dataset dimensions: ', user_movie_ratings.shape, '\n\nsubset example:')
user_movie_ratings.iloc[:6, :10]
```

dataset dimensions: (610, 9719)

Subset example:

title	'71 (2014)	'Hellboy': The Seeds of Creation (2004)	'Round Midnight (1986)	'Salem's Lot (2004)	'Til There Was You (1997)	'Tis the Season for Love (2015)	'burbs, The (1989)	'night Mother (1986)	(500) Days of Summer (2009)	*batteries not included (1987)
userId										
1	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
2	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
3	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
5	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
6	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

Figure 1: Data after Pre-processing

Model

To cluster this dataset the k-means model/algorithm was chosen to be used. This model was chosen because due to the fact the there are so many dimensions (movies) to this data it is hard to visualise, and it is therefore hard for me to know how this data should be classified due to this the simplest option was chosen being k-means. This method was also chosen to limit the number of movies recommendation per cluster if a method like GMM was chosen then multiple movies could belong to

many clusters if there is a lot of overlap movie recommendation could simply become the most highly rated and watched movies. k-means algorithms uses hard assignment and avoids this issue.

The algorithm uses all default values except for k or number of clusters. To find an optimal number of clusters for this problem a grid search like method was used creating 14 k-means models increasing the increment of cluster by 1 then calculating the Silhouette Score to determine how well users are clustered. A Silhouette Score measure of how similar an object is to its own cluster compared to other clusters. A higher value indicates that a user is more accurately depicted in a cluster on average. As seen in figure 2 below the Silhouette Score (Distortion) falls after 4 clusters however after running this grid search method a few times the location of this fall varies around 4 to 6 and therefore a k of 5 was chosen.

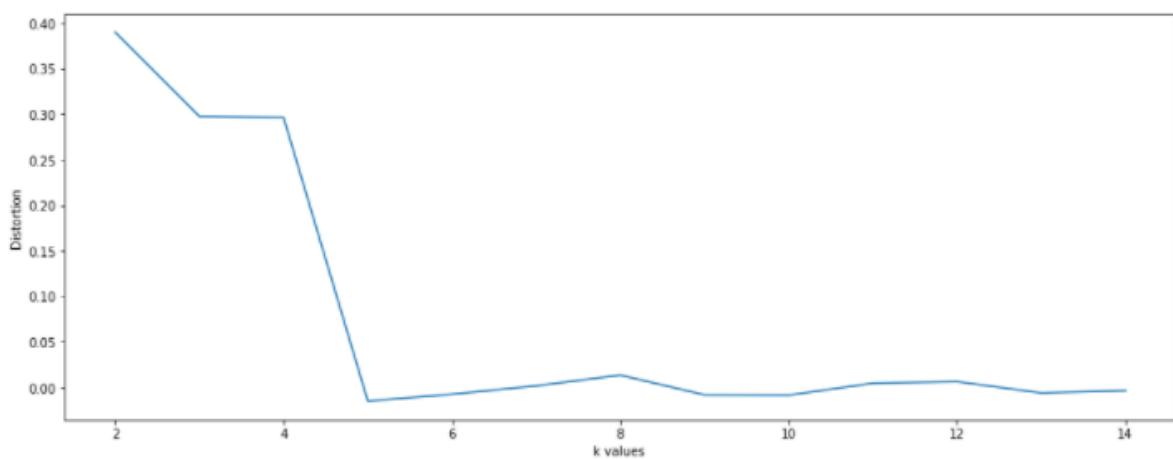


Figure 2: silhouette score of n clusters

How recommendations are formed

Using the model stated previously the model can generate predictions for users by locating the cluster of which an individual user belongs to then finding the highest rated movies that their cluster has rated that they themselves have not rated.

Results of the clustering

In figure 3 below it shows a visual of the resulting cluster that the model and dataset produce. Each of the plots in this figure shows the rating of movies (columns) in colour (5-star rating) by users (rows). Note that white means no rating given. There are clusters with both unproportionally huge and small groups showing the groups are not sorted very evenly. To determine how well this model fits the problem an analysis of individuals will be given below.

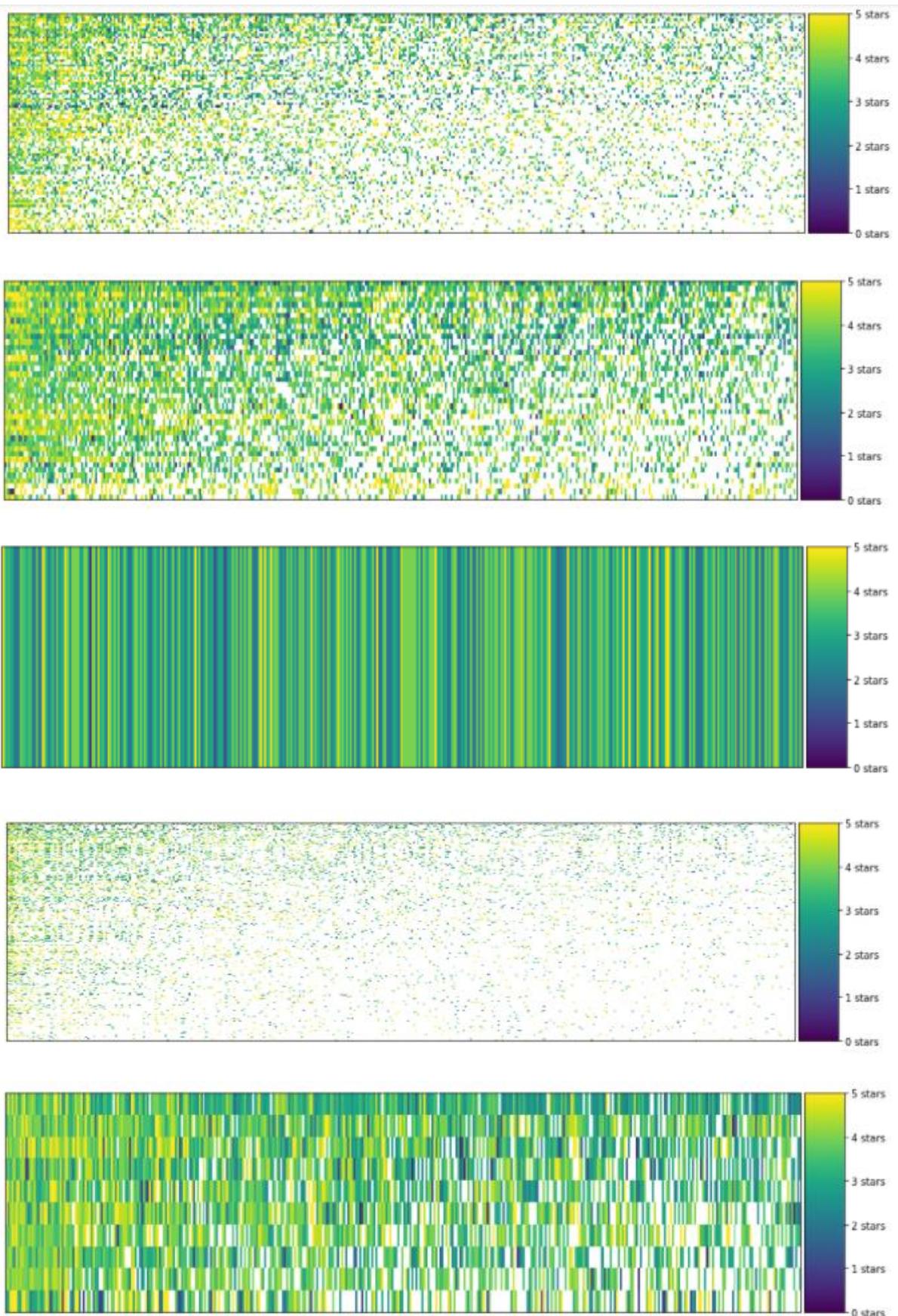


Figure 3: Resulting Clusters of Model.

Results/User Recommendations

User 4

In the below figure or figure 4 is shows users 4 most watched movies genres that was rated above 3. This user enjoys mostly drama, romance and comedy type movies.

Drama	15
Comedy	13
Comedy Drama	10
Comedy Drama Romance	6
Drama Romance	6
..	
Adventure Comedy Western	1
Children Comedy Fantasy Musical	1
Action Adventure Mystery Romance Thriller	1
Adventure Comedy	1
Adventure Fantasy Musical	1
..	

Figure 4: User 4's count of genres rated above 3.

In Figure 5 below it shows the top 5 recommendations the model has produced, and as seen Drama or Comedy is in the movie Genres. This show for User 4 the cluster can group similar movies to what user 4 has shown a liking towards.

movielid		title	genres
277	318	Shawshank Redemption, The (1994)	Crime Drama
733	953	It's a Wonderful Life (1946)	Children Drama Fantasy Romance
906	1204	Lawrence of Arabia (1962)	Adventure Drama War
929	1228	Raging Bull (1980)	Drama
1883	2502	Office Space (1999)	Comedy Crime

Figure 5: User 4's Top 5 recommendations

User 42

In the below figure or figure 6 it shows users 42 most watched movies genres that was rated above 3. This user enjoys mostly Drama and Comedy type movies and seems very similar to User 4 however it should be noted that this user was sorted into a different cluster. This may be because this user has shown a liking to other genres where User 4 seemed to be exclusively into comedy drama and Romance.

Comedy	33
Drama	22
Comedy Romance	18
Drama Romance	9
Action Adventure Thriller	9
..	
Comedy Drama War	1
Mystery Thriller	1
Crime Horror Mystery Thriller	1
Comedy Crime Drama Thriller	1
Action Adventure Mystery Thriller	1
..	

Figure 6: User 42's count of genres rated above 3.

In Figure 7 below it shows the top 5 recommendations the model has produced, and as seen the movies in the drama genre does show however unlike User 4 this cluster seems to have more action and adventure movies included in the selection.

movielid		title	genres
224	280	Star Wars: Episode IV - A New Hope (1977)	Action Adventure Sci-Fi
701	919	Wizard of Oz, The (1939)	Adventure Children Fantasy Musical
941	1242	Glory (1989)	Drama War
971	1272	Patton (1970)	Drama War
1291	1721	Titanic (1997)	Drama Romance

Figure 7: User 42's Top 5 recommendations

User 314

In the below figure or figure 4 is shows users 314 most watched movies genres that was rated above 3. This user enjoys mostly drama, comedy, and romance type movies. This preference for these types of movies genres to the point where they are the only genres that the user has rated above 5 more then once is very similar to User 4 and the model was able to recognise this as they are in the same cluster.

Comedy Romance	5
Drama	4
Comedy Drama Romance	2
Children Drama Fantasy Mystery	1
Action Adventure Thriller	1
Comedy	1
Drama Romance	1
Adventure Drama IMAX	1
Drama Horror Sci-Fi	1
Action Comedy Sci-Fi	1
Action Adventure Sci-Fi	1
Crime Drama	1
Action Thriller	1
Action Crime Drama War	1
Adventure Drama Sci-Fi	1
Comedy Drama Fantasy Romance Thriller	1
Crime Mystery Thriller	1
Action Drama War	1
Comedy Drama Romance War	1
Adventure Animation Children Drama Musical IMAX	1
Drama War	1
Action Sci-Fi	1
Comedy Fantasy Romance	1

Figure 8: User 314's count of genres rated above 3.

In Figure 9 below it shows the top 5 recommendations the model has produced, and as seen movies in the drama and romance genres are frequently suggested for this user showing a reasonable recommendation list however this is one suggestion in this list where the User has shown no interest in its genres. This was likely because movies in clusters are sorted by most highly rated and will therefore have a bias for movies with high average rating in their cluster.

movielid		title	genres
277	318	Shawshank Redemption, The (1994)	Crime Drama
686	904	Rear Window (1954)	Mystery Thriller
690	908	North by Northwest (1959)	Action Adventure Mystery Romance Thriller
975	1278	Cool Hand Luke (1967)	Drama
1917	2542	Lock, Stock & Two Smoking Barrels (1998)	Comedy Crime Thriller

Figure 9: User 314's Top 5 recommendations

Overall thoughts

Overall, I think the recommendation were reasonable as the model was able to recognise Users 4 and 314 with very similar tastes as they both clearly did not like much outside of 3 genres, they would only highly rate. User 34 also like those same genres but showed some interest in other genres and was able to be sorted into a group with a more diverse movie selection. As the model was able to easily categorized similar users the rating of each group would allow the model to find recommendation of user with similar preference. However there where some issues as there was one group with one user assigned to it. This is an issue as this user would not be able to find recommendation due to this issue. To fix this in the future the number of k would need to be selected to find a number that would have more then 1 amount of people per cluster.

Problem 2

Description of the algorithm

We eliminated the bad data and normalized the data. In the colour classification problem, we also performed data enhancement operations such as image flipping.

We use ResNet34 as our method and added a classifier to the last layer to get the results we want. The reason for using this network is that it has been proven to be very effective in many scenarios on image classification problems. In order to speed up the training speed, here we use pytorch's pre-trained initialization parameters when training the network, the loss function uses the crossentropy loss function, the optimization algorithm uses the Adam algorithm, and the learning rate is set to automatically decrease.

Our algorithm basically achieves an accuracy of 80%–90% on these test sets within 50 to 100 iteration steps. Due to the small amount of data, this effect is still very successful. If you want better results, you need more data. Regarding semantic search, we can use CNN to first extract the features we want from the sentence, and then use our image classification algorithm for matching search.

Data has been resized to 256,128 sizes and random horizontal flip for increasing the data augmentation.

Results

In this section, we list the results of classifying different traits on the testing data set. The number in the upper right corner of each picture represents the recognition result of our algorithm for each picture.

Gender

There are totally 20 Epoch has been used in training. The best train accuracy is 0.998066 and best valid accuracy is 0.664894. Training was completed in 8 mins.



Figure 1: Results of gender classification.



Figure 2: Results of gender classification.



Figure 3: Results of gender classification.



Figure 4: Results of gender classification.



Figure 5: Results of gender classification.

Torso Clothing Type

There are totally 20 Epoch has been used in training. The best train accuracy is 1.00 and best valid accuracy is 0.544231. Training was completed in 8m 22s.



Figure 6: Results of torso clothing type classification.



Figure 7: Results of torso clothing type classification.



Figure 8: Results of torso clothing type classification.



Figure 9: Results of torso clothing type classification.



Figure 10: Results of torso clothing type classification.

Torso Clothing Colour

There are totally 20 Epoch has been used in training. The best train accuracy is 0.829457 and best valid accuracy is 0.474490. Training was completed in 8m 22s.



Figure 11: Results of torso clothing colour classification.



Figure 12: Results of torso clothing colour classification.



Figure 13: Results of torso clothing colour classification.



Figure 14: Results of torso clothing colour classification.



Figure 15: Results of torso clothing colour classification.

Leg Clothing Type

There are totally 20 Epoch has been used in training. The best train accuracy is 1.00 and best valid accuracy is 0.84. Training was completed in 15m 36s.



Figure 16: Results of leg clothing type classification.



Figure 17: Results of leg clothing type classification.



Figure 18: Results of leg clothing type classification.



Figure 19: Results of leg clothing type classification.



Figure 20: Results of leg clothing type classification.

Leg Clothing Colour

There are totally 20 Epoch has been used in training. The best train accuracy is 0.998066 and best valid accuracy is 0.664894. Training completed in 5m 36s.



Figure 21: Results of leg clothing colour classification.

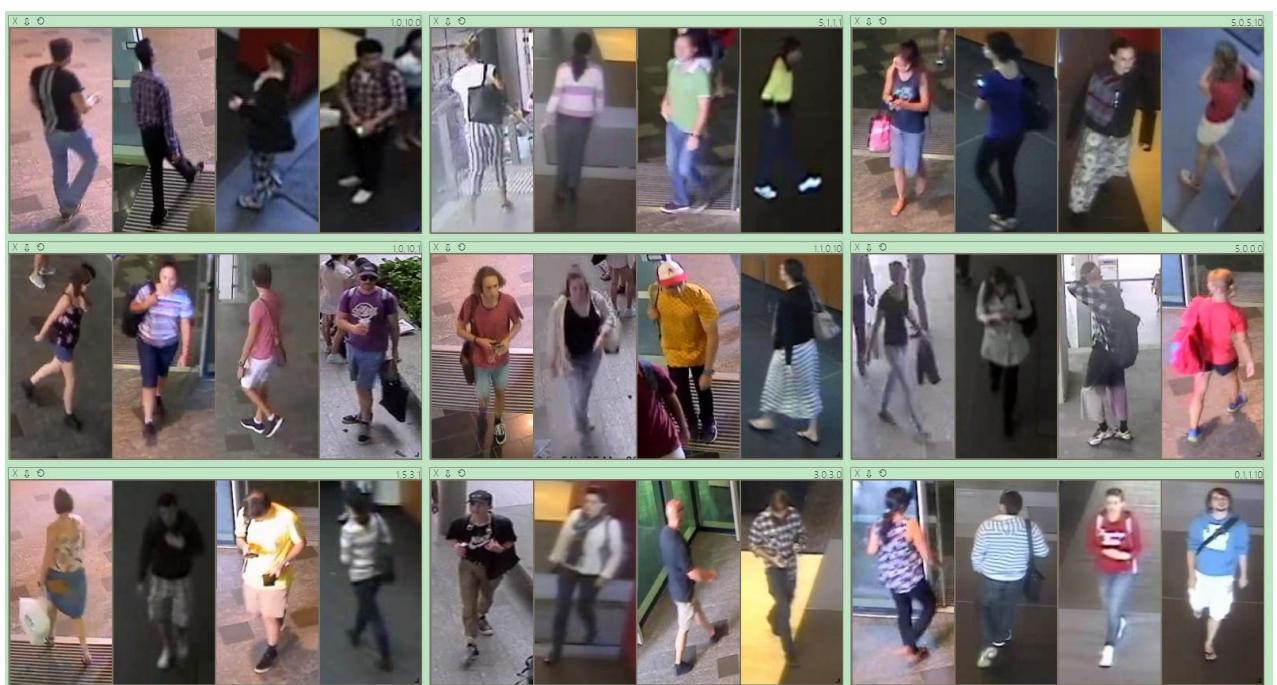


Figure 22: Results of leg clothing colour classification.



Figure 23: Results of leg clothing colour classification.



Figure 24: Results of leg clothing colour classification.



Figure 25: Results of leg clothing colour classification.

Luggage

There are totally 20 Epoch has been used in training. The best train accuracy is 0.863462 and best valid accuracy is 0.647059. Training was completed in 8m 5s.



Figure 26: Results of luggage classification.



Figure 27: Results of luggage classification.



Figure 28: Results of luggage classification.



Figure 29: Results of luggage classification.



Figure 30: Results of luggage classification.

Problems and solution

In this model, training accuracy has achieved significant accuracy which achieved 1.0 accuracy. The major reason which causes the problem is in some data set there are too many classifications but with small amount data, random horizontal flip has been applied to the model for data augmentation. Moreover, the model has overfitting in some points, thus learning rate decay scheduler is using for reducing the overfitting problem. The outcome is significant and achieved satisfied result within small amount of data.

Code

