

RECOMMENDER SYSTEM USING K-NN ALGORITHM

A SEMINAR REPORT

Submitted by

**ARPIT MANOCHA [Reg No: RA1811003010264]
RISHABH PRAKASH [Reg No: RA1811003010272]**

Under the guidance of

Mrs. Poornima S.E., Ph.D

(Professor, Department of Computer Science & Engineering)

in partial fulfillment for the award of the degree

of

BACHELOR OF TECHNOLOGY

in

COMPUTER SCIENCE

of

FACULTY OF ENGINEERING AND TECHNOLOGY



S.R.M. Nagar, Kattankulathur, Kancheepuram District

November 2020

SRM UNIVERSITY

(Under Section 3 of UGC Act, 1956)

BONAFIDE CERTIFICATE

Certified that this Industrial training report titled “**RECOMMENDER SYSTEM USING K-NN ALGORITHM**” is the bonafide work of “**ARPIT MANOCHA [Reg No: RA1811003010264], RISHABH PRAKASH [Reg No: RA1811003010272]**”, submitted for the course 18CSP103L Seminar-I. This report is a record of successful completion of the specified course evaluated based on literature reviews and the supervisor. No part of the Seminar Report has been submitted for any degree, diploma, title or recognition before.

SIGNATURE

Mrs. Poornima S.E., Ph.D
GUIDE
Professor
Dept. of Computer Science & Engineering

Signature of the Internal Examiner

SIGNATURE

Dr. B. Amutha
HEAD OF THE DEPARTMENT
Dept. of Computer Science

Signature of the External Examiner

ABSTRACT

Recommendation system can be defined as a system that produces individual recommendations (a personalized way of possible options) as an output based on their previous choices which are considered an input by the system. Most of the products that we use today are powered by recommendation system. A movie recommendation is important in our social life due to its strength in providing enhanced entertainment. Such a system can suggest a set of movies to users based on their interest, or the popularities of the movies. Although, a set of movie recommendation systems have been proposed, most of these either cannot recommend a movie to the existing users efficiently or to a new user by any means. In this project we propose a recommendation system using the K-NN algorithm that has the ability to recommend movies to a new user as well as the others. It mines movie databases to collect all the important information, such as, genres, keywords, required for recommendation. Also we have analysed how much weight should be given to which feature because it will increase the efficiency of the K-NN algorithm to give the most suitable result.

ACKNOWLEDGEMENTS

We would like to express our deepest gratitude to our guide, Mrs.Poornima S.E. her valuable guidance, consistent encouragement, personal caring, timely help and providing us with an excellent atmosphere for doing research. All through the work, in spite of her busy schedule, she has extended cheerful and cordial support to us for completing this research work.

Author

TABLE OF CONTENTS

ABSTRACT	iii
ACKNOWLEDGEMENTS	iv
LIST OF FIGURES	vii
1 INTRODUCTION	viii
1.1 Commercial Uses Of Recommendation Systems	viii
1.2 Why the Recommendation system?	ix
1.3 Our Objective	ix
2 Recommendation Systems	x
2.1 Content based filtering	x
2.1.1 Merits	xi
2.1.2 Demerits	xi
2.2 Collaborative filtering	xi
2.2.1 Memory based techniques	xii
2.2.2 Model-based techniques	xiii
2.2.3 Pros and Cons of collaborative filtering techniques	xiii
2.3 Hybrid filtering	xiv
3 Literature Review	xvi
3.1 K-Nearest Neighbors Algorithm	xvi
3.1.1 Definition	xvi
3.1.2 Mathematics related to the Algorithm	xvii
4 System Analysis	xx
4.1 Dataset	xx
4.2 Approach	xxi

5	Feature Selection	xxii
6	Coding	xxvii
6.1	Weight Analysis Python Code	xxvii
6.2	Recommender System Python Code	xxviii
7	Inferences/Conclusion	xxix
8	Future Enhancements	xxx
9	References	xxxi

LIST OF FIGURES

5.1	Impact of feature exclusion on Average Rank.	xxiii
5.2	Impact of feature exclusion on top 5 rate.	xxiii
5.3	Movie recommendations on exclusion of "Keyword" and "Cast". . .	xxiv
5.4	Movie recommendations on exclusion of "Keyword" and "Cast". . .	xxv
5.5	Average Rank for different weight scenarios.	xxv
5.6	Top 5 rate for different weight scenarios.	xxvi
6.1	Code of the Weight Analysis	xxvii

CHAPTER 1

INTRODUCTION

The explosive growth in the amount of available digital information and the number of visitors to the Internet have created a potential challenge of information overload which hinders timely access to items of interest on the Internet. Information Retrieval systems, such as Google, Devil Finder and Altavista have partially solved this problem but prioritization and personalization (where a system maps available content to user's interests and preferences) of information were absent. This has increased the demand for recommender systems more than ever before. Recommender systems are information filtering systems that deal with the problem of information overload by filtering vital information fragment out of large amount of dynamically generated information according to user's preferences, interest, or observed behavior about item. Recommender system has the ability to predict whether a particular user would prefer an item or not based on the user's profile.

1.1 Commercial Uses Of Recommendation Systems

Recommender systems are beneficial to both service providers and users. They reduce transaction costs of finding and selecting items in an online shopping environment. Recommendation systems have also proved to improve decision making process and quality. In e-commerce setting, recommender systems enhance revenue for the fact that they are effective means of selling more products. Recommender system can also be used to recommend different types of movies of different genres. Various entertainment websites such as Netflix, Amazon Prime, Hotstar, etc. uses recommendation system to predict the content which will be preferred by the user. It enhances the user's experience too.

1.2 Why the Recommendation system?

- Help item providers in delivering their items to the right user.
- Identity products that are most relevant to users.
- Personalized content.
- Help websites to improve user engagement.
- Benefits users in finding items of their interest.

1.3 Our Objective

Our aim with this project is to analyse the k-Nearest Neighbor Algorithm and develop a recommendation system that uses the key features with inclusion of their weights according to the importance of different features that can be used for recommendation.

In this paper, the key research contents is to help users to obtain user-interested movie automatically in the massive movie information data using KNN algorithm and content based filtering algorithm, and to develop a prototype of movie recommendation system based on KNN content based filtering algorithm.

In this model we have created a recommender engine which makes use of K-NN algorithm to find out the top 5 similar movies to the movie that the user has give as input. We have discussed in detail about the K-NN algorithm in this report.

CHAPTER 2

RECOMMENDATION SYSTEMS

Recommender system is defined as a decision making strategy for users under complex information environments. Also, recommender system was defined from the perspective of E-commerce as a tool that helps users search through records of knowledge which is related to users' interest and preference. Recommender system was defined as a means of assisting and augmenting the social process of using recommendations of others to make choices when there is no sufficient personal knowledge or experience of the alternatives . Recommender systems handle the problem of information overload that users normally encounter by providing them with personalized, exclusive content and service recommendations.

2.1 Content based filtering

Content-based technique is a domain-dependent algorithm and it emphasizes more on the analysis of the attributes of items in order to generate predictions. When documents such as web pages, publications and news are to be recommended, content-based filtering technique is the most successful. In content-based filtering technique, recommendation is made based on the user profiles using features extracted from the content of the items the user has evaluated in the past. Items that are mostly related to the positively rated items are recommended to the user.

CBF uses different types of models to find similarity between documents in order to generate meaningful recommendations. It could use Vector Space Model such as Term Frequency Inverse Document Frequency (TF/IDF) or Probabilistic models such as Naïve Bayes Classifier, Decision Trees or Neural Networks to model the relationship between different documents within a corpus. These techniques make recommendations by learning the underlying model with either statistical analysis or machine learning techniques.

Content-based filtering technique does not need the profile of other users since they do not influence recommendation. Also, if the user profile changes, CBF technique still has the potential to adjust its recommendations within a very short period of time. The major disadvantage of this technique is the need to have an in-depth knowledge and description of the features of the items in the profile. The information source that content-based filtering systems are mostly used with are text documents. A standard approach for term parsing selects single words from documents.

2.1.1 Merits

- There is no requirement of much of the user's data.
- We just need item data that enable us to start giving recommendations to users.
- A content-based recommender engine does not depend on the user's data, so even if a new user comes in, we can recommend the user as long as we have the user data to build his profile.
- It does not suffer from a cold start

2.1.2 Demerits

- Items data should be in good volume.
- Features should be available to compute the similarity.

2.2 Collaborative filtering

Collaborative filtering is a domain-independent prediction technique for content that cannot easily and adequately be described by metadata such as movies and music. Collaborative filtering technique works by building a database (user-item matrix) of preferences for items by users. It then matches users with relevant interest and preferences by calculating similarities between their profiles to make recommendations. Such users build a group called neighbourhood. An user gets recommendations to those items that

he has not rated before but that were already positively rated by users in his neighbourhood. Recommendations that are produced by CF can be of either prediction or recommendation. Prediction is a numerical value, R_{ij} , expressing the predicted score of item j for the user i , while Recommendation is a list of top N items that the user will like the most. The technique of collaborative filtering can be divided into two categories: memory-based and model-based.

2.2.1 Memory based techniques

The items that were already rated by the user before play a relevant role in searching for a neighbor that shares appreciation with him. Once a neighbour of a user is found, different algorithms can be used to combine the preferences of neighbours to generate recommendations. Due to the effectiveness of these techniques, they have achieved widespread success in real life applications. Memory-based CF can be achieved in two ways through user-based and item-based techniques. User based collaborative filtering technique calculates similarity between users by comparing their ratings on the same item, and it then computes the predicted rating for an item by the active user as a weighted average of the ratings of the item by users similar to the active user where weights are the similarities of these users with the target item. Item-based filtering techniques compute predictions using the similarity between items and not the similarity between users. It builds a model of item similarities by retrieving all items rated by an active user from the user-item matrix, it determines how similar the retrieved items are to the target item, then it selects the k most similar items and their corresponding similarities are also determined. Prediction is made by taking a weighted average of the active users rating on the similar items k . Several types of similarity measures are used to compute similarity between item/user. Similarity measure is also referred to as similarity metric, and they are methods used to calculate the scores that express how similar users or items are to each other. These scores can then be used as the foundation of user- or item-based recommendation generation. Depending on the context of use, similarity metrics can also be referred to as correlation metrics or distance metrics.

2.2.2 Model-based techniques

This technique employs the previous ratings to learn a model in order to improve the performance of Collaborative filtering Technique. The model building process can be done using machine learning or data mining techniques. These techniques can quickly recommend a set of items for the fact that they use pre-computed model and they have proved to produce recommendation results that are similar to neighbourhood-based recommender techniques. Examples of these techniques include Dimensionality Reduction technique such as Singular Value Decomposition (SVD), Matrix Completion Technique, Latent Semantic methods, and Regression and Clustering. Model-based techniques analyze the user-item matrix to identify relations between items; they use these relations to compare the list of top-N recommendations. Model based techniques resolve the sparsity problems associated with recommendation systems.

2.2.3 Pros and Cons of collaborative filtering techniques

Pros

Collaborative Filtering has some major advantages over CBF in that it can perform in domains where there is not much content associated with items and where content is difficult for a computer system to analyze (such as opinions and ideal).

CF technique has the ability to provide serendipitous recommendations, which means that it can recommend items that are relevant to the user even without the content being in the user's profile.

Cons

- Cold start: For a new user or item, there isn't enough data to make accurate recommendations.
- Scalability: In many of the environments in which these systems make recommendations, there are millions of users and products. Thus, a large amount of computation power is often necessary to calculate recommendations.

- **Sparsity:** The number of items sold on major e-commerce sites is extremely large. The most active users will only have rated a small subset of the overall database. Thus, even the most popular items have very few ratings.

2.3 Hybrid filtering

Most recommender systems now use a hybrid approach, combining collaborative filtering, content-based filtering, and other approaches. There is no reason why several different techniques of the same type could not be hybridized. Hybrid approaches can be implemented in several ways: by making content-based and collaborative-based predictions separately and then combining them; by adding content-based capabilities to a collaborative-based approach (and vice versa); or by unifying the approaches into one model (for a complete review of recommender systems). Several studies that empirically compare the performance of the hybrid with the pure collaborative and content-based methods and demonstrated that the hybrid methods can provide more accurate recommendations than pure approaches. These methods can also be used to overcome some of the common problems in recommender systems such as cold start and the sparsity problem, as well as the knowledge engineering bottleneck in knowledge-based approaches.

Netflix is a good example of the use of hybrid recommender systems. The website makes recommendations by comparing the watching and searching habits of similar users (i.e., collaborative filtering) as well as by offering movies that share characteristics with films that a user has rated highly (content-based filtering).

Some hybridization techniques include:

- **Weighted:** Combining the score of different recommendation components numerically.
- **Switching:** Choosing among recommendation components and applying the selected one.
- **Mixed:** Recommendations from different recommenders are presented together to give the recommendation.
- **Feature Combination:** Features derived from different knowledge sources are com-

bined together and given to a single recommendation algorithm.

- Feature Augmentation: Computing a feature or set of features, this is then part of the input to the next technique.
- Cascade: Recommenders are given strict priority, with the lower priority ones breaking ties in the scoring of the higher ones.
- Meta-level: One recommendation technique is applied and produces some sort of model, which is then the input used by the next technique.

CHAPTER 3

LITERATURE REVIEW

Rishabh Ahuja, Arun Solanki, Anand Nayyar [1] showed in their paper that k-NN algorithm is capable of predicting movies on the basis of the various features selection about which we have given details in the further chapter, how we have implemented in our project.

Content based [2], [3] collaborative [4] and hybrid [5] are the different approaches used by past researcher for the development of recommender system. In 2007 a web-based movie recommendation system using hybrid filtering methods is presented by the authors [6]. In 2013 a Bayesian network and Trust model based movie recommendation system is proposed, the Bayesian network is imported for user preference modeling and trust model is used to filter the recommending history data and enable the system to tolerant the noisy data [7]

3.1 K-Nearest Neighbors Algorithm

3.1.1 Definition

- KNN is a non-parametric method used for classification and regression. In pattern recognition, the k-nearest neighbors algorithm (k-NN) is a non-parametric method proposed by Thomas Cover used for classification and regression. In both cases, the input consists of the k closest training examples in the feature space. The output depends on whether k-NN is used for classification or regression
- k-NN is a type of instance-based learning, or lazy learning, where the function is only approximated locally and all computation is deferred until function evaluation. Since this algorithm relies on distance for classification, normalizing the training data can improve its accuracy dramatically.
- Both for classification and regression, a useful technique can be to assign weights to the contributions of the neighbors, so that the nearer neighbors contribute more to the average than the more distant ones. For example, a common weighting scheme consists in giving each neighbor a weight of $1/d$, where d is the distance to the neighbor.

- The neighbors are taken from a set of objects for which the class (for k-NN classification) or the object property value (for k-NN regression) is known. This can be thought of as the training set for the algorithm, though no explicit training step is required.

3.1.2 Mathematics related to the Algorithm

- Assume a dataset $[(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)]$ where $X_i \in R^d$ and $Y_i \in R$. We define X_i as the feature set and Y_i , as the class label.
- Given a new data point, $x \in R^d$. Let,

$$d(X_1, x) < d(X_2, x) < d(X_3, x) < d(X_4, x) \dots < d(X_n, x) \quad (3.1)$$

$$\text{where, } d(A, B) = 1 - \frac{A \cdot B}{\|A\| \|B\|} \quad (3.2)$$

we can produce a label y for x such that,

$$y = \text{mode}(Y_1, Y_2, Y_3, \dots, Y_n) \quad (3.3)$$

- In our problem, we don't have class labels (Y) and it is a rank ordering problem where we are ordering our movies, books or any other stuff by comparing the features of the input data point to all the movies in the dataset.
- We define $d(a, b)$ as the distance between the two coordinates in the plot of the features that we are going to use. It can be calculated using any suitable distance measure. The three approaches that we studied are:
 - Cosine Distance
 - Jaccard Similarity
 - Euclidean Distance
- We have chosen cosine distance in our implementation since this is a popularly used distance measure which can be used to compare text documents by looking for common words between the documents. In our implementation, the features are textual in nature and thus, we assessed this measure to be the most appropriate.

Euclidean Distance

The Euclidean distance between two points, the length of a line segment between the two points. It can be calculated from the Cartesian coordinates of the points using the equation:

$$d(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_i - q_i)^2 + \dots + (p_n - q_n)^2} \quad (3.4)$$

Cosine Distance

We can calculate the similarity between two vectors using the cosine of the angle between the two vectors, vectors of A and B. If the vectors are closer, then small will be the angle and large will be the cosine. Since the cosine similarity ranges $[0, 1]$, we can calculate the distance using $1 - \text{Cosine Similarity}$. Cosine Distance can be computed using:

$$d(A, B) = 1 - \frac{A \cdot B}{\|A\| \|B\|} \quad (3.5)$$

Jaccard Similarity

The Jaccard Similarity measures similarity between finite sample sets, and is defined as the size of the intersection divided by the size of the union of the sample sets. We can calculate Jaccard Similarity as:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \quad (3.6)$$

Pros and Cons of the Algorithm

Pros

- K-NN is pretty intuitive and simple: K-NN algorithm is very simple to understand and equally easy to implement. To classify the new data point K-NN algorithm reads through whole dataset to find out K nearest neighbors.
- K-NN has no assumptions: K-NN is a non-parametric algorithm which means there are no assumptions to be met to implement K-NN. Parametric models like linear regression has lots of assumptions to be met by data before it can be implemented which is not the case with K-NN.
- No Training Step: K-NN does not explicitly build any model, it simply tags the new data entry based learning from historical data. New data entry would be tagged with majority class in the nearest neighbor.
- It constantly evolves: Given it's an instance-based learning; k-NN is a memory-based approach. The classifier immediately adapts as we collect new training data. It allows the algorithm to respond quickly to changes in the input during real-time use.

- Very easy to implement for multi-class problem: Most of the classifier algorithms are easy to implement for binary problems and needs effort to implement for multi class whereas K-NN adjust to multi class without any extra efforts.
- Can be used both for Classification and Regression: One of the biggest advantages of K-NN is that K-NN can be used both for classification and regression problems.
- One Hyper Parameter: K-NN might take some time while selecting the first hyper parameter but after that rest of the parameters are aligned to it.

Cons

- K-NN slow algorithm: K-NN might be very easy to implement but as dataset grows efficiency speed of algorithm declines very fast.
- Curse of Dimensionality: KNN works well with small number of input variables but as the numbers of variables grow K-NN algorithm struggles to predict the output of new data point.
- K-NN needs homogeneous features: If you decide to build k-NN using a common distance, like Euclidean or Manhattan distances, it is completely necessary that features have the same scale, since absolute differences in features weight the same, i.e., a given distance in feature 1 must means the same for feature 2.
- Optimal number of neighbors: One of the biggest issues with K-NN is to choose the optimal number of neighbors to be consider while classifying the new data entry.
- Imbalanced data causes problems: k-NN doesn't perform well on imbalanced data. If we consider two classes, A and B, and the majority of the training data is labeled as A, then the model will ultimately give a lot of preference to A.
- This might result in getting the less common class B wrongly classified. Outlier sensitivity: K-NN algorithm is very sensitive to outliers as it simply chose the neighbors based on distance criteria.
- Missing Value treatment: K-NN inherently has no capability of dealing with missing value problem.

CHAPTER 4

SYSTEM ANALYSIS

4.1 Dataset

- Source of our dataset is: <https://github.com/codeheroku/Introduction-to-Machine-Learning/blob/master/Building>
- This dataset contains the following features of a movie:
 - index
 - budget
 - genre
 - homepage
 - id
 - keywords for the movie
 - original language
 - original title
 - overview
 - popularity
 - production companies
 - production countries
 - revenue
 - movie length
 - spoken language
 - tagline
 - vote average
 - vote count
 - cast
 - crew
 - director
- The movie dataset contains a dataset of 5808 movies.
- The dataset is cleaned too before using it for the project to avoid error due to dataset.

4.2 Approach

- For our movie prediction we have used 4 features which are Genres, Director ,cast and keywords.
- We have also analysed that which feature have how much effect on the movie prediction by the recommendation engine, by changing the weightage of each of the features one by one and taking various cases each giving different weightage to different features.
- The algorithm we have used for this recommendation problem i.e. a rank ordering problem, is KNN Algorithm in which the similarity between points is calculated using cosine distance.
- Our recommendation model finds out the nearest 10 movies to the movie entered as input by the user using coordinates of different movies(calculated using their features) and print them in order.

CHAPTER 5

FEATURE SELECTION

- We have selected 4 features from the dataset of the movie that are Genres, Director, Cast and keywords based on which we are trying to predict the similar movies to the movie the user has entered.
- We have tested the importance of the selected features for the KNN algorithm.
- To test the most suitable weightage that should be given to different features we did some observations and analysis on some movie sample.
- In our observation we selected a sample movie from the dataset and created a list of 3 movies recommendations based on subjective oversight and google search results.
- For Example, we took a sample movie named: Man Of Steel which have the following different features:
 - Title: Man Of Steel
 - Genre: Action Adventure Fantasy Science Fiction
 - Director Name: Zack Snyder
 - Cast: Henry Cavill Amy Adams Michael Shannon Kevin Costner Diane Lane
 - Keywords: saving the world dc comics superhero based on comic book superhuman
- Expected Recommendations: Superman, Batman Vs Superman, Watchmen. These movies were selected as expected due to their similarity with the movie, Man Of Steel and also they were recommended by the google recommendation Engine.
- We expect the above movies to be included in the top 10 recommendations by the model, especially among the top 5. So, we have defined the following error metrics to test the model performance,

$$Average Rank = \frac{Rank(Movie_1) + Rank(Movie_2) + Rank(Movie_3)}{3} \quad (5.1)$$

$$Top 5 Rate = \frac{Number\ of\ expected\ movies\ in\ top\ 5 * 100}{3} \quad (5.2)$$

We define a default rank of 100 if a movie isn't selected in top 10 to penalize the recommendation.

- We excluded each feature and analyzed the impact on Average Rank and Top 5 rate. It was observed that both the metrics declined on exclusion of these features except for the feature, "Cast". The results can be observed in figures 5.1 and 5.2.

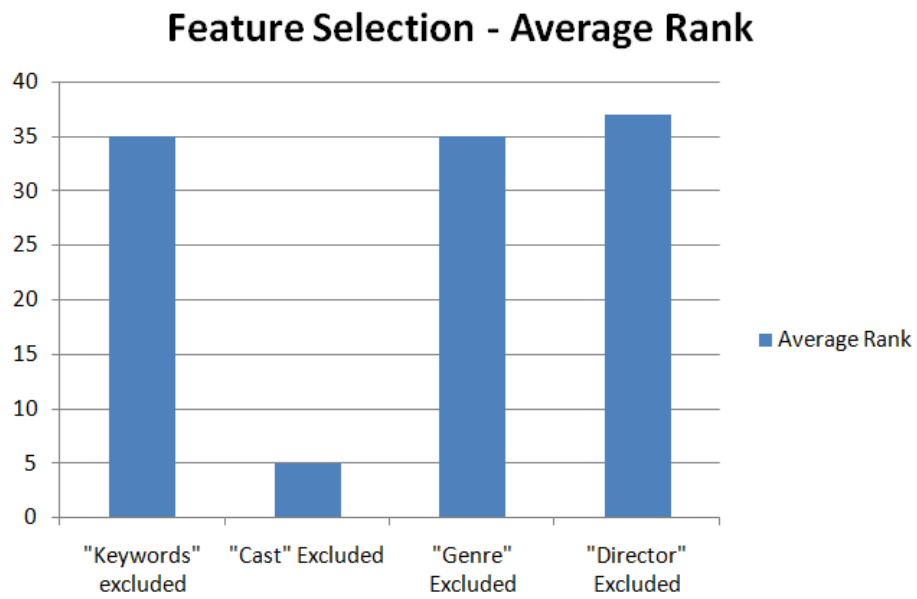


Figure 5.1: Impact of feature exclusion on Average Rank.

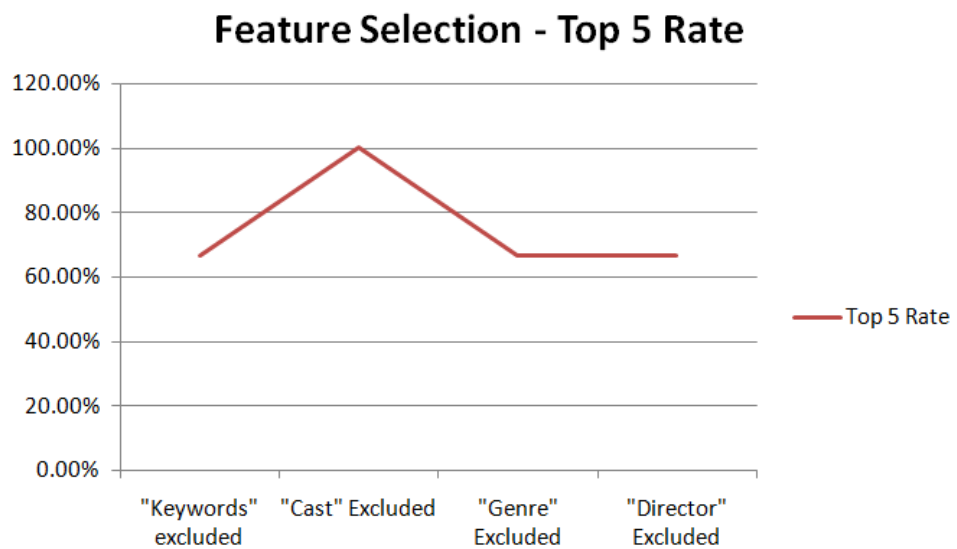


Figure 5.2: Impact of feature exclusion on top 5 rate.

```

Enter your favourite movie name: > Man of steel
Top 5 movies that you might love to watch:
Man of Steel
Batman v Superman: Dawn of Justice
Watchmen
Sucker Punch
300
Dawn of the Dead
=====
Top 5 movies that you might love to watch:
Man of Steel
Batman v Superman: Dawn of Justice
Watchmen
Superman II
Superman Returns
X-Men
=====

```

Figure 5.3: Movie recommendations on exclusion of "Keyword" and "Cast".

- The model recommendations on exclusion of these features are also provided for reference. Assessment of model recommendations are provided below:
 - Exclusion of "Keyword": The unique recommendations for this scenario, "300", "Sucker Punch" and "Dawn of the Dead" were primarily similar to the test movie in terms of the features, "Director" and "Genre". Due to the exclusion, these features had higher weights than normal. Please refer to figure 6.2 for the list of recommendations.
 - Exclusion of "Cast": The recommendations on this exclusion are impacted equally by "Director", "Keyword" and "Genre". Please refer to figure 6.2 for the list of recommendations.
 - Exclusion of "Genre": This scenario lead to the inclusion of the movie, "Legend of the Guardians: The Owls of Ga'Hoole", which we have assessed to be noisy. Even though it's a fantasy film, it's different from the sample movie. Please refer to figure 5.4 for the list of recommendations.
 - Exclusion of "Director": Recommended movies are impacted heavily by genre and keywords. Please refer to figure 5.4 for the list of recommendations.
- We created different scenarios by weighing each feature differently. The results from the analysis can be observed in 5.5 and 5.6. We have observed that the features "Keyword" and "Director" have higher than usual impact on the error metrics.
- Focusing on the conclusion of analysis we found out the features namely: "keyword" and "director" have higher than usual impact on the K-nearest neighbor algorithm and therefore they should be assigned higher weights. The impact of exclusion of the feature "Cast" didn't seem to have a major impact but we have chosen to keep this feature and assign a lower weight to it.


```

=====
Top 5 movies that you might love to watch:
Man of Steel
Batman v Superman: Dawn of Justice
Watchmen
Legend of the Guardians: The Owls of Ga'Hooie
Sucker Punch
300
=====
Top 5 movies that you might love to watch:
Man of Steel
Superman II
Batman v Superman: Dawn of Justice
Superman Returns
Ant-Man
X-Men: Days of Future Past

```

Figure 5.4: Movie recommendations on exclusion of "Keyword" and "Cast".

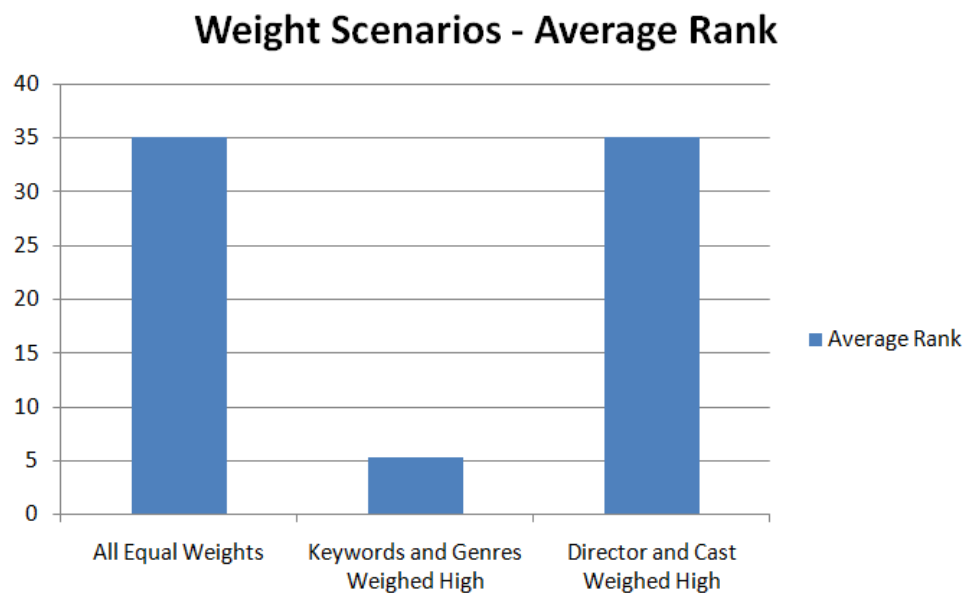


Figure 5.5: Average Rank for different weight scenarios.

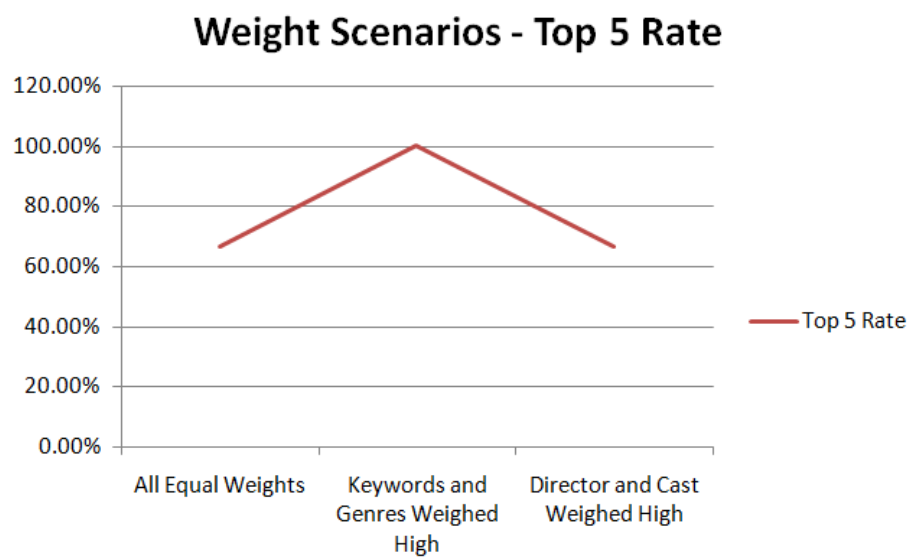


Figure 5.6: Top 5 rate for different weight scenarios.

CHAPTER 6

CODING

6.1 Weight Analysis Python Code

```
default_rank = 100
ranks=[]
weights = [[0, 0.33, 0.33, 0.33], [0.33, 0, 0.33, 0.33], [0.33, 0.33, 0, 0.33], [0.33, 0.33, 0.33, 0],
           [0.25, 0.25, 0.25, 0.25], [0.4, 0.1, 0.4, 0.1], [0.1, 0.4, 0.1, 0.4]]
for weight in weights:
    similarity_scores = []
    for i in range(df.shape[0]):
        candidate = df.iloc[i]
        similarity_score = calculate_similarity_score(movie_input, candidate, np.array(weight))
        similarity_scores.append([i, similarity_score])

    # Step 7: Get a list of similar movies in descending order of similarity score
    sorted_similar_movies = sorted(similarity_scores, key=lambda x:x[1], reverse=True)

    # Step 8: Print titles of first 4 movies
    j = 0
    print("Top 5 movies that you might love to watch: ")
    for element in sorted_similar_movies:
        print(get_title_from_index(element[0]))
        j += 1
        if j > 5:
            break
    print("=====")

    count = 0
    average_rank = 0
    for m in range(10):
        movie = get_title_from_index(sorted_similar_movies[m][0])
        if movie in ["Superman", "Batman v Superman: Dawn of Justice",
                    "Watchmen"]:
            count += 1
            average_rank += m+1
    if count == 2:
        average_rank += default_rank
    elif count == 1:
        average_rank += 2*default_rank
    elif count == 0:
        average_rank += 3*default_rank
    average_rank /= 3.0
    ranks.append([weight, average_rank])
```

Figure 6.1: Code of the Weight Analysis

6.2 Recommender System Python Code

```
import pandas as pd
import numpy as np
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics.pairwise import cosine_similarity
##### helper functions. Use them when needed #####
def get_title_from_index(index):
    return df[df.index == index]["original_title"].values[0]

def get_index_from_title(title):
    return int(df[df.title == title]["index"].values[0])
#####

##Step 1: Read CSV file
df = pd.read_csv("movie_dataset.csv")
movies = df["title"]
movies = movies.fillna('', inplace=False)
for i in range(movies.shape[0]):
    movies.iloc[i] = movies.iloc[i].lower()

df['title'] = movies
#print df.columns
##Step 2: Select Features

features = ['keywords', 'cast', 'genres', 'director']
##Step 3: Create a column in DF which combines all selected features
for feature in features:
    df[feature] = df[feature].fillna('')

def combine_features(row):
    try:
        return row['keywords'] + " " + row['cast'] + " " + row['genres'] + " " + row['director'] + " " + row['original_language']
    except:
        print("Error:", row)

df["combined_features"] = df.apply(combine_features, axis=1)

#print "Combined Features:", df["combined_features"].head()

##Step 4: Create count matrix from this new combined column
cv = CountVectorizer()

count_matrix = cv.fit_transform(df["combined_features"])

##Step 5: Compute the Cosine Similarity based on the count_matrix
cosine_sim = cosine_similarity(count_matrix)

inp = input("Enter your favourite movie name: ")
inp = inp.lower()
movie user likes= inp

## Step 6: Get index of this movie from its title
movie_index = get_index_from_title(movie_user_likes)

similar_movies = list(enumerate(cosine_sim[movie_index]))

## Step 7: Get a list of similar movies in descending order of similarity score
sorted_similar_movies = sorted(similar_movies, key=lambda x: x[1], reverse=True)

## Step 8: Print titles of first 5 movies
i=0
print("Top 5 movies that you might love to watch: ")
for element in sorted_similar_movies:
    print(get_title_from_index(element[0]))
    i=i+1
    if i>5:
        break
```

CHAPTER 7

INFERENCES/CONCLUSION

- The best choice of k depends upon the data; generally, larger values of k reduces effect of the noise on the classification, but make boundaries between classes less distinct.
- K-NN Algorithm is one of the most suitable approach for rank-ordering problems and also it is widely accepted.
- Feature Selection is one of the most important steps involved in building a machine learning algorithm as we analysed in our model.
- We observed a significant impact to the model output on inclusion/exclusion of certain features.
- We also applied weights to the features to understand their relative importance. This helped in improving the model accuracy as we saw in our graph.
- Distance Measures: We studied about the different distance measures for the purpose of this project.

CHAPTER 8

FUTURE ENHANCEMENTS

- K-NN slow algorithm: K-NN might be very easy to implement but as dataset grows efficiency the speed of algorithm declines very fast.
- Optimal number of neighbors: One of the biggest issues with K-NN is to choose the optimal number of neighbors to be consider while classifying the new data entry.
- Imbalanced data causes problems: k-NN doesn't perform well on imbalanced data. If we consider two classes, A and B, and the majority of the training data is labeled as A, then the model will ultimately give a lot of preference to A. This might result in getting the less common class B wrongly classified.
- Cosine similarity works best when there are a great many (and likely sparsely populated) features to choose from. Under these conditions, Euclidean methods tail off in terms of their sensitivity. Likewise, cosine similarity is less optimal under spaces of lower dimension.
- The selection of features is the most important step for this type of model, if there is any error in selection of features that might decrease the efficiency of the model.
- It is very important to assign correct weightage to each of the feature being used as it affects the prediction strongly.
- For weightage selection we need to analyse some sample data like we did in our project hence it takes little longer time to assign the correct weightage.

CHAPTER 9

REFERENCES

1. "Movie Recommender System Using K-Means Clustering AND K-Nearest Neighbor", Rishabh Ahuja, Arun Solanki, Anand Nayyar
2. Lops P., de Gemmis M., Semeraro G. (2011) Content-based Recommender Systems: State of the Art and Trends. In: Ricci F., Rokach L., Shapira B., Kantor P. (eds) Recommender Systems Handbook. Springer, Boston, MA
3. Hatami, M., Pashazadeh, S. (2014) Improving results and performance of collaborative filtering-based recommender systems using cuckoo optimization algorithm, Int J Comput Appl Volume 88
4. Z. Huang, D. Zeng, H. Chen (2007) A comparison of collaborative-filtering algorithms for e-commerce IEEE Intell Syst, 22, pp
5. R. Burke (2007) Hybrid web recommender systems, Adapt Web, pp.377-408, 10.1007/978-3-540-72079-912
6. Nguyen N.T., Rakowski M., Rusin M., Sobecki J., Jain L.C. (2007) Hybrid Filtering Methods Applied in Web-Based Movie Recommendation System. In: Apolloni B., Howlett R.J., Jain L. (eds) Knowledge-Based Intelligent Information and Engineering Systems. KES 2007. Lecture Notes in Computer Science, vol 4692. Springer, Berlin, Heidelberg
7. Wei D., Junliang C. (2013) The Bayesian Network and Trust Model Based Movie Recommendation System. In: Du Z. (eds) Intelligence Computation and Evolutionary Computation. Advances in Intelligent Systems and Computing, vol 180. Springer, Berlin, Heidelberg
8. D.A. Adeniyi, Z. Wei, Y. Yongquan, Automated web usage data mining and recommendation system using K-Nearest Neighbor (KNN) classification method, Saudi Computer Society, King Saud University, October 2014
9. Vibhor Kanta, Kamal K. Bharadwaj, Enhancing Recommendation Quality of Content-based Filtering through Collaborative Predictions and Fuzzy Similarity Measures, Procedia Engineering Volume 38, 2012, Pages 939-942
10. J. Ben Schafer, Dan Frankowski, Jon Herlocker, and Shilad Sen: Collaborative Filtering Recommender Systems
11. YuanHuang: A news recommendation engine driven by collaborative reader behavior
12. Yibo Wang, Mingming Wang, and Wei Xu: A Sentiment-Enhanced Hybrid Recommender System for Movie Recommendation

13. Sang-Ki Ko, Sang-Min Choi, Hae-Sung Eom , Jeong-Won Cha, Hyunchul Cho, Laehyum Kim, and Yo-Sub Han : A Smart Movie Recommendation System
14. Mark Rethana :Building a Song Recommendation System using Cosine Similarity and Euclidian Distance
15. Mohit Soni , Shivam Bansal: Movie Recommendation System