

# Deep Neural Networks for YouTube Recommendations

## Main objective:

The main objective of the paper is to improve the recommendation system of YouTube using deep learning techniques.

The authors aim to develop a deep neural network architecture that can handle large-scale datasets, integrate different types of data, and outperform the current collaborative filtering approach in terms of accuracy and scalability.

The paper also discusses the challenges and limitations of the proposed approach and suggests directions for future research.

## What are they trying to solve?

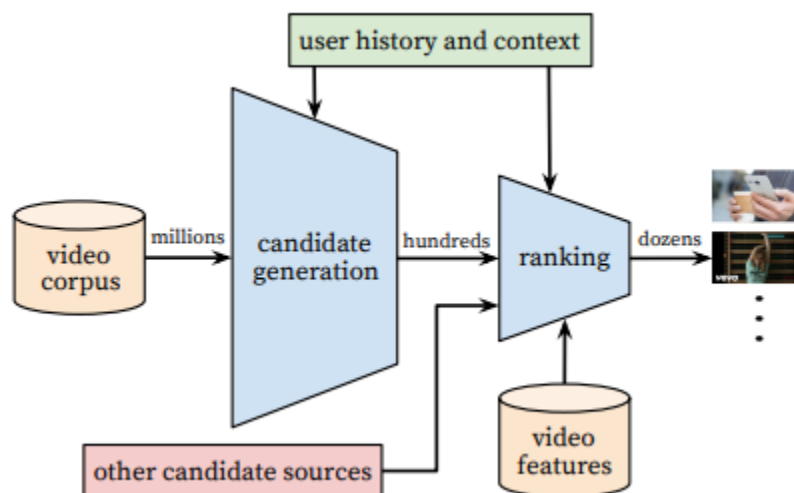
The authors of the paper are trying to solve the limitations of the current collaborative filtering approach used in the recommendation system of YouTube.

This approach has limitations in terms of scalability and the ability to incorporate different types of data.

The authors aim to develop a deep neural network architecture that can handle large-scale datasets and integrate different types of data such as video content, user behavior, and contextual information.

The proposed architecture is expected to improve the accuracy and scalability of the recommendation system and provide a more personalized and diverse set of recommendations to the users.

## Short summary:



The article describes the use of deep learning techniques for improving the recommendation system of YouTube.

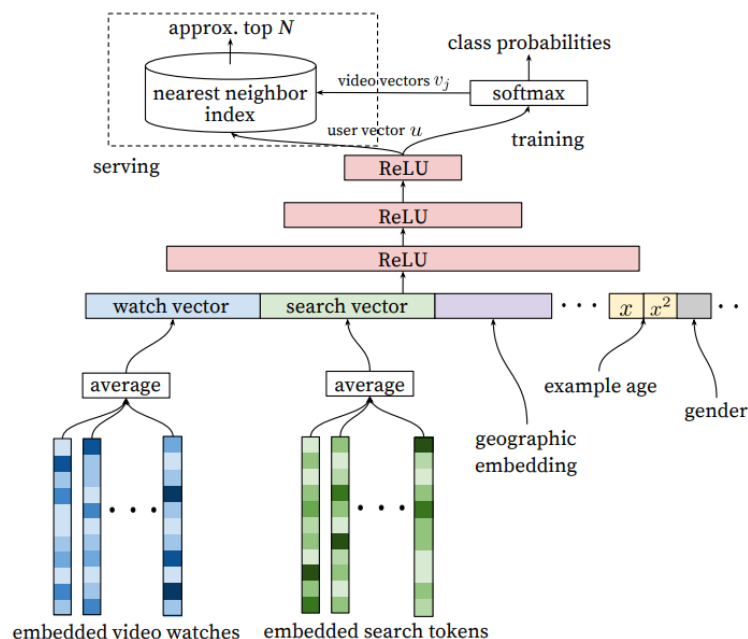
The authors explain how the current recommendation system is based on a collaborative filtering approach that uses matrix factorization techniques. However, this approach has limitations in terms of scalability and the ability to incorporate various types of data.

As we mention before, to overcome these limitations, the authors propose a deep neural network architecture that can handle large-scale datasets and can integrate different types of data such as video content, user behavior, and contextual information. The proposed architecture is based on a multi-layer feedforward neural network that uses embeddings to represent users and items (videos) in a low-dimensional space. The neural network is trained using a combination of supervised and unsupervised learning techniques.

The results of the experiments show that the proposed neural network architecture outperforms the current collaborative filtering approach in terms of accuracy and scalability. The authors also discuss the challenges and limitations of the proposed approach and suggest some directions for future research.

Overall, the article provides an in-depth explanation of the proposed deep neural network architecture for YouTube recommendations and presents experimental evidence to support its effectiveness.

#### Candidate generation process:

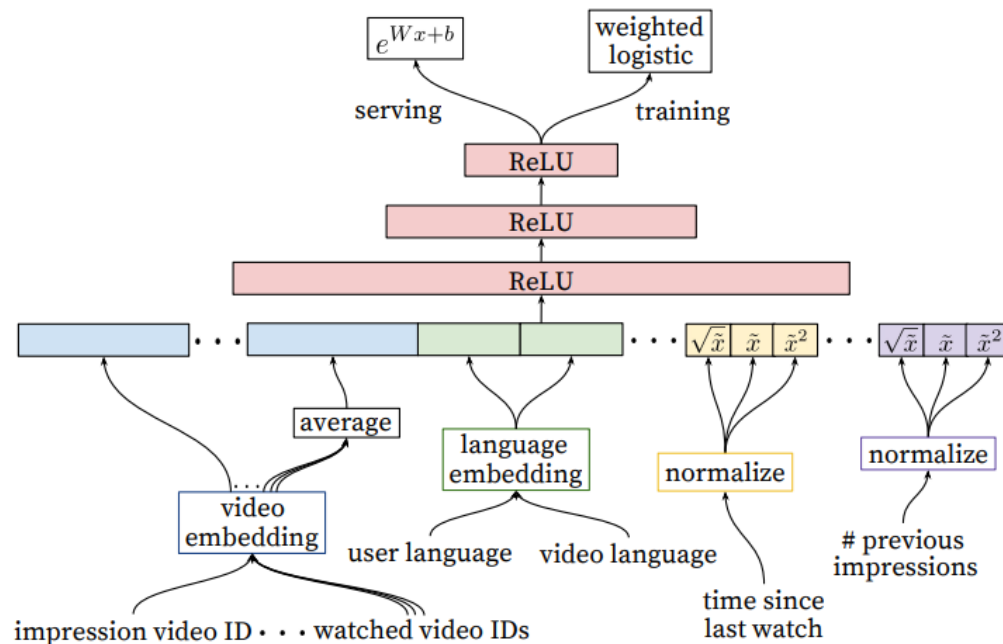


The Candidate generation process in YouTube recommendations involves the selection of a set of videos that are potentially relevant to the user's interests. This process consists of two main steps:

**Retrieval of candidate videos:** In this step, a large pool of candidate videos is retrieved based on various factors such as the user's search history, watch history, and demographic information. The retrieved

videos are usually filtered based on some criteria such as the relevance to the user's interests and the popularity of the video.

### Ranking of candidate videos:



In this step, the retrieved videos are ranked based on their relevance to the user's interests. This involves the use of machine learning algorithms that take into account various features of the video such as its title, description, and tags, as well as the user's watch history and engagement with similar videos.

The candidate generation process is an important step in the recommendation system as it determines the set of videos that will be presented to the user. The accuracy and relevance of the candidate generation process are crucial in ensuring user satisfaction and engagement with the platform.

### Suggesting an improvement:

For the improvement version of the model we decided to use Monte Carlo Dropout technique that allows uncertainty estimation in deep learning models. It involves performing multiple forward passes with dropout enabled at test time to obtain a distribution of predictions for each input. The fact that dropout is enabled during inference time cause the model to give us slightly different prediction each time.

There are several reasons why using Monte Carlo Dropout might improve the algorithm's performance. Firstly, it can help prevent overfitting by adding more regularization to the model during inference time. Secondly, it can provide a measure of the model's uncertainty, which can help make more informed decisions in cases where the model's prediction is less certain. Finally, it can help capture more complex interactions between the input features, which can improve the accuracy of the model. However, it's worth noting that using Monte Carlo Dropout can also increase the computational complexity of the model, so it's important to consider the tradeoff between performance and computational cost.

As a first step, we added dropout to our model, which already showed improvement over the anchor version, and as a second step, we implemented Monte Carlo Dropout.

### **Well-known algorithm for comparison:**

the authors compare their proposed deep neural network architecture with a well-known baseline algorithm called matrix factorization. Matrix factorization is a collaborative filtering approach commonly used in recommendation systems that factorizes the user-item rating matrix into low-dimensional user and item embeddings.

The similarity between the user and item embeddings is used to make personalized recommendations to the users. The authors use matrix factorization as a baseline algorithm to compare the performance of their proposed deep neural network architecture in terms of accuracy and scalability.

According to the experimental results reported in the paper the proposed deep neural network architecture outperforms the baseline matrix factorization approach in terms of accuracy and scalability.

The deep neural network architecture can incorporate different types of data and handle large-scale datasets better than the matrix factorization approach.

The authors also report that the deep neural network architecture provides a more diverse set of recommendations to the users. Therefore, based on the results presented in the paper, the proposed deep neural network architecture is considered better than the baseline matrix factorization approach for YouTube recommendations.

### **Evaluating the algorithms you choose:**

#### **Datasets**

The dataset was provided by MovieLens, a movie recommendation service. It includes the movies and the rating scores made for these movies. contains. It contains 100,000 ratings (1–5) from 943 users on 1682 movies.

#### **Hyperparameter optimization**

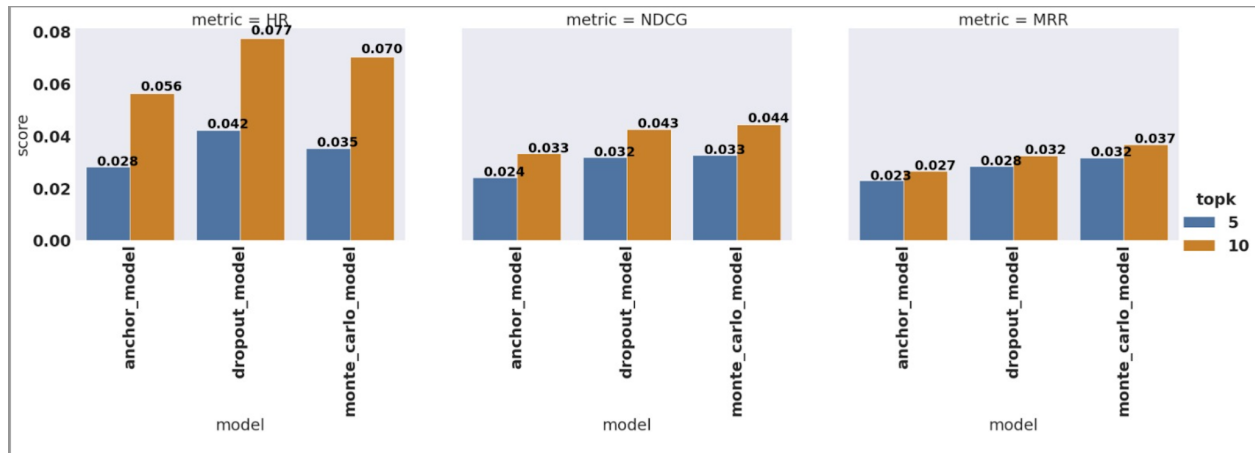
##### **1. Anchor model :**

candidate best generation: number of hidden units: 208, best number of learning rate: 0.01

ranking: best number of hidden units: 176, best number of learning rate: 0.01

##### **2. Monte Carlo model: best number of hidden units: 176, best number of learning rate: 0.01**

## Performance metrics for evaluation



state of the art metrics:

```
Epoch: 887   best rmse: 0.9058916
Epoch: 966   best mae: 0.7128238
Epoch: 204   best ndcg: 0.900631663508036
```

---

## Reporting your conclusions

we see that our model is better than the anchor version but still there are many improvements to make. It is our belief that balancing the data would improve the model performance more if we had more time. As an example, if a user has a very negative attitude towards giving feedback, we would like to normalize those cases as well as the opposite case where users have a very generous attitude toward giving feedback

we think that handling and cleaning the data will do most of the work regarding movielens 100k recommendation.