

Configuration Manual

MSc Research Project Data Analytics

Muhammad Imran Shaikh Student ID: x17119308

School of Computing National College of Ireland

Supervisor: Dr. Muhammad Iqbal

National College of Ireland Project Submission Sheet School of Computing



Student Name:	Muhammad Imran Shaikh
Student ID:	x17119308
Programme:	Data Analytics
Year:	2020
Module:	MSc Research Project
Supervisor:	Dr. Muhammad Iqbal
Submission Due Date:	17/08/2020
Project Title:	Configuration Manual
Word Count:	1458
Page Count:	12

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	
Date:	25th September 2020

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).□Attach a Moodle submission receipt of the online project submission, to
each project (including multiple copies).□You must ensure that you retain a HARD COPY of the project, both for□

your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Muhammad Imran Shaikh x17119308

1 Introduction

The main reason to write this configuration manual is to demonstrate the configuration of system setup, software, and hardware compatibility to run and implement the programming language code which will help to design the Research project and report. This manual will cover sections like System Configuration, Project Development, Codes Implementation, and Experiments with different machine learning models.

2 System Configuration

2.1 Hardware

Processor:3rd Generation Intel Core i5-3320M (2.6 GHz, 3MB L3 cache, 2 cores)1 Up to 3.30 GHz, Ram:16gb,System type:64-bit OS, Graphics: NVIDIA Quadro K2000M with 2GB dedicated DDR3, Windows:10pro(2019)

2.2 Software

Microsoft Office 356: Microsoft Word (For all the professional written stuff), Microsoft Excel(For storing dataset as CSV and Excel format and for visualization purpose too), Microsoft PowerPoint(presentation slides).

Python coding Language: Loading Libraries, Data Cleaning, Data preprocessing and engineering, Initial Data Analysis, Training and Test data splitting, Models Implementation, Hyperparameter Tuning, and Evaluation.Python IDE: Jupyter Notebook and PyCharm.

3 Project Development

The main steps involved in the project development phase are the selection for ideal IDE(Jupyter Notebook) to perform our coding task, Loading suitable libraries, Data cleaning(checking null values and imputations with aggregations), and Data preprocessing and engineering(grouping and joining datasets, Data merging, Describing columns, Removal of unnecessary features, Removal of pipes with string split into Genre columns).Initial data visualization(Word Cloud, several bar charts).

Preparation of separate class for movie dataset to implement Recommendation system techniques(Content-based filtering and Collaborative filtering), Splitting train and test set to fit in various models of our recommendation engine. Getting Top-N movies results from our Recommendation machine learning models by utilizing different techniques. Models Hyperparameter tuning (for extracting the best parameters for optimal results). Models Evaluation(K-fold cross-validation,LOO(Leave One Out) Cross-validation with splits to get better Accuracy results from our machine learning models)Multiple Evaluation plots are plotted to get visual analysis.

3.1 Data Extraction and Pre-processing

The dataset is collected and generated by GroupLens Research Group, Dataset (ml-latestsmall)¹ consist of almost 100k ratings by different users with over 1200 movie tags from 9125 movies. 'ml' stands for movie lens. Each selected user had at least rated 20 movies. There are 4 files included in this dataset named 'movies.csv', 'ratings.csv', 'tags.csv', 'links.csv' but for recommendation purpose, we are considering only two dataset files i.e. 'movies.csv' and 'rating.csv' The entire coding is done in the python programming language. Various python Libraries are imported based on the implementation of different recommendation techniques. All Recommendation system models are imported from 'Surprise Library' which is an official python recommendation system library can be seen in Figure 1. Data preprocessing and engineering are obtained by grouping and joining



Figure 1: Importing python libraries

datasets, Data merging, Describing columns, deletion of unnecessary features, Removal of pipes with string split in Genre columns, Create a function that counts the number of times each genre appear can be seen in Figure 2.

4 Implementation of Recommendation Engine Machine Learning Models

Our recommendation engine is tested with 3 different recommendation system techniques i.e. (Content-based filtering, Collaborative filtering, and Matrix Factorization) by resulting Top-N movies results for users and even for movies in content-based filtering. A separate Movie class is generated to combine both movies and ratings CSV files by grouping them by users and movie ids numbers can be seen in the figure. Different machine

¹https://grouplens.org/datasets/movielens/latest/



Figure 2: Create a function that counts the number of times each genre appears

learning models are implemented with their defined parameters to get better recommendations. Surprise library which is an official python recommendation system library has been utilized to implement machine learning models for our recommendation engine Figure 1.Following are the steps that are taken to create and evaluate our recommendation engine.

4.1 Experiment with Movielens Dataset Analysis

The Analysis on movielens dataset based on the merging of two dataset files i.e. 'movies.csv' and 'rating.csv' as an inner join by movies Ids in fig. With the help of this merged dataset, we are visualizing movie genres by the creation of world cloud and histogram to analyses which movie genres are most popular can be seen in Figure 3, Figure 4. Moreover, top 25 movies with the highest ratings are also plotted to analyze which movie is rated highest by different users in Figure 5.

4.2 Experiment with Content Based Filtering

Content-based filtering works on the phenomena of users' interest in different items. So, based on that interest, similar items are recommended to the users. Recommendation action becomes more accurate if the user provides more input. In our Content-based recommendation engine, we are finding 10 nearest neighbors of the movie of our interest by implementing the KNNBaseline algorithm with a similarity metric of Pearson baseline Figure 6.



Figure 3: Word cloud to analyze which Genre are the popular ones



Figure 4: Histogram to analyze which Genre are the popular ones



Figure 5: Top 25 movies with highest ratings

Using pearson_baseline as a similarity measure in KnnBasic algorithm



Figure 6: Content based filtering technique by using KNNBaseLine algorithm

4.3 Experiment with User and Item(Memory) Based Collaborative Filtering

User and Item-based collaborative filtering are one of the most extensively used techniques in the recommendation system, it works by finding a group of similar users who have given the similar reactions to the item of your interest. The rating matrix is created to find similar users and items based on ratings that are given by the user.KNNwithMean machine learning algorithm along with similarity metric of Cosine is utilized to get Top-10 nearest neighbor movies for specific user Figure 7, Figure 8.





Item based Collaborative Filtering





4.4 Experiment with Matrix Factorization(Model)Based Collaborative Filtering

In Matrix factorization or Model-based collaborative filtering is the Dimensionality Reduction technique just like Principal Component Analysis(PCA). Matrix factorization breaks down the user-item large matrix into a smaller matrix. The hidden features are defined by latent factors that are created by item and user column and row matrix. In our Matrix factorization method, we are implementing two matrix factorization algorithms i.e. SVD(Singular Value Decomposition) and SVDpp(Singular Value Decomposition plus plus) to get Top-10 recommendation results shown in Figure 9, Figure 10.

```
Matrix Factorization
In [230]:
           1 # using above defined trainedset
              trainset = trainset
            2
            3
              def GetAntiTestSetForUser(testSubject):
            4
                      fill = trainset.global_mean
                      anti_testset = []
            5
                       u = trainset.to_inner_uid(str(testSubject))
            6
            7
                       user_items = set([j for (j, _) in trainset.ur[u]])
                       anti_testset += [(trainset.to_raw_uid(u), trainset.to_raw_iid(i), fill) for
            8
           9
                                               i in trainset.all_items() if
                                                i not in user_items]
           10
           11
                       return anti_testset
In [231]:
           1 # Defining svd model and fitting training set
            3
              model svd = SVD()
              model_svd.fit(trainset)
            4
           5
              User = 85
            6
              testset = GetAntiTestSetForUser(User)
            7 predictions_svd = model_svd.test(testset)
In [232]:
           1 recommendations_svd = []
              print("\nMatrix factorization results using SVD for user: {}\n".format(User))
            2
            3
              for userID, movieID, actualRating, estimatedRating, _ in predictions_svd:
            4
                  intMovieID = int(movieID)
                  recommendations_svd.append((intMovieID, estimatedRating))
            5
            6
                   recommendations_svd.sort(key=lambda x: x[1], reverse=True)
            7
            8
              for ratings in recommendations_svd[:10]:
            9
                   print(ml.getMovieName(ratings[0]))
          Matrix factorization results using SVD for user: 85
          Eternal Sunshine of the Spotless Mind (2004)
          Matrix, The (1999)
          Dark Knight, The (2008)
          Maltese Falcon, The (1941)
          Wallace & Gromit: A Close Shave (1995)
          Her (2013)
          Godfather, The (1972)
          Lock, Stock & Two Smoking Barrels (1998)
          Indiana Jones and the Last Crusade (1989)
          Roger & Me (1989)
```

Figure 9: Matrix Factorization technique by using SVD algorithm

SVD++ Algorithm



Figure 10: Matrix Factorization technique by using SVD++ algorithm

5 Experiment with Models Hyperparameter Tuning

For Models Hyperparameter tuning we are considering Grid Search CV from Python Surprise Library, which provides us the best parameters to get optimal value from our machine learning models when we are training our dataset. The main parameters which are considered for KNN and Matrix factorization-based algorithms are different K-Values, Epoches no, learning rate, Similarity options, and accuracy measures(RMSE and MAE) can be seen in Figure 11, Figure 12.

In [34]:	<pre>1 param_grid_K&M_UB = ('k': [10, 20, 50], 'n_epochs': [1, 3, 5],</pre>	
	Computing the mad similarity matrix Done computing similarity matrix Computing the cosine similarity matrix Done computing similarity matrix	•
In [35]:	<pre>1 # We can now use the algorithm that yields the best rmse: 2 algo KM4/UB - gg KM4/UB.best estimator['rmse'] # pass the best model to algo 3 print(gg.KM4/UB.best.params['mse']) the best RNSE score 5 print(gg.KM4/UB.best.params['mse']) 6 cross_validate(algo KM4/UB, data, measures-['RMSE', 'MAE'], cv-5, verbose-True} 0.250512465135269 {'k': 50, 'n_epochs': 1, 'learning_rate': 1e-05, 'sim_options': {'name': 'msd', 'min_support': 1, 'user_based': True}} Done computing time ad similarity matrix. Done computing similarity matrix.</pre>	





Figure 12: Hyperparameter tuning with GridSearchCV for different params of Matrix Factorization models

6 Experiment with Models Evaluation

For Models Cross-Validation, we are considering two cross validators to test the accuracy of our recommendation engine models in different splits. Those two Cross validators are K-Fold cross Validator and LOOCV leave one out cross validator can be seen in Figure 13, Figure 14

K-F	old Cross Validator
# defi	ne a cross-validation iterator
kf = K	Fold(n_splits=3)
# Defi	ning models
<pre># KNW sim_op algo_k</pre>	Baseline Algo for Content based Recommendation tion_UNN + ('name': 'peerson_baseline', 'user_based': False) MLCE = NUNBaseline(sin_option=sim_option_UNN)
# KNW sim_op algo_K	with mean Algo for User based Collaborative Filtering tions = {'rame': 'cosine', 'user_based : Trme} Wills = NUMHINTERNIK-50, sing potion=sing potions)
<pre># KNW sim_op algo_k</pre>	with mean Algo for Item bosed Collaborative Filtering tions_ID = { "name": 'cosine','user_based': False) Will E = NNMMithems(kc30, sim_options=Sim_options_ID)
# SVD(algo_s	Singular Value Decomposition) model for Matrix Factorization vd = SVD()
# SVDp algo_s	p(Singular Value Decomposition)plus plus model for Matrix Factorization vdpp = SVDpp()
<pre># Dict rmse_K mae_KM rmse_K mae_KM rmse_K mae_SV rmse_s mae_sv</pre>	ionary for folds accuracy N = [] N = [] Wi [] = [] Wi [] = [] Wi [] = [] Vi = [] Vi = [] Vi = [] Vi = []
# Trai for tr	ning and testing data on basis of K-Folds ainset, testset in kf.split(døta):
#	Fitting Train and test sets in models for predictions
# al pr	XXW Baseline Algo for Content based Recommendation go.XXW_CG.Fit(trainset) go.XXW_CG.Fit(tos_UNM_CB.test(testset)
# al pr	KWW with mean Algo for User based Collaborative Filtering go,KWH_UB.fit(trainset) Bittions_GMU_UB = algo_XBM_UB.test(testset)
# al pr	KMW with mean Algo for Item based Collaborative Filtering go,SWDLD.Fit(trainset) dictions_SWDLB = algo_SWP_IB.test(testset)
# al pr	SVD(Singular Value Decomposition) model for Matrix Factorization gg_svd.fit(trainset) editions_vd= = algg_svd.test(testset)
# al	SVDpp(Sinpular Value Decomposition)plus plus model for Matrix Factorization go_svdpp.fit(trainset)

Figure 13: Models Evaluation with K-fold Cross Validation

Leave One Out Cross Validator

```
In [38]: # define a Leave one out cross-validation iterator
                  Looc = LeaveOneOut(n_splits=3)
                  # Defining models
                  # KNW BaseLine Algo for Content based Recommendation
sim_option_KNW = {'name': 'pearson_baseline', 'user_based': False}
algo_KNN_CB = KNNBaseline(sim_options=sim_option_KNN)
                  # KNN with mean Algo for User based Collaborative Filtering
sim_options = {'name': 'cosine','user_based': True}
algo_KNM_UB = KNNWithMeans(k=50, sim_options=sim_options)
                  # KNW with mean Algo for Item based Collaborative Filtering
sim_options_IB = {'name': 'cosine','user_based': False}
algo_KuM_IB = KNNWithMeans(k=50, sim_options=sim_options_IB)
                  # SVD(Singular Value Decomposition) model for Matrix Factorization
algo_svd = SVD()
                  # SVDpp(Singular Value Decomposition)plus plus model for Matrix Factorization
algo_svdpp = SVDpp()
                  # Dictionary for folds Of LOOCV accuracy
                 # Dictionary for fo
rmse_KNN_L= []
nmae_KNM_UB_L = []
rmse_KNM_UB_L = []
rmse_KNM_UB_L = []
rmse_KNM_IB_L = []
rmse_svd_L = []
rmse_svd_L = []
rmse_svdpp_L = []
                  # Training and testing data on basis of Leave One Out CV
for trainset, testset in Looc.split(data):
                         # Fitting Train and test sets in models for predictions
                         # KNN Baseline Algo for Content based Recommendation
algo_KNN_CB.fit(trainset)
predictions_KNN = algo_KNN_CB.test(testset)
                          # KNW with mean Algo for User based Collaborative Filtering
algo_KNM_UB.fit(trainset)
predictions_KNM_UB = algo_KNM_UB.test(testset)
                          # KNN with mean Algo for Item based Collaborative Filtering
algo_KNM_IB.fit(trainset)
predictions_KNM_IB = algo_KNM_IB.test(testset)
                          # SVD(Singular Value Decomposition) model for Matrix Factorization
algo_svd.fit(trainset)
predictions_svd = algo_svd.test(testset)
                          # SVDpp(Singular Value Decomposition)plus plus model for Matrix Factorization
algo_svdpp.fit(trainset)
predictions_svdpp = algo_svdpp.test(testset)
```

Figure 14: Models Evaluation with LOO(Leave One Out) Cross Validation