f III	MAtriculation nr313301 Lab 10 Q2 and Q3 Exercise 2: Implementing basic matrix factorization (MF) technique or recommender systems this task you are required to implement a matrix factorization (MF) technique for recommender systems. You are given a rating matrix Rn×m and you have to learn latent matrices P n×k and Qk×m, where n is the number of users, m is the number of items and k the latent dimer ou can solve the MF problem by implementing Stochastic Gradlent Descent (SGD) learning algorithms (Algorithm LearnLatentFactors on slide 29). Measure the prediction quality (the RMSE score) on the validation and test dataset. You can set 10%/10% of ratings aside for valid import numpy as no import tensorf low as tf from tensorf low as not from yellow import tensor singer tensors as no from yellow import tensor as ans from tensorf low import tensors as no import tensor as ans from plate import tensor flow import tensors as no import tensor flow import flow imp
	import matplotlib.pyplot as plt import matplotlib import rc from google.colab import train_test_split from sklearn.model_selection import train_test_split from collections import counter, defaultdict from sklearn.metrics import accuracy_score import matplotlib.ticker as ticker from math import sqrt from sklearn.metrics import mean_squared_error rcParams['figure.figsize'] = 14, 8 RANDOM_SEED = 3116 np.random.seed(RANDOM_SEED)
	#Route of at content of the main data of
1	9997 76 1990 1 879795543 9998 13 225 2 88239156 9999 12 203 3 879959583 00000 rows × 4 columns # loading the item part of the movie lens dataset item_data = pd.read_csv('content/drive/MyDrive/Recommender/u.item', sep="[", names=["movieid", "movietitle","releasedate","videoreleasedate", "IMDURL", "conedy", "Crime", "Motory", "Adventure", "Animation", "Children's", "Comedy", "Crime", "Documentary", "Drama", "Fantasy", "Film-Noir", "Motory", "
:	COLUMNS="1-TentId" values="rating" status="status="rating" status="sta
	94
S	### supplication of the dataset into 80% train, 10% validation and 10%test dataset train, val, test = np.split(ratings_df.sample(frac=1, random_state=3116), [int(.8*len(ratings_df)), int(.9*len(ratings_df))]) train_list=np.array(train) validate_list=np.array(val) test_list=np.array(train) np.av(np.stann(train_list)) val=val.replace(np.nan, 0) validate_list=np.array(val) test_test_replace(np.nan, 0) test_test_test_replace(np.nan, 0) test_test_test_test_test_test_test_test
: [# defining a RMSE function which gives mean squared error comparing the ground truth and predicted values def rmse(prediction, ground_truth): prediction = prediction[ground_truth.nonzero()].flatten() ground_truth = ground_truth[ground_truth.nonzero()].flatten() return sqrt(mean_squared_error(prediction, ground_truth)) defining a class which initialises a function, normalises the P and Q matrices and implementing stochastic gradient descent class Recommender: definit(self, n_epochs=2, n_latent_features=3, lmbda=0.1, learning_rate=0.001): self.n_epochs = n_epochs self.n_latent_features = n_latent_features
	<pre>self.lmbda = lmbda self.learning_rate = learning_rate def predictions(self, P, Q): return np.dot(P.T, Q) def fit(self, X_train, X_val): m, n = X_train.shape #Mormalising the dataset of movie lens and scaling it such that specific users don't always rate high or low self.P=np.random.normal(scale=1./self.n_latent_features, size=(self.n_latent_features, m)) self.Q=np.random.normal(scale=1./self.n_latent_features, size=(self.n_latent_features, n)) # self.P = 3 * np.random.rand(self.n_latent_features, m) # self.Q = 3 * np.random.rand(self.n_latent_features, n) self.train_error = [] self.val_error = []</pre>
	<pre>users, items = X_train.nonzero() #users, items = X_train.shape #print(users, items) #Stochastic gradient descent for epoch in range(self.n_epochs): for u, in zip(users, items): error = X_train[u, i] - self.predictions(self.P[:,u], self.Q[:,i]) self.P[:, u] += self.learning_rate * (error * self.P[:, u] - self.lmbda * self.Q[:, i]) self.P[:, i] += self.learning_rate * (error * self.P[:, u] - self.lmbda * self.Q[:, i]) train_rmse = rmse(self.predictions(self.P, self.Q), X_train) val_rmse = rmse(self.predictions(self.P, self.Q), X_val) self.train_error.append(train_rmse) self.val_error.append(val_rmse) print("the final test rmse with epochs=", self.n_epochs, "latent dimensions=", self.n_latent_features, "lambda=", self.lmbda, "alpha=" ,self.learning_rate, "parameters is", val_rmse) return self def predict(self, X_train, user_index):</pre>
: : : : : : : : : : : : : : : : : : :	y_hat = self.predictions(self.P, self.Q) predictions_index = np.where(X_train[user_index, :] == 0)[0] return y_hat[user_index, predictions_index].flatten() 2) optimize the hyper-parameters i.e. λ regularization constant, α learning rate, k latent dimensions deeping the epochs to 25 as fixed and varying the alpha, lambda and latent dimensions i.e the hyperparameter optimisation. combinations used:-alpha=[0.0001,0.01,0.01,0.1] lambda=[0.01,0.1] latent dimensions=[3,5,8] recommender1 = Recommender() recommender1init(n_epochs=25, n_latent_features=3, lmbda=0.01, learning_rate=0.0001) return y_hat [user_index, predictions index]. the final test rmse with epochs= 25 latent dimensions= 3 lambda= 0.01 alpha= 0.0001 parameters is 3.7085368315574447 <_mainRecommender at 0x7f2b39579190> recommender2 = Recommender() recommender2init(n_epochs=25, n_latent_features=3, lmbda=0.1, learning_rate=0.0001)
: : [recommender2.fit(train_list, validate_list) the validation rnse with epochs= 25 latent dimensions= 3 lambda= 0.1 alpha= 0.0001 parameters is 2.278171679475042 <_mainRecommender3 = Recommender() recommender3init(n_epochs=25, n_latent_features=5, lmbda=0.01, learning_rate=0.0001) recommender3fit(train_list, validate_list) the final test rmse with epochs= 25 latent dimensions= 5 lambda= 0.01 alpha= 0.0001 parameters is 4.058583280050861 <_mainRecommender at 0x7f2b3ae65ad0> recommender4init(n_epochs=25, n_latent_features=5, lmbda=0.1, learning_rate=0.0001) recommender4init(n_epochs=25, n_latent_features=5, lmbda
:	recommender5 = Recommender() recommender5init(n_epochs=25, n_latent_features=5, lmbda=0.1, learning_rate=0.001) recommender5.fit(train_list, validate_list) the validation rmse with epochs= 25 latent dimensions= 5 lambda= 0.1 alpha= 0.001 parameters is 1.3784116075020219 <_mainRecommender at 0x7fad481c5310> recommender6 = Recommender() recommender6init(n_epochs=25, n_latent_features=5, lmbda=0.01, learning_rate=0.001) recommender6.fit(train_list, validate_list) the validation rmse with epochs= 25 latent dimensions= 5 lambda= 0.01 alpha= 0.001 parameters is 1.4689259939462378 <_mainRecommender at 0x7fad48184310> recommender7 = Recommender()
: [recommender7init(n_epochs=25, n_latent_features=5, lmbda=0.01, learning_rate=0.01) recommender7init(n_epochs=25, n_latent_features=5, lmbda=0.01 alpha= 0.01 parameters is 1.2090601952853912 <mainrecommender 0x7fad48184dd0="" at=""> recommender8 = Recommender() recommender8init(n_epochs=25, n_latent_features=5, lmbda=0.1, learning_rate=0.01) recommender8.fit(train_list, validate_list) the validation rmse with epochs= 25 latent dimensions= 5 lambda= 0.1 alpha= 0.01 parameters is 1.147437249408723 <mainrecommender 0x7fad481c5f10="" at=""> recommender11 = Recommender() recommender11init(n_epochs=25, n_latent_features=8, lmbda=0.01, learning_rate=0.0001)</mainrecommender></mainrecommender>
:	recommender11.fit(train_list, validate_list) the validation rmse with epochs= 25 latent dimensions= 8 lambda= 0.01 alpha= 0.0001 parameters is 4.804520272102495 <_mainRecommender at 0x7fad4818ee50> recommender12 = Recommender() recommender12init(n_epochs=25, n_latent_features=8, lmbda=0.1, learning_rate=0.0001) recommender12.fit(train_list, validate_list) the validation rmse with epochs= 25 latent dimensions= 8 lambda= 0.1 alpha= 0.0001 parameters is 4.038553881960411 <_mainRecommender at 0x7fad482cbb10> recommender13 = Recommender() recommender13init(n_epochs=25, n_latent_features=8, lmbda=0.01, learning_rate=0.0001) recommender13init(n_epochs=25, n_latent_features=8, lmbda=0.01, learning_rate=0.0001) recommender13init(n_epochs=25, n_latent_features=8, lmbda=0.01, learning_rate=0.0001) recommender13init(n_epochs=25, n_latent_features=8, lmbda=0.01, learning_rate=0.0001)
: [the validation rmse with epochs= 25 latent dimensions= 8 lambda= 0.01 alpha= 0.001 parameters is 1.9484016280136647 <pre></pre>
:	the validation rmse with epochs= 25 latent dimensions= 8 lambda= 0.01 alpha= 0.01 parameters is 1.3090677151549854 <_mainRecommender at 0x7fad4818ef90> recommender17 = Recommender() recommender17init(n_epochs=25, n_latent_features=8, lmbda=0.1, learning_rate=0.1) recommender17.fit(train_list, validate_list) //usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:46: RuntimeWarning: overflow encountered in multiply //usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:45: RuntimeWarning: overflow encountered in multiply //alueError
	2 recommender17init_(n.epochs=25, n_latent_features=8, lmbda=0.1, learning_rate=0.1)
	437 """ 438
F	799 if force_all_finite:> 800
T F	Training Data Training Data Training Data Training Data Training Data
	2.25 2.00 WE 1.75 1.50 1.25 1.00 0 5 10 15 20 25 Number of Epochs
\ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \	low, looking at the various combinations of different learning rate, lambda and latent dimensions, we can clearly see that the rmse is minimum when n_latent_features=5, lmbda=0.1, learning_rate=0.01. Therefore, we shall calculate our test rmse finally on these set of hyperparar Ve calculate our final rmse results on the test dataset from the optimal hyperparameters we get on the validation set recommender_final = Recommender() recommender_finalinit(n_epochs=25, n_latent_features=5, lmbda=0.1, learning_rate=0.01) recommender_final.fit(train_list, test_list) the final test rmse with epochs= 25 latent dimensions= 5 lambda= 0.1 alpha= 0.01 parameters is 1.1209035594368975 <mainrecommender 0x7fad47d18990="" at=""> 2) Compute the test RMSE so, the final test rmse is 1.12090</mainrecommender>
	Exercise 3: Recommender Systems using matrix factorization sckitlearn In this task you are required to use off-the-shelf libraries such as libmf or sckit-learn. You have to learn a matrix factorization model. Optimize the hyper parameters and perform a 3-fold cross validation. Compare your results with the results in task 1. List in detail which/how you us braries?, what it solves?, and why it is selected? Present your results in form of plots and tables. Since surprise is a library from scikit learn, I am using this for the question (and as i can use any off the shelf library). Please find another implementation of library at the end) I plip install scikit-surprise -q I plip install scikit-surprise (setup.py) done import pandas as pd from surprise import Reader from surprise import Dataset
]: [from surprise import NormalPredictor from surprise import NormalPredictor from surprise import NMF from surprise import NMF from surprise import nMF from surprise import accuracy import rmse from surprise import accuracy from surprise import accuracy from surprise.model_selection import train_test_split from surprise.model_selection import GridSearchCV all_data.drop(['timestamp'], axis=1, inplace=True) #drop the timestamp column reader = Reader(rating_scale=(1, 5)) #create a reader to load the dataset following data = Dataset.load_from_df(all_data[['userId', 'itemId', 'rating']], reader) axplaining how this methodology works:-Answers which, how questions??
T T V a	Surprise library reduces the user-item interaction into the lower dimensional space latent matrix. However, there is no dimensionality reduction technique like SVD or PCA applied under the hood. This method estimates the latent factor matrix and the bias termed minimizing the loss between the original explicit rating and the reconstructed prediction rating. Therefore, that is why it is called approach model-based. To make a predicted rating, we can multiply the estimated latent factors matrix and add the bias term for the user or item like the below figure. $\hat{r}_{ui} = \mu + b_u + b_i + q_i^T p_u$ Where the u is denoting user index, and i is denoting the item index. (hat)ui is the predicted rating of the user u on the item integration of the user u on the item integration.
b c p s	is the bias term (baseline) for the user or item (paseline) f
F N C	why is this library selected? Provides various ready-to-use prediction algorithms such as baseline algorithms, neighborhood methods, matrix factorization-based (SVD, PMF, SVD++, NMF), and many others. Also, various similarity measures (cosine, MSD, pearson) are built-in. Adakes it easy to implement new algorithm ideas. Provide tools to evaluate, analyse and compare the algorithms' performance. Cross-validation procedures can be run very easily using powerful CV iterators (inspired by scikit-learn excellent tools), as well as exhaustive search over a set of parameters. benchmark = [] rmse_list=[] # Iterate over the below two kinds of algorithms # I tried out two most themost used methodologies in Matrix factorisation which are SVD and NMF
	<pre>algorithms = [SVD(), NMF()] print ("Attempting: ", str(algorithms), '\n\n\n') for algorithm in algorithms: print("Starting: ", str(algorithm)) # Perform 3 fold cross validation and only use RMSE as performance metric results = cross_validate(algorithm, data, measures=['RMSE'], cv=3, verbose=False) # results = cross_validate(algorithm, data, measures=['RMSE', cv=3, verbose=False) rmse_list.append(results['test_rmse']) # Get results & append algorithm name tmp = pd.DataFrame.from_dict(results).mean(axis=0) tmp = tmp.append(pd.Series([str(algorithm).split(' ')[e].split(',')[-1]], index=['Algorithm']))</pre>
	benchmark.append(tmp) print("Done: ",str(algorithm), "\n\n") print ('\n\tDONE\n') Attempting: [<surprise.prediction_algorithms.matrix_factorization.svd 0x7f851ecc6ed0="" at="" object="">, <surprise.prediction_algorithms.matrix_factorization.svd 0x7f851ecc6ed0="" at="" object=""> Starting: <surprise.prediction_algorithms.matrix_factorization.svd 0x7f851ecc6ed0="" at="" object=""> Done: <surprise.prediction_algorithms.matrix_factorization.svd 0x7f852ecc6ed0="" at="" object=""> Starting: <surprise.prediction_algorithms.matrix_factorization.nmf 0x7f852eb195d0="" at="" object=""> Done: <surprise.prediction_algorithms.matrix_factorization.nmf 0x7f852eb195d0="" at="" object=""> Done: <surprise.prediction_algorithms.matrix_factorization.nmf 0x7f852eb195do="" at="" object=""></surprise.prediction_algorithms.matrix_factorization.nmf></surprise.prediction_algorithms.matrix_factorization.nmf></surprise.prediction_algorithms.matrix_factorization.nmf></surprise.prediction_algorithms.matrix_factorization.svd></surprise.prediction_algorithms.matrix_factorization.svd></surprise.prediction_algorithms.matrix_factorization.svd></surprise.prediction_algorithms.matrix_factorization.svd>
F	rmse_list [array([0.94345155, 0.94563508, 0.94394977]), array([0.97427259, 0.97688421, 0.97477224])] Plotting a table of the test rmse and its time of running against both the selected algorithms. surprise_results = pd.DataFrame(benchmark).set_index('Algorithm').sort_values('test_rmse') test_rmse
:	Algorithm SVD 0.944345 4.103123 0.379911 NMF 0.975310 4.988396 0.347002 # smaller grid for lesser running time param_grid = { "n_epochs": [10, 20, 40, 75], "!r_all": [0.002, 0.005], #alphalearning_rate "reg_all": [0.02] #lambda } gs = GridSearchCV(SVD, param_grid, measures=["rmse", "mae"], refit=True, cv=3) gs.fit(data)
	training_parameters = gs.best_params["rmse"] print("BEST RMSE: \t", gs.best_score["rmse"]) print("BEST MAE: \t", gs.best_score["mae"]) print("BEST params: \t", gs.best_params["rmse"]) BEST RMSE: 0.9432340469672061 BEST MAE: 0.7440505288284266 BEST params: {'n_epochs': 20, 'lr_all': 0.005, 'reg_all': 0.02} results_df = pd.DataFrame.from_dict(gs.cv_results) results_df results_df
	Part
:	rmse_split_cols=[0.946623,0.942259,0.949082] plt.xlabel("3 fold splits") plt.ylabel("rmse per split") plt.plot([1,2,3],rmse_split_cols) [<matplotlib.lines.line2d 0x7f851a884650="" at="">] 0.949- 0.948-</matplotlib.lines.line2d>
	0.947 - 0.945 - 0.944 - 0.943 -
ī.	200 1/25 1/50 1/75 2/00 2/25 2/50 2/75 3/00 Comparison of results from task 2 against task 3 The test rmse achieved using the hand made function of matrix factorisation gave a RMSE of 1.12 whereas using the surprise library we got better results of 0.9434. Hence, the library does implement the algorithm in a better fashion than the crude algorithm developed from scratce implementation of the same question as above but with a different library pip install matrix_factorization
	<pre>from matrix_factorization import BaselineModel, KernelMF, train_update_test_split import pandas as pd from sklearn.metrics import mean_squared_error cols = ["user_id", "item_id", "rating", "timestamp"] movie_data = pd.read_csv("/content/drive/MyDrive/Recommender/u.data", names=cols, sep="\t", usecols=[0, 1, 2], engine="python") X = movie_data[["user_id", "item_id"]] y = movie_data["rating"] # Prepare data for online learning (</pre>
	<pre>y_train_initial, x_train_update, y_train_update, X_test_update, y_test_update,) = train_update_test_split(movie_data, frac_new_users=0.2) # Initial training matrix_fact = KernelMF(n_epochs=2, n_factors=100, verbose=1, lr=0.001, reg=0.005) matrix_fact.fit(X_train_initial, y_train_initial) # Update model with new users atrix_fact.update_users(X_train_update, y_train_update, lr=0.001, n_epochs=20, verbose=1) pred = matrix_fact.predict(X_test_update) rmse = mean_squared_error(y_test_update, pred, squared=False)</pre>
	rmse = mean_squared_error(y_test_update, pred, squared=False) print(f"\nTest RMSE: {rmse:.4f}") # Get recommendations user = 200 items_known = X_train_initial.query("user_id == @user")["item_id"] matrix_fact.recommend(user=user, items_known) Epoch 1 / 2 - train_rmse: 1.072446827939706 Epoch 2 / 2 - train_rmse: 1.03233436500082 Epoch 1 / 2 - train_rmse: 1.038233436500082 Epoch 1 / 2 - train_rmse: 1.06775766358668 Epoch 2 / 2 - train_rmse: 1.06775766358668 Epoch 3 / 20 - train_rmse: 1.06775766358668 Epoch 4 / 20 - train_rmse: 1.0584399074213695 Epoch 5 / 20 - train_rmse: 1.0584399074213695 Epoch 6 / 20 - train_rmse: 1.054039074213695 Epoch 6 / 20 - train_rmse: 1.054039074213695 Epoch 6 / 20 - train_rmse: 1.054039074213695 Epoch 7 / 20 - train_rmse: 1.0540319374239423
: _	usr id item. id item. id rating_pred 116 200 112 4.2272 157 200 801 4.08305 48 200 181 4.053374 15 200 127 4.03001 200 183 4.02665 87 200 183 4.02665 80 201 4.02400 10 4.02400 10 4.02400
	679 200 697 4.005560 References:- 1) https://towardsdatascience.com/a-complete-guide-to-recommender-system-tutorial-with-sklearn-surprise-keras-recommender-5e52e8ceace1
3	thttps://githubhelp.com/Quang-Vinh/matrix-factorization thttps://www.quuxlabs.com/blog/2010/09/matrix-factorization-a-simple-tutorial-and-implementation-in-python/ thttps://towardsdatascience.com/comprehensive-data-explorations-with-matplotlib-a388be12a355