

دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

Information Retrieval Course

Fall 2023

IR-HW03

Koorosh Khavari Moghaddam

koorosh.k@aut.ac.ir

koorosh.khavari.m@gmail.com

Student ID: 401131023

Contents

۳	Preprocessing Data
۴	Matrix Factorization
۵	Item Representation Using BERT
۶	Concatenation Of Bert & MF
۶	Evaluation
۷	Evaluation Methods
۸	BERT Model Evaluation

Making Representation Vectors

```
1 # clean the text
2 def clean_text(text):
3     return gsp.remove_stopwords(gsp.strip_tags(gsp.strip_punctuation(gsp.strip_numeric(text))))
4 D_train['review_text'] = D_train['review_text'].astype('string')
5 D_train['review_text'] = D_train['review_text'].fillna('')
6 D_train['review_text'] = D_train['review_text'].apply(clean_text)
7 # tokenize and transform the text
8 tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
9 model = BertModel.from_pretrained('bert-base-uncased')
10 reviews = [torch.tensor([tokenizer.encode(sentence, add_special_tokens=True)][:-512] for sentence in list(D_train['review_text']))]
11 with torch.no_grad():
12     results = []
13     for review in tqdm(reviews):
14         try: results += [model(review).last_hidden_state.mean(dim=1).squeeze()]
15         except: print(f"[error] : {review}")
16     reviews_vec = torch.stack(results)
17 # store the texts
18
19 with open ("bert_review_vectors.pth", "wb") as f:
20     pickle.dump(reviews_vec, f)
```

✓ 115m 28.6s
100% ██████████ 23038/23038 [1:54:29<00:00, 3.35it/s]

Preprocessing Data

In order to tokenize the comments inside the dataset, we must clean the texts and reduce them to 512 tokens, said procedure is done using gensim library, also the processed representations are stored in the bert_review_vectors.pth file.

[1:6] we cleared the text of the comments using gensim library by iterating through all of them.

[7:16] we tokenize the cleansed texts and represent them as described in the homework problem (centroid of the transformed tokens in the latent space of pretrained BERT model)

[17:20] storing the result of the model for the later use.

Matrix Factorization

Learn U & I Matrices (As Ground Truth For Evaluation)

```
1 Latent_Dimension = 64
2  $\lambda$  = 1e-3
3  $\tau$  = 1e-3
4 U = np.random.uniform(low=0, high=1, size=(Users.shape[0],Latent_Dimension))
5 I = np.random.uniform(low=0, high=1, size=(Items.shape[0],Latent_Dimension))
6 dU, dI = np.zeros_like(U), np.zeros_like(I)
7 N = D_train.shape[0]
8 def E(R,I,U): return np.sum(R*((R-U@I.T)**2))/N+ $\lambda$ *(np.linalg.norm(U)+np.linalg.norm(I))
9 def _dU(R,I,U): return np.sum(-2*(R-U@I.T)@I)/N+ $\lambda$ *2*U
10 def _dI(R,I,U): return np.sum(-2*(R-U@I.T).T@U)/N+ $\lambda$ *2*I
11 _ $\Delta$  = []
12 for i in tqdm(range(100)):
13     # calculate loss
14      $\Delta$  = E(R,I,U)
15     _ $\Delta$  += [ $\Delta$ .item()]
16     # calculate gradient vectors
17     dU = _dU(R,I,U)
18     dI = _dI(R,I,U)
19     # calculate descent direction
20     dU = dU/np.linalg.norm(dU)
21     dI = dI/np.linalg.norm(dI)
22     # bisection to find step size
23     while  $\Delta$ <E(R,I,U-dU): dU=dU/2
24     while  $\Delta$ <E(R,I-dI,U): dI=dI/2
25     # update factors
26     U -= dU
27     I -= dI
28     # early stopping
29     if (np.linalg.norm(dI)< $\tau$ )&(np.linalg.norm(dU)< $\tau$ ): break
30
31 plt.plot(list(range(len(_ $\Delta$ ))), _ $\Delta$ )
32 plt.show()
33
34 # real vs approximate utility matrix at (0,0)
35 R1 = U@I.T
```

✓ 21.1s

99%|██████████| 99/100 [00:20<00:00, 4.72it/s]

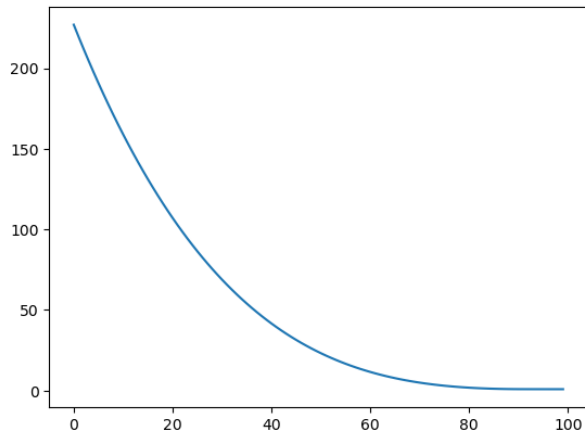
I changed the matrix naming to [U]ser and [I]tem instead of Q and P for simpler reading.

[1:7] we defined the requireid variables for the training procedure. τ parameter is an threshold used to early stop for the training procedure. dU and dI will used to store the gradient of the loss with respect to each of the User and Item Matrices. N is a normalization variable equal to the size of the trainig dataset, used in the calculation of the error term to make it stable.

[8:10] required functions are defined, E is the error term as defined in the homework problem, _dU and _dI is the gradient of the error term with respect to matrices U and I respectively.

[11] _ Δ is a variable used to store the history of the loss during the training.

[12:29] main training procedure, here we first calculate and stored the error[13:15], then calculated the descent direction by normalizing the resulted gradient vector of the loss with respect to each of the U and I matrices [16:21], then we find the appropriate step size using bisection method [22:24] and update current U and I matrices towards it [25:27], we also checked the value of the error to early stop the training process if needed in [28-29].



The final training plot is shown in the side figure.

Item Representation Using BERT

Part 3.A:: Replacing I Matrix

```

1 # load representation vectors
2 with open ("bert_review_vectors.pth", "rb") as f:
3     reviews_vec = pickle.load(f).numpy()
4 # make 64 dimensional representation of comments
5 pca = PCA(n_components=64)
6 pca.fit(reviews_vec)
7 reviews_pca_vec = pca.transform(reviews_vec)
8 # make Item Comment matrix
9 Item_Comment = []
10 PCA_Item_Comment = []
11 for item in Items: PCA_Item_Comment += [np.mean(reviews_pca_vec[(D_train.loc[D_train['item_id'] == item].index)[0].astype(int)], axis=0)]
12 for item in Items: Item_Comment += [np.mean(reviews_vec[(D_train.loc[D_train['item_id'] == item].index)[0].astype(int)], axis=0)]
13 bert_I = np.array(Item_Comment)
14 pca_I = np.array(PCA_Item_Comment)
15 R2 = U@pca_I.T

```

✓ 4.6s

In the second part of the homework, we used the BERT model to represent each item in the latent space.

[1:3] we load the stored representation which are prepared in the first part.

[4:7] we transform each comment to 64-dimensional space using PCA module of the sklearn.

[8:15] for each item, we represent it as the centroid of the all-related comments in the 64-dimensional space. Therefore the result is the new Item Matrix with the same shape of the older I matrix, and as told in the homework, we used the U matrix of the matrix factorization part alongside it.

Concatenation Of Bert & MF

Part 3.B:: Concatinating I Matrix

```
1 concatenated_I = np.concatenate((I, bert_I), axis=-1)
2 pca = PCA(n_components=64)
3 pca.fit(concatinated_I)
4 pca_concatinated_I = pca.transform(concatinated_I)
5 R3 = U@pca_concatinated_I.T
```

✓ 0.4s

[1] we concatenated the matrix factorization item matrix with the BERT representation which resulted in 832-dimentional Item matrix.

[2:5] we used the PCA to reduce the item matrix dimension to the 64-dimensional space.

Evaluation

```
1 I0,I1 = np.where(R > 0)
2 ground_truth, bert_method, concatinete_method = [],[],[]
3
4 for I,r1 in enumerate(np.argsort(-R1,axis=-1)):
5     ground_truth += [r1[np.in1d(r1, I1[np.where(I0==I)]), invert=True]][:20]]
6 for I,r2 in enumerate(np.argsort(-R2,axis=-1)):
7     bert_method += [r2[np.in1d(r2, I1[np.where(I0==I)]), invert=True]][:20]]
8 for I,r3 in enumerate(np.argsort(-R3,axis=-1)):
9     concatinete_method += [r3[np.in1d(r3, I1[np.where(I0==I)]), invert=True]][:20]]
10
11 ground_truth = np.array(ground_truth)
12 bert_method = np.array(bert_method)
13 concatinete_method = np.array(concatinete_method)
```

✓ 0.8s

Given the test dataset, we can't find out the user ratings for the 20 items which is required to calculate the value of the mentioned evaluation metrics inside the homework, because there are different number of ratings of each user for the items which are fewer than required 20 items.

Important Not Valid but Necessary Assumption: I used the Rating Matrix resulted from the Matrix Factorization section as the ground truth. After dropping the locations filled using the training values, I sorted each user rating for the items and kept only the 20 top rated ones. We are going to use these lists as the ground truth lists to calculate each of the evaluation metrics in the homework.

[1] location of the training data inside training set

[2] three matrices used to store the rating matrices of the Matrix Factorization part, Pure BERT representation part and concatenated BERT and MF part

[4:9] calculating lists of top 20 rating items for each user given each of the methods.

[11:13] converting the lists to the NumPy arrays

Evaluation Methods

Evaluation

```
1 def Recall(y_true,y_pred,k): return np.intersect1d(y_true, y_pred).shape[0]/k
2 def DCG(X): return X/(np.log2(1+np.arange(X.shape[0]))+np.eye(1,X.shape[0]))
3 def NDCG(y_true,y_pred,k): return np.sum((DCG(y_pred)/DCG(y_true))[:k])
4 def Rank_Correlation(y_true,y_pred): return spearmanr(y_true, y_pred)
5 def evaluate(y_true, y_pred, label, y_true_matrix, y_pred_matrix):
6     _R,_N,_C = [],[],[]
7     for ind,y in enumerate(zip(y_true,y_pred)):
8         y_t,y_p = y[0],y[1]
9         _R += [Recall(y_t,y_p,20)]
10        _N += [NDCG(y_true_matrix[ind,y_t],y_pred_matrix[ind,y_p],20)]
11        _C += [Rank_Correlation(y_t,y_p)]
12    _R = np.array(_R).mean()
13    _N = np.array(_N).mean()
14    _C = np.array(_C).mean()
15    print(f"Evaluation Results Of {label}:")
16    print(f"[Recall]           : {_R:0.2f}")
17    print(f"[NDCG]             : {_N:0.2f}")
18    print(f"[Rank_Correlation]: {_C:0.2f}")
```

Each of the mentioned method inside the homework is implemented in this section as a function that receives two sorted sets and returns the evaluation metric.

[1] Recall metric: will simply assumes whole True set is the top 20 rating of the matrix factorization part and counts the intersection of the provided list with this list as the True Positive value, then it will normalize this value by the size of the True set (which is 20) and return the Recall value

[2:3] NDCG metric: it required two sorted sets. We used the Matrix Factorization rating matrix as the ground truth list, after calculating the DCG value for the true list and predicted list, the method will return the sum of normalized DCG up to k'th element (which is 20 as told in the homework)

[4] Rank Correlation Method: we used the implementation of the Spearman's correlation provided in the SciPy library. Given presented mismatch in the two list, the result of this list is not reliable to begin with.

BERT Model Evaluation

BERT Method Result

```
1 evaluate(ground_truth,bert_method,"Pure BERT Representation", R1, R2)
✓ 0.7s
```

Evaluation Results Of Pure BERT Representation:

[Recall] : 0.03
[NDCG] : 23.00
[Rank_Correlation]: 0.25

```
1 evaluate(ground_truth,concatinate_method,"Concatinated Representation", R1, R3)
✓ 0.8s
```

Evaluation Results Of Concatinated Representation:

[Recall] : 0.03
[NDCG] : 23.57
[Rank_Correlation]: 0.26

Assuming the Matrix Factorization ratings matrix as the ground truth, the metric results are shown above. It is not a reliable result though, because incorporating the MF ratings with the BERT model's representation in the concatenation part will definitely injects more information from the MF part inside the rating matrix and hence make the representation more like the MF part which is CONSIDERED to be ground truth but not the real one, therefor we can't say improvement in these metrics will make the concatenation method better.