

QUORA QUESTION PAIR SIMILARITIES USING NBC, LSTM AND BERT

NORTHEASTERN UNIVERSITY – MS IN INFORMATION SYSTEMS

INFO 7610 – NATURAL LANGUAGE ENGINEERING METHODS AND TOOLS

AKSHAY SURESH BHOSALE (001001091)

QUORA QUESTION PAIR SIMILARITIES

1. NBC
 1. LEMMATIZATION, PREPROCESSING,
 2. VECTORIZATION
 3. MODELING
2. LSTM-SCRATCH
 1. PREPROCESSING
 2. TOKENIZATION
 3. PADDING
 4. NEURAL NETWORK MODELING
3. LSTM USING GOOGLE NEWS VECTORS
 1. INDEXING VECTORS
 2. PREPROCESSING
 3. TOKENIZATION
 4. PADDING & EMBEDDING
 5. NEURAL NETWORK MODELING
4. BERT
 1. PREPROCESSING
 2. SEMANTIC DATA GENERATOR
 3. DISTRIBUTION STRATEGY SCOPE
 4. MODELING

NAÏVE BAYES CLASSIFIER

1. IMPORT DATASET
2. EXPLORE DATASET
 1. DROP NULL VALUES
3. PREPROCESSING THE TEXT
 1. LEMMATIZE
 2. REMOVE STOP WORDS
 3. CONCAT QUESTIONS1 AND QUESTIONS2
4. VECTORIZING USING TF-IDF
5. SPLIT DATASET INTO TRAIN-TEST
6. MODEL TRAINING
7. FIT MODEL ON TEST SET

NAÏVE BAYES CLASSIFIER

Lemmatize and Removing Stop Words

```
#preprocessing
def preprocess(series):
    #remove characters other than alphabets & numerics
    words = re.sub("[^A-Za-z0-9]", " ",series).lower().split()

    #lemmatize words
    lemm = WordNetLemmatizer()
    stpwords = stopwords.words('english')
    lemmitized = [lemm.lemmatize(word) for word in words if word not in stpwords]
    sent = ' '.join(lemmitized)
    return sent
```

Combined Question 1 & Question 2

	id	qid1	qid2	question1	question2	is_duplicate	combine
0	0	1	2	step step guide invest share market india	step step guide invest share market	0	step step guide invest share market india step...
1	1	3	4	story kohinoor koh noor diamond	would happen indian government stole kohinoor ...	0	story kohinoor koh noor diamond would happen i...
2	2	5	6	increase speed internet connection using vpn	internet speed increased hacking dns	0	increase speed internet connection using vpn i...
3	3	7	8	mentally lonely solve	find remainder math 23 24 math divided 24 23	0	mentally lonely solve find remainder math 23 2...
4	4	9	10	one dissolve water quickly sugar salt methane c...	fish would survive salt water	0	one dissolve water quickly sugar salt methane c...
5	5	11	12	astrology capricorn sun cap moon cap rising say	triple capricorn sun moon ascendant capricorn say	1	astrology capricorn sun cap moon cap rising sa...
6	6	13	14	buy tiago	keep children active far phone video game	0	buy tiago keep children active far phone video...

NAÏVE BAYES CLASSIFIER

Vectorization using TF-IDF

```
cv = TfidfVectorizer(max_features=50000)#Word to Vectors using Tf-Idf
#Take combine questions data as X
x = cv.fit_transform(train['combine'])
y = np.array(train['is_duplicate'])
print(x.shape)

#Train-Test split
X_train, X_test, y_train, y_test = train_test_split(x,y,test_size=0.05)
print(X_train.shape,X_test.shape)

(404287, 50000)
(384072, 50000) (20215, 50000)
```

Naive Bayes Modeling

```
naive_model = MultinomialNB()#Training
naive_model.fit(X_train,y_train)

#Predictions
y_pred_train = naive_model.predict(X_train)
y_pred_test = naive_model.predict(X_test)
```

Fitting the model on testing set

```
accuracy_train = sum((y_pred_train == y_train).astype(int))/len(y_train)
accuracy_test = sum((y_pred_test == y_test).astype(int))/len(y_test)
print(accuracy_train,accuracy_test)

0.7519241183944677 0.7406381399950531
```

Accuracy on Test Set = 74.06%

LSTM (LONG SHORT-TERM MEMORY): SCRATCH

1. EXPLORE DATASET
 1. SHAPE, SIZE AND DUPLICATES
 2. FILL MISSING VALUES USING 'FILLNA()' METHOD
2. DEFINE FUNCTION FOR TEXT PREPROCESSING:
 1. REMOVE PUNCTUATIONS
 2. REMOVE STOPWORDS
 3. USE REGULAR EXPRESSIONS LIBRARY TO SORT THROUGH MISSPELLED WORDS
 4. STEMMING OF TEXT
3. TEXT TO TOKENS OR TOKENIZATION OF TEXT AND REPLACE THEM
4. ADD PADDING SEQUENCES USING KERAS PREDEFINED METHODS
5. NEURAL NETWORK MODELING
 1. DEFINE LSTM MODEL
 2. TRAIN MODEL
 3. CHECK VALIDATION ACCURACY FOR OPTIMUM RESULT
6. PREDICT ON TEST SET

LSTM (LONG SHORT-TERM MEMORY)

[contd.]

- PRE-DEFINED VARIABLES USED:
 - MAXIMUM TEXT LENGTH = 50
 - MAXIMUM TOKEN LENGTH = $\text{MAX}[\text{MAX}(\text{TRAINING TOKENS}), \text{MAX}(\text{TESTING TOKENS})]$ USING 'NP.MAX'
 - BATCH SIZE FOR NEURAL NETWORK TRAINING = 128
 - EPOCHS = 16 (RANDOMLY SELECTED)
 - TRAIN SIZE : VALIDATION SIZE = 82 : 18

LSTM (Model Architecture & Layers)

[contd.]

- SEQUENTIAL MODEL
- EMBEDDING LAYER DIMENSIONS : MAX_TOKEN +1000 , 32
- DROPOUT AT 0.3
- LSTM LAYER WITH SHAPE NONE, 32
- DROPOUT LAYER 2 AT 0.3
- DENSE LAYER USING SIGMOID ACTIVATION FUNCTION
- LOSS FUNCTION = BINARY CROSS ENTROPY
- OPTIMIZER = 'RMSPROP'

Model: "sequential_1"

Layer (type)	Output Shape	Param #
<hr/>		
embedding_1 (Embedding)	(None, None, 32)	2930560
dropout_2 (Dropout)	(None, None, 32)	0
lstm_1 (LSTM)	(None, 32)	8320
dropout_3 (Dropout)	(None, 32)	0
dense_1 (Dense)	(None, 1)	33
<hr/>		

Total params: 2,938,913

Trainable params: 2,938,913

Non-trainable params: 0

LSTM (RESULTS)

[contd.]

```
loss: 0.4742 - val_accuracy: 0.7864
Epoch 11/16
2594/2594 [=====] - 52s 20ms/step - loss: 0.3656 - accuracy: 0.8368 - val_
loss: 0.4666 - val_accuracy: 0.7830
Epoch 12/16
2594/2594 [=====] - 51s 20ms/step - loss: 0.3639 - accuracy: 0.8367 - val_
loss: 0.4748 - val_accuracy: 0.7777
Epoch 13/16
2594/2594 [=====] - 52s 20ms/step - loss: 0.3634 - accuracy: 0.8370 - val_
loss: 0.4621 - val_accuracy: 0.7877
Epoch 14/16
2594/2594 [=====] - 52s 20ms/step - loss: 0.3607 - accuracy: 0.8389 - val_
loss: 0.4640 - val_accuracy: 0.7873
Epoch 15/16
2594/2594 [=====] - 51s 20ms/step - loss: 0.3595 - accuracy: 0.8395 - val_
loss: 0.4679 - val_accuracy: 0.7855
Epoch 16/16
2594/2594 [=====] - 52s 20ms/step - loss: 0.3579 - accuracy: 0.8395 - val_
loss: 0.4681 - val_accuracy: 0.7867
```

In [39]: `results.head()`

Out[39]:

	test_id	is_duplicate
0	0	0.221046
1	1	0.327794
2	2	0.859311
3	3	0.069816
4	4	0.367610

Accuracy on Validation Set
78.67%

LSTM USING GOOGLE NEWS VECTORS EMBEDDINGS

1. PREDEFINE PARAMETERS
2. INDEX WORD VECTORS
3. TEXT PREPROCESSING
 1. REMOVE STOPWORDS
 2. CONVERT TEXT TO LOWER CASE
 3. REMOVE PUNCTUATIONS
 4. STEMMING TEXTS
4. TOKENIZATION: CONVERT TEXT TO SEQUENCES
5. PAD SEQUENCES TO DETERMINE SHAPE TO BUILD NEURAL NETWORK TENSOR
6. PREPARE EMBEDDING MATRIX
7. SPLIT TRAINING SET INTO TRAIN AND VALIDATION SET
8. DEFINE NEURAL NETWORK
 1. SET UP EMBEDDING LAYER
 2. SET UP LSTM LAYER
 3. SET UP SEQUENTIAL LAYERS
 4. MERGE LAYERS
 5. CHECK WHETHER TO RE-WEIGHT CLASSES TO FIT 17.5% SHARE IN TEST SET
 6. COMPILE AND TRAIN THE MODEL
 7. CHECK VALIDATION ACCURACY
9. PREDICT ON TEST SET

LSTM USING GOOGLE NEWS VECTORS EMBEDDINGS

[contd.]

- PRE-DEFINED VARIABLES USED:
 - MAXIMUM TEXT LENGTH = 30
 - EMBEDDING DIMENSION = 300
 - BATCH SIZE FOR NEURAL NETWORK TRAINING = 2048
 - EPOCHS = 27 (RANDOMLY SELECTED)
 - RANDOM INPUT TO LSTM LAYER BETWEEN 175 - 275
 - RANDOM INPUT TO DENSE LAYER BETWEEN 100 - 150
 - TRAIN SIZE : VALIDATION SIZE = 90 : 10

LSTM using Google News Vectors Embeddings (Model Architecture & Layers)

[contd.]

- SEQUENTIAL MODEL: 2 INPUT LAYERS
- EMBEDDING LAYER DIMENSIONS : 30, 300
- LSTM: WITH SHAPE NONE, 188
- CONCATENATE FOLLOWED BY DROPOUT LAYER
- BATCH NORMALIZATION FOLLOWED BY DENSE LAYER WITH SHAPE NONE, 119
- SECOND DROPOUT LAYER
- BATCH NORMALIZATION LAYER
- SECOND DENSE LAYER
- LOSS FUNCTION = BINARY CROSS ENTROPY
- OPTIMIZER = 'N-ADAM'

Layer (type)	Output Shape	Param #	Connected to
input_1 (InputLayer)	[None, 30]	0	[]
input_2 (InputLayer)	[None, 30]	0	[]
embedding (Embedding)	(None, 30, 300)	36150000	['input_1[0][0]', 'input_2[0][0]']
lstm (LSTM)	(None, 188)	367728	['embedding[0][0]', 'embedding[1][0]']
concatenate (Concatenate)	(None, 376)	0	['lstm[0][0]', 'lstm[1][0]']
dropout (Dropout)	(None, 376)	0	['concatenate[0][0]']
batch_normalization (BatchNormalization)	(None, 376)	1504	['dropout[0][0]']
dense (Dense)	(None, 119)	44863	['batch_normalization[0][0]']
dropout_1 (Dropout)	(None, 119)	0	['dense[0][0]']
batch_normalization_1 (BatchNormalization)	(None, 119)	476	['dropout_1[0][0]']
dense_1 (Dense)	(None, 1)	120	['batch_normalization_1[0][0]']

LSTM (RESULTS)

[contd.]

```
356/356 [=====] - 886s 2s/step - loss: 0.2480 - acc: 0.8120 - val_loss: 0.2705 - val_acc: 0.8041
Epoch 20/27
356/356 [=====] - ETA: 0s - loss: 0.2447 - acc: 0.8153WARNING:tensorflow:`evaluate()` received a value
for `sample_weight`, but `weighted_metrics` were not provided. Did you mean to pass metrics to `weighted_metrics` in `compile
()`? If this is intentional you can pass `weighted_metrics=[]` to `compile()` in order to silence this warning.
356/356 [=====] - 889s 2s/step - loss: 0.2447 - acc: 0.8153 - val_loss: 0.2654 - val_acc: 0.8055
Epoch 21/27
356/356 [=====] - ETA: 0s - loss: 0.2421 - acc: 0.8177WARNING:tensorflow:`evaluate()` received a value
for `sample_weight`, but `weighted_metrics` were not provided. Did you mean to pass metrics to `weighted_metrics` in `compile
()`? If this is intentional you can pass `weighted_metrics=[]` to `compile()` in order to silence this warning.
356/356 [=====] - 890s 2s/step - loss: 0.2421 - acc: 0.8177 - val_loss: 0.2661 - val_acc: 0.8105
Epoch 22/27
356/356 [=====] - ETA: 0s - loss: 0.2394 - acc: 0.8203WARNING:tensorflow:`evaluate()` received a value
for `sample_weight`, but `weighted_metrics` were not provided. Did you mean to pass metrics to `weighted_metrics` in `compile
()`? If this is intentional you can pass `weighted_metrics=[]` to `compile()` in order to silence this warning.
356/356 [=====] - 883s 2s/step - loss: 0.2394 - acc: 0.8203 - val_loss: 0.2675 - val_acc: 0.8099
Epoch 23/27
356/356 [=====] - ETA: 0s - loss: 0.2372 - acc: 0.8220WARNING:tensorflow:`evaluate()` received a value
for `sample_weight`, but `weighted_metrics` were not provided. Did you mean to pass metrics to `weighted_metrics` in `compile
()`? If this is intentional you can pass `weighted_metrics=[]` to `compile()` in order to silence this warning.
356/356 [=====] - 880s 2s/step - loss: 0.2372 - acc: 0.8220 - val_loss: 0.2714 - val_acc: 0.8161
```

In [19]:	results.head()	
Out[19]:	test_id	is_duplicate
	0	0.001727
	1	0.106776
	2	0.403705
	3	0.002177
	4	0.515709

Accuracy on Validation Set : 81.61%

BERT – TRANSFER LEARNING USING TRANSFORMERS TOKENIZERS

1. PREDEFINE PARAMETERS
2. INDEX WORD VECTORS
3. DROP NULL VALUES
4. SPLIT DATASET INTO TRAIN, TEST AND VALIDATION SET
5. CHECK TARGET DISTRIBUTION
6. DEFINE PREPROCESSING AND BERT SEMANTIC DATA GENERATOR FUNCTION
 1. LOAD BERT TOKENIZER (BERT-BASE-UNCASED)
 2. DEFINE NUMBER OF BATCHES PER EPOCH
 3. ENCODE BOTH QUESTIONS TOGETHER SEPARATED BY SEP TOKEN
 4. CONVERT ENCODED FEATURE TO NUMPY ARRAY IN BATCHES
 5. SHUFFLE INDEXES
7. CREATE MODEL UNDER DISTRIBUTION STRATEGY SCOPE
8. LOAD PRETRAINED BERT MODEL
 1. SET UP ENCODED TOKEN IDS FROM BERT TOKENIZER
 2. SET UP ATTENTION MASKS
 3. SET UP TOKEN TYPE IDS
 4. LOADING PRETRAINED BERT MODEL
 5. FREEZE THE BERT MODEL TO REUSE THE PRETRAINED FEATURES WITHOUT MODIFYING THEM.
 6. ADD TRAINABLE LAYERS
 7. COMPILE & TRAIN MODEL
9. EVALUATE ON TEST SET

BERT – TRANSFER LEARNING USING TRANSFORMERS TOKENIZERS

[contd.]

- PRE-DEFINED VARIABLES USED:
 - MAXIMUM TEXT LENGTH = 50
 - EMBEDDING DIMENSION = 300
 - BATCH SIZE FOR NEURAL NETWORK TRAINING = 168
 - EPOCHS = 2 (SELECTED ON THE BASIS OF COMPUTING RESOURCE AND TIME)
 - TRAIN SIZE : TEST SIZE : VALIDATION SIZE = 70 : 15 : 15

BERT – Transfer Learning using Transformers Tokenizers

[contd.]

- INPUT LAYER: NONE, 50
- ATTENTION MASK : NONE, 50
- TOKEN TYPE IDs: NONE, 50
- TF_BERT_MODEL: NONE, 50, 768
- IDIRECTIONAL LSTM: NONE, 50, 256
- GLOBAL AVERAGE POOLING: NONE, 256
- GLOBAL MAX POOLING: NONE, 256
- CONCATENATE: NONE, 512
- DROPOUT: NONE, 512
- DENSE: NONE, 2
- LOSS FUNCTION = CATEGORICAL CROSS ENTROPY
- OPTIMIZER = 'ADAM'

Model: "model"				
Layer (type)	Output Shape	Param #	Connected to	
input_ids (InputLayer)	[(None, 50)]	0	[]	
attention_masks (InputLayer)	[(None, 50)]	0	[]	
token_type_ids (InputLayer)	[(None, 50)]	0	[]	
tf_bert_model (TFBertModel)	TFBaseModelOutputWithPoolingAndCrossAttentions(last_hidden_state=(None, 50, 768), pooler_output=(None, 768), past_key_values=None, hidden_states=None, attentions=None, cross_attentions=None)	109482240	['input_ids[0][0]', 'attention_masks[0][0]', 'token_type_ids[0][0]']	
bidirectional (Bidirectional)	(None, 50, 256)	918528	['tf_bert_model[0][0]']	
global_average_pooling1d (GlobalAveragePooling1D)	(None, 256)	0	['bidirectional[0][0]']	
global_max_pooling1d (GlobalMaxPooling1D)	(None, 256)	0	['bidirectional[0][0]']	
concatenate (Concatenate)	(None, 512)	0	['global_average_pooling1d[0][0]', 'global_max_pooling1d[0][0]']	
dropout_37 (Dropout)	(None, 512)	0	['concatenate[0][0]']	
dense (Dense)	(None, 2)	1026	['dropout_37[0][0]']	
Total params: 110,401,794				
Trainable params: 110,401,794				
Non-trainable params: 0				

BERT (RESULTS)

[contd.]

17390/17684 [=====>.] - ETA: 50s - loss: 0.2132 - accuracy: 0.9115

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

17392/17684 [=====>.] - ETA: 49s - loss: 0.2132 - accuracy: 0.9115

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

17393/17684 [=====>.] - ETA: 49s - loss: 0.2132 - accuracy: 0.9115

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

17394/17684 [=====>.] - ETA: 49s - loss: 0.2132 - accuracy: 0.9115

17684/17684 [=====] - 3268s 185ms/step - loss: 0.2133 - accuracy: 0.9115 - val_loss: 0.2426 - val_accuracy: 0.8996

3765/3785 [=====>.] - ETA: 1s - loss: 0.2427 - accuracy: 0.8991

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

3768/3785 [=====>.] - ETA: 1s - loss: 0.2427 - accuracy: 0.8991

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

3769/3785 [=====>.] - ETA: 1s - loss: 0.2426 - accuracy: 0.8992

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

3770/3785 [=====>.] - ETA: 1s - loss: 0.2426 - accuracy: 0.8992

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

3771/3785 [=====>.] - ETA: 0s - loss: 0.2426 - accuracy: 0.8992

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

3772/3785 [=====>.] - ETA: 0s - loss: 0.2426 - accuracy: 0.8992

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

3773/3785 [=====>.] - ETA: 0s - loss: 0.2426 - accuracy: 0.8992

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

3774/3785 [=====>.] - ETA: 0s - loss: 0.2426 - accuracy: 0.8992

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

Be aware, overflowing tokens are not returned for the setting you have chosen, i.e. sequence truncation strategy. So the returned list will always be empty even if some tokens have been

3775/3785 [=====>.] - ETA: 0s - loss: 0.2426 - accuracy: 0.8992

3785/3785 [=====] - 261s 69ms/step - loss: 0.2426 - accuracy: 0.8991

[0.24255122244358063, 0.8991413712501526]

Accuracy on validation Set : 89.91%

REFERENCES:

1. [HTTPS://KERAS.IO/EXAMPLES/NLP/SEMANTIC_SIMILARITY_WITH_BERT/](https://keras.io/examples/nlp/semantic_similarity_with_BERT/)
2. [HTTPS://WWW.KAGGLE.COM/COMPETITIONS/QUORA-QUESTION-PAIRS](https://www.kaggle.com/competitions/quora-question-pairs)
3. [HTTPS://WWW.KAGGLE.COM/DATASETS/LEADBEST/GOOGLENEWSVECTORSNEGATIVE300](https://www.kaggle.com/datasets/leadbetter/google-news-vectors-negative300)
4. [HTTPS://HUGGINGFACE.CO/DATASETS?TASK_IDS=TASK_IDS:SEMANTIC-SIMILARITY-CLASSIFICATION](https://huggingface.co/datasets?task_ids=task_ids:semantic-similarity-classification)
5. [HTTPS://HUGGINGFACE.CO/BERT-BASE-UNCASED](https://huggingface.co/bert-base-uncased)
6. [HTTPS://HUGGINGFACE.CO/TFTRANSFORMERS/BERT-BASE-UNCASED](https://huggingface.co/tftransformers/bert-base-uncased)

THANK YOU SO MUCH