

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×

You have **2** free stories left this month. Sign up and get an extra one for free.

# BERT: Multilabel Text Classification



Zuzanna Deutschman

Follow

Feb 9 · 3 min read ★

## Introduction

In my previous article, I introduced various machine learning methods that enable assigning a set of relevant genres for a single movie description (please visit the article for dataset). The best F1 score = 0.43 was obtained for Classifier Chain model. My idea to be verified is to train neural network with BERT embeddings.

BE To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#),  
re including cookie policy. ×

It's main innovation is applying the bidirectional training of Transformer, a popular attention model, to language modelling. This results in a deeper sense of language context and flow than single-direction language models.

## Code

Bert\_serving enables using BERT model as a sentence encoding service for mapping a variable-length sentence to a fixed-length.

To find the best bunch of parameters I used sacred module. Sacred is a tool to help you configure, organize, log and reproduce experiments in order to:

- keep track of all the parameters of your experiment
- easily run your experiment for different settings
- save configurations for individual runs in a database

```
1 from bert_serving.server import BertServer
2 from bert_serving.server.helper import get_args_parser
3 from bert_serving.server.helper import get_shutdown_parser
```

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×

```
7  from sklearn.preprocessing import MultiLabelBinarizer
8  import json
9  import os
10 import tensorflow as tf
11 from keras import backend as K
12 import collections
13 import numpy as np
14 from sacred import Experiment
15 from sacred.observers import MongoObserver
16 ex = Experiment()
17 ex.observers.append(MongoObserver(
18     #url='mongodb://mongo_user:mongo_password@localhost:27017/?authMechanism=SCRAM-SHA-
19
20 base_path = ''
21
22
23 def prepare_data():
24     trainfilename = 'movies_train.csv'
25     evalfilename = 'movies_eval.csv'
26
27     if os.path.exists(trainfilename) and os.path.exists(evalfilename):
28         train_df = pd.read_pickle(trainfilename)
29         eval_df = pd.read_pickle(evalfilename)
30
31         return train_df, eval_df
32
33     #data preprocessing
34     meta = pd.read_csv(os.path.join(
35         base_path, "movie.metadata.tsv"), sep='\t', header=None)
36     meta.columns = ["movie_id". 1. "movie name". 3. 4. 5. 6. 7. "genre"]
```

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×

```
39 plots = pd.read_csv(os.path.join(
40     base_path, "plot_summaries.txt"), sep='\t', header=None)
41 plots.columns = ["movie_id", "plot"]
42
43 genres['movie_id'] = genres['movie_id'].astype(str)
44 plots['movie_id'] = plots['movie_id'].astype(str)
45 movies = pd.merge(plots, genres, on='movie_id')
46
47 genres_lists = []
48 for i in movies['genre']:
49     genres_lists.append(list(json.loads(i).values()))
50 movies['genre'] = genres_lists
51 multilabel_binarizer = MultiLabelBinarizer()
52 multilabel_binarizer.fit_transform(movies['genre'])
53 # transform target variable
54 y = multilabel_binarizer.transform(movies['genre'])
55 for idx, genre in enumerate(multilabel_binarizer.classes_):
56     movies[genre] = y[:, idx]
57 movies.to_csv('movies.csv')
58 movies_new = pd.read_csv('movies.csv')
59 movies = movies_new
60 movies_columns = movies.columns
61
62 del movies['Unnamed: 0']
63 del movies['movie_name']
64 del movies['genre']
65
66 df = pd.DataFrame()
67 df['id'] = movies['movie_id']
68 df['labels'] = list(map(list, zip(*[movies[col] for col in movies
```

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy.

```
72     TRAIN_VAL_RATIO = 0.8
73     LEN = df.shape[0]
74     SIZE_TRAIN = int(TRAIN_VAL_RATIO*LEN)
75
76     train_df = df[:SIZE_TRAIN].drop(labels='id', axis=1)
77     eval_df = df[SIZE_TRAIN:]
78
79     train_df.to_pickle(trainfilename)
80     eval_df.to_pickle(evalfilename)
81
82     return train_df, eval_df
83
84
85
86 def encode_with_bert(train_df, eval_df, max_seq_len = 50):
87     #bert-serving-start -model_dir L-12_H-768_A-12/ -num_worker=2
88     filename = f'max_seq_len_{max_seq_len}'
89     trainfile = f'{filename}_train.npy'
90     testfile = f'{filename}_test.npy'
91     if not os.path.exists(trainfile) or not os.path.exists(testfile):
92         args = get_args_parser().parse_args(['-model_dir', 'uncased_L-12_H-768_A-12/',
93                                             '-max_seq_len', f'{max_seq_len}',
94                                             '-port', '5555',
95                                             '-fp16', '-xla',
96                                             '-num_worker', '1'])
97
98         server = BertServer(args)
99         server.start()
100
101         bc = BertClient()
102         x_train = bc.encode(list(train_df.text))
```

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×

```

105         ['-ip', 'localhost', '-port', '5555', '-timeout', '5000'])
106         BertServer.shutdown(shut_args)
107
108         np.save(trainfile, x_train)
109         np.save(testfile, x_test)
110     else:
111         x_train = np.load(trainfile)
112         x_test = np.load(testfile)
113
114     return x_train, x_test
115
116
117 def focal_loss(gamma=2., alpha=.25):
118     def focal_loss_fixed(y_true, y_pred):
119         pt_1 = tf.where(tf.equal(y_true, 1), y_pred, tf.ones_like(y_pred))
120         pt_0 = tf.where(tf.equal(y_true, 0), y_pred, tf.zeros_like(y_pred))
121         return -K.mean(alpha * K.pow(1. - pt_1, gamma) * K.log(pt_1)) - K.mean((1 - alpha)
122         * K.pow(pt_0, gamma) * K.log(pt_0))
123     return focal_loss_fixed
124
125 class MultiLabelClassifier:
126     def __init__(self):
127         self.model = tf.keras.Sequential()
128         self.model.add(tf.keras.Input(shape=(768,)))
129         self.model.add(tf.keras.layers.Dense(512, activation='relu'))
130         self.model.add(tf.keras.layers.Dropout(rate=0.8))
131         self.model.add(tf.keras.layers.Dense(363, activation='sigmoid'))
132         self.model.summary()
133
134     def train(self, x_train, y_train, x_test, y_test, optimizer='adam', loss=[focal_loss])

```

1 To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#),  
 1 including cookie policy.

```

138         metrics = [tf.keras.metrics.AUC()]
139         self.model.compile(optimizer, loss, metrics)
140
141         history = self.model.fit(
142             x_train, y_train, batch_size, epochs, validation_data=(x_test, y_test))
143         results = self.model.evaluate(x_test, y_test, batch_size)
144         for name, value in zip(self.model.metrics_names, results):
145             print("%s: %.3f" % (name, value))
146         return history, results
147
148     def predict(self, x_test):
149         return self.model.predict(x_test)
150
151     def test_sample(self, test_data, actual):
152         print("actual ground truth={}, predicted={}".format(
153             actual, self.model.predict(test_data)))
154
155 @ex.config
156 def my_config():
157     max_seq_len = 256
158     batch_size = 128
159     gamma = 2
160
161 @ex.automain
162 def train_and_evaluate(max_seq_len, batch_size, gamma):
163     train_df, eval_df = prepare_data()
164     y_train = np.array(train_df['labels'].tolist())
165     y_test = np.array(eval_df['labels'].tolist())
166
167     x_train, x_test = encode_with_bert(train_df, eval_df, max_seq_len=max_seq_len)
  
```

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×

```
170     history, results = classifier.train(x_train,
171                                         y_train,
172                                         x_test,
173                                         y_test,
174                                         optimizer='adam',
175                                         loss=[focal_loss(alpha=.10, gamma=gamma)],
176                                         metrics='auc',
177                                         batch_size=batch_size,
178                                         epochs=100
179                                         )
180
181     predictions = classifier.model.predict(x_test)
182
183     t = 0.20
184     predicted = []
185     predicted = (predictions >= t).astype(int)
186     #print(predicted)
187     f1 = f1_score(y_test, predicted, average='micro')
188     print("F1 of BERT model is:", f1)
189
190     loss, auc = results
191
192     ex.log_scalar("test_f1", f1)
193     ex.log_scalar("train_loss", loss)
194     ex.log_scalar("train_auc", auc)
195
```

bert\_multilabel\_extract.py hosted with ❤ by GitHub

[view raw](#)



Medium To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×

1. **max\_seq\_len** — maximum length of sequence. The best found value is **256** and it required GPU to be used
2. **batch\_size** — number of samples that will be propagated through the network. Chosen number is **128**
3. **gamma** — the focusing parameter in focal loss that smoothly adjusts the rate at which easy examples are down-weighted. The focal loss is designed to address class imbalance by down-weighting inliers (easy examples) such that their contribution to the total loss is small even if their number is large. It focuses on training a sparse set of hard examples. The most optimal value of gamma in our example is **2**

```
1 from bert_multiclass_torch_extract import ex
2 import itertools
3
4 max_seq_len_values = [256]
5 batch_size_values = [128]
6 gamma_values = [2]
7
8 for max_seq_len, batch_size, gamma in itertools.product(max_seq_len_values, batch_size_values, gamma_values):
9     ex.run(config_updates={'max_seq_len': max_seq_len, 'batch_size': batch_size, 'gamma': gamma})
```

run\_sacred.py hosted with ❤ by GitHub

[view raw](#)

Ol To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×

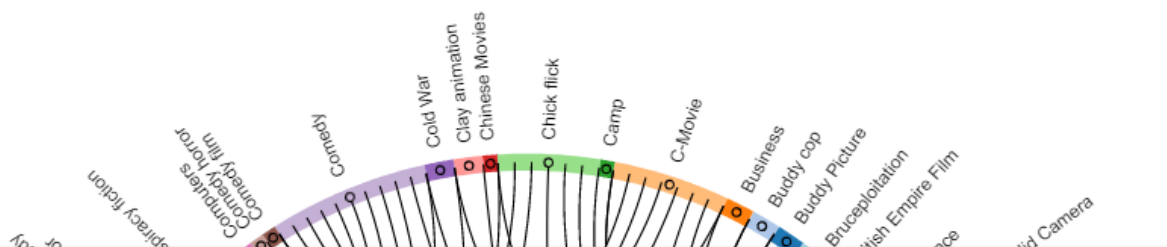
## Labels co-occurrences

We need to remember that there are multiple labels in our dataset and sometimes one label indicate occurrence of the other one. In order to check this percentage dependency, I created matrix of co-occurrences. Now we know that, for example:

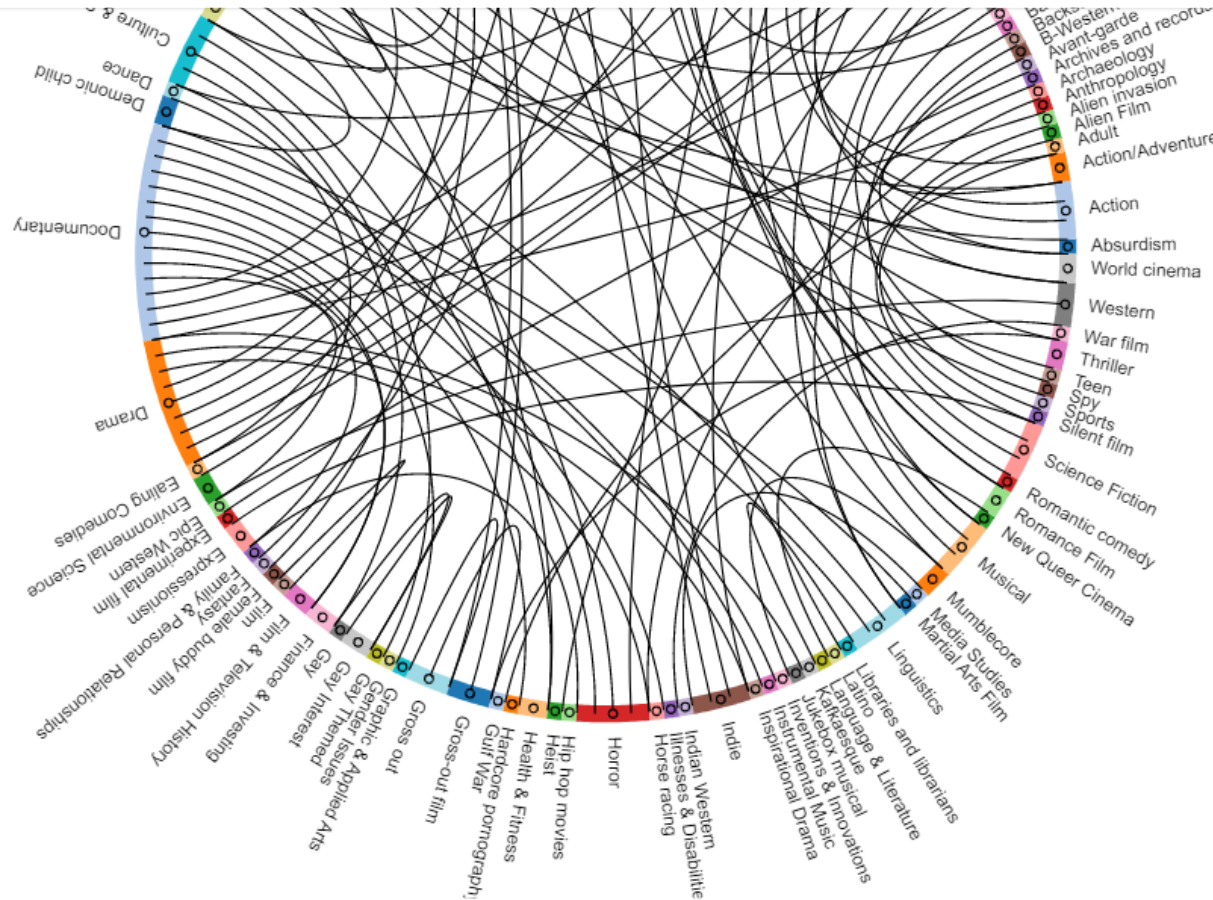
```
Action Comedy -> Action, 0.9722222222222222
Action Comedy -> Comedy, 0.9629629629629629
Action Thrillers -> Action, 0.9492537313432836
Action Thrillers -> Thriller, 0.9253731343283582
Adventure Comedy -> Comedy, 0.9101123595505618
```

The top 100 strongest dependencies represented on graph below.

**Directed Graph**



To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×



Route chord of labels cooccurrences

With gained knowledge, we can try to modify our predictions now. After few trials I decided to change value of resulting label for 1, only if its co-occurrence with indicating label was greater than or equal 0.9.

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×

```
4     meta = pd.read_csv(os.path.join(
5         base_path, "movie.metadata.tsv"), sep='\t', header=None)
6     meta.columns = ["movie_id", 1, "movie_name", 3, 4, 5, 6, 7, "genre"]
7     genres = meta[["movie_id", "movie_name", "genre"]]
8
9     plots = pd.read_csv(os.path.join(
10         base_path, "plot_summaries.txt"), sep='\t', header=None)
11     plots.columns = ["movie_id", "plot"]
12
13     genres['movie_id'] = genres['movie_id'].astype(str)
14     plots['movie_id'] = plots['movie_id'].astype(str)
15     movies = pd.merge(plots, genres, on='movie_id')
16
17     genres_lists = []
18     for i in movies['genre']:
19         genres_lists.append(list(json.loads(i).values()))
20     movies['genre'] = genres_lists
21     multilabel_binarizer = MultiLabelBinarizer()
22     multilabel_binarizer.fit_transform(movies['genre'])
23     # transform target variable
24     y = multilabel_binarizer.transform(movies['genre'])
25     for idx, genre in enumerate(multilabel_binarizer.classes_):
26         movies[genre] = y[:, idx]
27
28     #del movies['Unnamed: 0']
29     del movies['movie_name']
30     del movies['movie_id']
31     del movies['plot']
32     del movies['genre']
33
```

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×

```

36     LEN = df.shape[0]
37     SIZE_TRAIN = int(TRAIN_VAL_RATIO*LEN)
38
39     train_df = df[:SIZE_TRAIN]
40
41     return train_df
42
43
44 def cooccurrences(t):
45     df = prepare_train_data_for_cooccurrences()
46     cooccurrence_matrix = np.dot(df.transpose(),df)
47     cooccurrence_matrix_diagonal = np.diagonal(cooccurrence_matrix)
48     with np.errstate(divide='ignore', invalid='ignore'):
49         cooccurrence_matrix_percentage = np.nan_to_num(np.true_divide(cooccurrence_matrix,
50         tuples = {}
51         result = np.where(cooccurrence_matrix_percentage > t)
52         listOfCoordinates= list(zip(result[0], result[1]))
53         for cord in listOfCoordinates:
54             if cord[0] != cord[1]:
55                 tuples[(df.columns.tolist()[cord[0]],
56                        df.columns.tolist()[cord[1]])] = cooccurrence_matrix_percentage
57
58 movies_columns = get_columns_names()
59 pred_df = pd.DataFrame(predicted, columns = movies_columns)
60 coocs = cooccurrences(t=0.9)
61 ind = [i[0] for i in coocs]
62 conc = [i[1] for i in coocs]
63
64 for i, label in enumerate(ind):
65     for index, row in pred_df.iterrows():
66         if row[label] == 1:

```

```
6 To make Medium work, we log user data. By using Medium, you agree to our Privacy Policy,  
6 including cookie policy.  
69 f1_updated = f1_score(y_test, pred_df.to_numpy(), average='micro')  
70 print("Updated F1 of BERT model is:", f1_updated)
```

cooccurrences.py hosted with ❤ by GitHub [view raw](#)

Updated F1 score is **0.5**, which is small improvement.

## Results

Sample results look as follows.

First Platoon centers around Rock Brannigan ([[Scott Gibson and his ragtag squad of ex-military zombie hunters trying to make a living in the desert Southwest two years after the zombie apocalypse. Along the way they encounter the grizzled Pa Jericho , and the eccentric Rex Necro .

**Action, Comedy, Horror, Parody, Science Fiction**

The life of the S&M-theme artist and author Seiu Ito is depicted in the film. His artistic life and Sadian philosophy, inspired by his torturing of his two wives and Tae, his favorite prostitute, are portrayed as shown in his journalistic writings. Tae is eventually driven insane due to Ito's attentions.

**Drama, Japanese Movies, World cinema**

**C** To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×

Usage of BERT embeddings enabled to gain 7.7% F1 score improvement, (50% overall).

In future work, I think that good idea would be to reduce number of labels keep only the main ones. As 'Action Comedy' is in 100% of cases 'Action' and in 100% 'Comedy', maybe we don't really need this category.

Thank you for reading.

[Deep Learning](#)[Data Science](#)[NLP](#)[Bert](#)[Machine Learning](#)

## Discover Medium

Welcome to a place where words matter. On Medium, smart voices and original ideas take center stage - with no ads in sight. Watch

## Make Medium yours

Follow all the topics you care about, and we'll deliver the best stories for you to your homepage and inbox. Explore

## Become a member

Get unlimited access to the best stories on Medium — and support writers while you're at it. Just \$5/month. Upgrade

[About](#)[Help](#)[Legal](#)

To make Medium work, we log user data. By using Medium, you agree to our [Privacy Policy](#), including cookie policy. ×