

Movie Genre Text Classification

This notebook uses a Kaggle dataset containing pieces of about 22,500 scripts. These script chunks will be used to predict the main movie genre for the given script.

Data set can be found at : https://www.kaggle.com/datasets/lykin22/movie-genre-data?select=kaggle_movie_train.csv

Code inspired by Karen Mazidi's GitHub notebooks:

https://github.com/kjmazidi/NLP/blob/master/Part_6-Deep%20Learning/Chapter_23_Keras/Keras_imdb_1_Dense_Sequential.ipynb

https://github.com/kjmazidi/NLP/blob/master/Part_6-Deep%20Learning/Chapter_24_DL_variations/Keras_imdb_2_RNN.ipynb

https://github.com/kjmazidi/NLP/blob/master/Part_6-Deep%20Learning/Chapter_25_Embeddings/Embedding%20Layer.ipynb

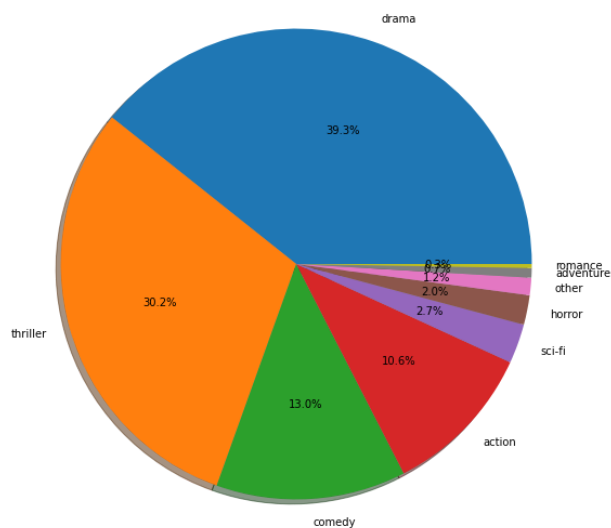
Import Data

```
1 from google.colab import drive
2 import pandas as pd
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

	text	genre
0	eady dead, maybe even wishing he was. INT. 2ND...	thriller
1	t, summa cum laude and all. And I'm about to l...	comedy
2	up Come, I have a surprise.... She takes him ...	drama
3	ded by the two detectives. INT. JEFF'S APARTME...	thriller
4	nd dismounts, just as the other children reach...	drama

Display Target Distribution

```
1 import matplotlib.pyplot as plt
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```



Text Preprocessing

There are a lot of potential outcomes for this dataset- many of which aren't representative of much of the data. I want to drop all genres that aren't drama or thriller. I feel like the other genres can fit in these categories anyway.

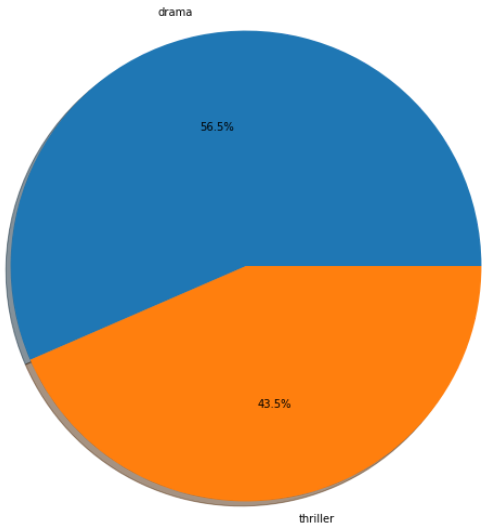
```
1 movies = movies[movies.genre != 'romance']
2 movies = movies[movies.genre != 'adventure']
3 movies = movies[movies.genre != 'other']
4 movies = movies[movies.genre != 'horror']
5 movies = movies[movies.genre != 'sci-fi']
6 movies = movies[movies.genre != 'comedy']
7 movies = movies[movies.genre != 'action']
```

1 movies

	text	genre
0	eady dead, maybe even wishing he was. INT. 2ND...	thriller
2	up Come, I have a surprise.... She takes him ...	drama
3	ded by the two detectives. INT. JEFF'S APARTME...	thriller
4	nd dismounts, just as the other children reach...	drama
5	breadth of the bluff. Gabe pulls out his ancie...	thriller
...
22571	watching us all. SWISH PAN TO INT. TANGIERS E...	drama
22572	HER and TWO COOKS are standing in a row waitin...	thriller
22574	n in the world to decide what I'm going to do ...	drama
22575	shards. BOJO LAZI! Laz pushes Deke back through...	drama
22576	OTTIE You've got a thing about Ernie's, haven'...	thriller

15697 rows × 2 columns

```
1 genre_labels = movies['genre'].value_counts().index.tolist()
2 plt.figure(1, figsize=(20,10))
3 plt.pie(movies['genre'].value_counts(), labels = genre_labels, shadow = True, autopct='%1.1f%%')
4 plt.show()
```



The updated data looks much easier to work with since there are no longer 9 potential target values. The model should be able to classify the chunk of movie script into one of the remaining 4 categories: action, drama, thriller, comedy.

```
1 # convert labels to numeric categorical
2 # dictionaries to track what is what
3 movies['genre_id'] = movies['genre'].factorize()[0]
4 genre_id_df = movies[['genre', 'genre_id']].drop_duplicates().sort_values('genre_id')
5 genre_to_id = dict(genre_id_df.values)
6 id_to_genre = dict(genre_id_df [['genre_id', 'genre']].values)

1 import nltk
2 nltk.download('stopwords')
3 from nltk.corpus import stopwords
4 from nltk.tokenize import word_tokenize
5 #from tensorflow.keras.preprocessing.text import Tokenizer
6 #from tensorflow.keras import preprocessing
7 nltk.download('punkt')
8 stop_words = set(stopwords.words('english'))
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
```

```
1 def clean_text(x):
2     words = [word.lower() for word in word_tokenize(x) if word.isalpha() and not word in stop_words and len(word)<10]
3     words = ' '.join(words)
4     return words

1 for index in movies.index:
2     movies['text'][index]= clean_text(movies['text'][index])

<ipython-input-13-d9c933a3f669>:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
movies['text'][index]= clean_text(movies['text'][index])
```

▼ Split Data

```
1 from sklearn.model_selection import train_test_split

1 X = movies['text'].values
2 y = movies['genre_id'].values

1 train_text,test_text, train_labels, test_labels = train_test_split(X, y, test_size=0.2, random_state=0)
```

▼ Sequential Model

The simple sequential model is being used as a baseline.

```
1 import tensorflow as tf
2 import numpy as np
3 from tensorflow.keras import layers, models
4 from sklearn.feature_extraction.text import CountVectorizer

1 # variables
2 vocab_size = 10000
3 dimensions = vocab_size

1 vectorizer = CountVectorizer()
2 vectorizer.fit(train_text)
3
4 X_train = vectorizer.transform(train_text).toarray()
5 X_test = vectorizer.transform(test_text).toarray()

1 y_train = np.asarray(train_labels).astype('float32')
2 y_test = np.asarray(test_labels).astype('float32')

1 from keras.backend import clear_session

1 input_dim = X_train.shape[1]
2 clear_session()
3 model = models.Sequential()
4 model.add(layers.Dense(16, activation='relu', input_shape=(input_dim,)))
5 model.add(layers.Dense(16, activation='relu'))
6 model.add(layers.Dense(1, activation='sigmoid')) # for binary output

1 # compile
2 model.compile(optimizer='adam', loss='binary_crossentropy',metrics=['accuracy']) # use adam since NOT 1 hot encoded

1 # evaluate
2 history = model.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=20, batch_size=512)

Epoch 1/20
25/25 [=====] - 7s 239ms/step - loss: 0.5661 - accuracy: 0.7077 - val_loss: 0.4007 - val_accuracy: 0.8822
Epoch 2/20
25/25 [=====] - 5s 209ms/step - loss: 0.2405 - accuracy: 0.9627 - val_loss: 0.2062 - val_accuracy: 0.9443
Epoch 3/20
25/25 [=====] - 5s 211ms/step - loss: 0.0826 - accuracy: 0.9931 - val_loss: 0.1416 - val_accuracy: 0.9529
Epoch 4/20
25/25 [=====] - 5s 210ms/step - loss: 0.0333 - accuracy: 0.9988 - val_loss: 0.1245 - val_accuracy: 0.9551
Epoch 5/20
25/25 [=====] - 5s 212ms/step - loss: 0.0163 - accuracy: 0.9998 - val_loss: 0.1167 - val_accuracy: 0.9580
Epoch 6/20
25/25 [=====] - 9s 359ms/step - loss: 0.0094 - accuracy: 1.0000 - val_loss: 0.1138 - val_accuracy: 0.9557
Epoch 7/20
25/25 [=====] - 10s 407ms/step - loss: 0.0061 - accuracy: 1.0000 - val_loss: 0.1130 - val_accuracy: 0.9554
Epoch 8/20
25/25 [=====] - 9s 378ms/step - loss: 0.0042 - accuracy: 1.0000 - val_loss: 0.1135 - val_accuracy: 0.9564
Epoch 9/20
25/25 [=====] - 8s 300ms/step - loss: 0.0031 - accuracy: 1.0000 - val_loss: 0.1141 - val_accuracy: 0.9551
Epoch 10/20
25/25 [=====] - 5s 208ms/step - loss: 0.0024 - accuracy: 1.0000 - val_loss: 0.1143 - val_accuracy: 0.9554
```

▼ RNN

```
1 history = model.fit(train_data, y_train, epochs=20, batch_size=128, validation_split=0.2)
```

4/7

```

79/79 [=====] - 2s 23ms/step - loss: 0.6823 - accuracy: 0.5725 - val_loss: 0.6885 - val_accuracy: 0.5629
Epoch 15/20
79/79 [=====] - 2s 23ms/step - loss: 0.6827 - accuracy: 0.5728 - val_loss: 0.6876 - val_accuracy: 0.5641
Epoch 16/20
79/79 [=====] - 2s 23ms/step - loss: 0.6828 - accuracy: 0.5717 - val_loss: 0.6851 - val_accuracy: 0.5621
Epoch 17/20
79/79 [=====] - 2s 23ms/step - loss: 0.6821 - accuracy: 0.5728 - val_loss: 0.6847 - val_accuracy: 0.5637
Epoch 18/20
79/79 [=====] - 2s 23ms/step - loss: 0.6826 - accuracy: 0.5727 - val_loss: 0.6855 - val_accuracy: 0.5633
Epoch 19/20
79/79 [=====] - 2s 23ms/step - loss: 0.6830 - accuracy: 0.5719 - val_loss: 0.6843 - val_accuracy: 0.5649
Epoch 20/20
79/79 [=====] - 2s 23ms/step - loss: 0.6827 - accuracy: 0.5726 - val_loss: 0.6846 - val_accuracy: 0.5641

```

```

1 from sklearn.metrics import classification_report
2
3 pred = model.predict(test_data)
4 pred = [1.0 if p>= 0.5 else 0.0 for p in pred]
5 print(classification_report(y_test, pred))

99/99 [=====] - 1s 4ms/step
      precision    recall  f1-score   support

      0.0         0.67       0.01       0.03       1391
      1.0         0.56       0.99       0.72       1749

 accuracy         0.56       3140
 macro avg       0.61       0.50       0.37       3140
 weighted avg    0.61       0.56       0.41       3140

```

The results are pretty sad. RNN took quite a long time to train compared to the sequential model and the accuracy is much lower at 56%. This could be due the vanishing gradient problem, so I want to build another model.

I want to retry with LSTM and see if that helps.

```

1 clear_session()
2 model = models.Sequential()
3 model.add(layers.Embedding(max_features, 32))
4 model.add(layers.LSTM(32))
5 model.add(layers.Dense(1, activation='sigmoid')) # sigmoid for binary classification

1 # compile
2 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

1 history = model.fit(train_data, y_train, epochs=10, batch_size=128, validation_split=0.2) # fewer epochs because slower

Epoch 1/10
79/79 [=====] - 7s 53ms/step - loss: 0.6857 - accuracy: 0.5655 - val_loss: 0.6862 - val_accuracy: 0.5621
Epoch 2/10
79/79 [=====] - 4s 45ms/step - loss: 0.6838 - accuracy: 0.5686 - val_loss: 0.6856 - val_accuracy: 0.5621
Epoch 3/10
79/79 [=====] - 6s 70ms/step - loss: 0.6838 - accuracy: 0.5686 - val_loss: 0.6859 - val_accuracy: 0.5621
Epoch 4/10
79/79 [=====] - 6s 72ms/step - loss: 0.6841 - accuracy: 0.5686 - val_loss: 0.6854 - val_accuracy: 0.5621
Epoch 5/10
79/79 [=====] - 4s 46ms/step - loss: 0.6838 - accuracy: 0.5686 - val_loss: 0.6855 - val_accuracy: 0.5621
Epoch 6/10
79/79 [=====] - 4s 45ms/step - loss: 0.6839 - accuracy: 0.5686 - val_loss: 0.6855 - val_accuracy: 0.5621
Epoch 7/10
79/79 [=====] - 4s 45ms/step - loss: 0.6837 - accuracy: 0.5686 - val_loss: 0.6857 - val_accuracy: 0.5621
Epoch 8/10
79/79 [=====] - 4s 46ms/step - loss: 0.6839 - accuracy: 0.5686 - val_loss: 0.6854 - val_accuracy: 0.5621
Epoch 9/10
79/79 [=====] - 4s 45ms/step - loss: 0.6838 - accuracy: 0.5686 - val_loss: 0.6854 - val_accuracy: 0.5621
Epoch 10/10
79/79 [=====] - 4s 45ms/step - loss: 0.6839 - accuracy: 0.5686 - val_loss: 0.6854 - val_accuracy: 0.5621

```

```

1 pred = model.predict(test_data)
2 pred = [1.0 if p>= 0.5 else 0.0 for p in pred]
3 print(classification_report(y_test, pred))

99/99 [=====] - 1s 7ms/step
      precision    recall  f1-score   support

      0.0         0.00       0.00       0.00       1391
      1.0         0.56       1.00       0.72       1749

 accuracy         0.56       3140
 macro avg       0.28       0.50       0.36       3140
 weighted avg    0.31       0.56       0.40       3140

```

```

/usr/local/lib/python3.8/dist-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in lab
_warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.8/dist-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in lab
_warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.8/dist-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in lab
_warn_prf(average, modifier, msg_start, len(result))

```

The results are still extremely unimpressive.

Embedding Trial

```
1 from tensorflow.keras.layers.experimental.preprocessing import TextVectorization
2 from tensorflow import keras

1 embedding_dim = 128
2 max_seq = 200

1 vectorizer = TextVectorization(max_tokens=20000, output_sequence_length=200)
2 text_ds = tf.data.Dataset.from_tensor_slices(train_text).batch(128)
3 vectorizer.adapt(text_ds)

1 voc = vectorizer.get_vocabulary()
2 word_index = dict(zip(voc, range(len(voc))))

1 # set up embedding layer
2 embedding_layer = layers.Embedding(len(word_index) + 1, embedding_dim, input_length=max_seq)

1 int_sequences_input = keras.Input(shape=(None,), dtype="int64")
2 embedded_sequences = embedding_layer(int_sequences_input)
3 x = layers.Conv1D(128, 5, activation="relu")(embedded_sequences)
4 x = layers.MaxPooling1D(5)(x)
5 x = layers.Conv1D(128, 5, activation="relu")(x)
6 x = layers.MaxPooling1D(5)(x)
7 x = layers.Conv1D(128, 5, activation="relu")(x)
8 x = layers.GlobalMaxPooling1D()(x)
9 x = layers.Dense(128, activation="relu")(x)
10 x = layers.Dropout(0.5)(x)
11 preds = layers.Dense(1, activation="sigmoid")(x)
12 model = keras.Model(int_sequences_input, preds)
13 model.summary()
14
```

Model: "model_2"

Layer (type)	Output Shape	Param #
input_4 (InputLayer)	[(None, None)]	0
embedding_1 (Embedding)	(None, None, 128)	2560128
conv1d_9 (Conv1D)	(None, None, 128)	82048
max_pooling1d_6 (MaxPooling 1D)	(None, None, 128)	0
conv1d_10 (Conv1D)	(None, None, 128)	82048
max_pooling1d_7 (MaxPooling 1D)	(None, None, 128)	0
conv1d_11 (Conv1D)	(None, None, 128)	82048
global_max_pooling1d_3 (GlobalMaxPooling1D)	(None, 128)	0
dense_7 (Dense)	(None, 128)	16512
dropout_3 (Dropout)	(None, 128)	0
dense_8 (Dense)	(None, 1)	129

Total params: 2,822,913
Trainable params: 2,822,913
Non-trainable params: 0

```
1 X_train = vectorizer(np.array([s for s in train_text])).numpy()
2 X_test = vectorizer(np.array([s for s in test_text])).numpy()
3
4 y_train = np.array(train_labels)
5 y_test = np.array(test_labels)

1 model.compile(loss="binary_crossentropy", optimizer="adam", metrics=["accuracy"])
2 model.fit(X_train, y_train, batch_size=128, epochs=10, validation_data=(X_test, y_test))

Epoch 1/10
99/99 [=====] - 51s 506ms/step - loss: 0.0590 - accuracy: 0.9806 - val_loss: 0.1301 - val_accuracy: 0.9548
Epoch 2/10
99/99 [=====] - 46s 470ms/step - loss: 0.0087 - accuracy: 0.9981 - val_loss: 0.1788 - val_accuracy: 0.9541
Epoch 3/10
99/99 [=====] - 46s 466ms/step - loss: 7.9227e-04 - accuracy: 0.9998 - val_loss: 0.2363 - val_accuracy: 0.9529
Epoch 4/10
99/99 [=====] - 45s 457ms/step - loss: 9.1496e-05 - accuracy: 1.0000 - val_loss: 0.2453 - val_accuracy: 0.9535
Epoch 5/10
99/99 [=====] - 65s 657ms/step - loss: 4.2321e-05 - accuracy: 1.0000 - val_loss: 0.2626 - val_accuracy: 0.9532
Epoch 6/10
99/99 [=====] - 56s 562ms/step - loss: 2.7528e-05 - accuracy: 1.0000 - val_loss: 0.2768 - val_accuracy: 0.9529
Epoch 7/10
99/99 [=====] - 49s 494ms/step - loss: 1.2811e-05 - accuracy: 1.0000 - val_loss: 0.2927 - val_accuracy: 0.9532
Epoch 8/10
```

https://colab.research.google.com/drive/1TtX72Fweb0_RUBUDSBp2Nkv8_JneB8NV#scrollTo=LzYHloyW2O0s&printMode=true

```
99/99 [=====] - 58s 592ms/step - loss: 1.4771e-05 - accuracy: 1.0000 - val_loss: 0.3053 - val_accuracy: 0.9529
Epoch 9/10
99/99 [=====] - 49s 498ms/step - loss: 7.3442e-06 - accuracy: 1.0000 - val_loss: 0.3139 - val_accuracy: 0.9525
Epoch 10/10
99/99 [=====] - 59s 601ms/step - loss: 5.9494e-06 - accuracy: 1.0000 - val_loss: 0.3210 - val_accuracy: 0.9525
<keras.callbacks.History at 0x7fcc53366520>
```

▼ Analysis

I was initially very excited about my dataset choice. I was excited at the thought of determining genre based upon a script snippet. I played with the original data for a long time before determining that 7 potential categories were too many for me to overcome as a beginner. The highest accuracy I could muster was about 16% after mulitple hours. So I narrowed the potential genres to only drama and thriller. I preprocessed my text to be normalized and remove stop words and some of the nonsense words.

Now that multinomial classificaiton problem was a binary classificaiton problem, I felt more confident. Upon my first sequential model run, I acheived over 90% accuracy. I used a sigmoid output node activation for the binary output expected. I changed the optimizer to Adam since I didn't use One Hot encoding, and used binary-cross entropy for loss and acheived a max of 95%.

My next model was going to be a simple RNN. I had to preprocess the text slightly differently in this case. I kept the same preprocessing from before but added padding so everything was the same length. This model took much longer to run and gave horrendous accuracy. I then tried to overcome any potential vanishing gradient issues by using an LSTM. The accuracy for this was also disappointing.

The last model I built dealt with embeddings. The embedding process felt very complicated and the model took seemingly forever to run (this could be due to a slow internet connection), but the accuracy is pretty good at 95% and it learned very fast. I don't know if it is worth all of that fuss however for comparable results to the oringinal sequential model.