

1. Preprocessing method:

- I. Turn to lowercase
- II. Remove stopwords
- III. Apply Lemmatization
- IV. Remove punctuation
- V. Remove HTML Tags

Preprocessing method 效果如下

```
text: One of the other reviewers has mentioned that after watching just 1 oz episode you'll be hooked. They are right, as this is exactly what happened with me.<br /><br />The fi
lowercase_text: one of the other reviewers has mentioned that after watching just 1 oz episode you'll be hooked. they are right, as this is exactly what happened with me.<br /><br />
removed_stopword_text: one reviewers mentioned watching 1 oz episode ' hooked. right , exactly happened me.<br /><br />the first thing struck oz brutality unflinching scenes vi
lemmatized_and_removed_punctuation_text: one reviewer mentioned watching 1 oz episode hooked. right exactly happened me.<br /><br />the first thing struck oz brutality unflinching
removed_html_tags_text: one reviewer mentioned watching 1 oz episode hooked. right exactly happened me.the first thing struck oz brutality unflinching scene violence set right wo
```

VI. 討論:

雖然處理方法眾多，但在測試後，發現只有 Turn to lowercase 有較好的效果，其餘方法反而造成 language model 有更差的表現，在使用 IMDB Dataset 與運用 Bigram+Sklearn.GaussianNB model 的詳細數據如下表(feature_num=500)。

	F1-Score	Precision	Recall
Without Preprocessing	0.7057	0.7088	0.7065
Only turn to lowercase	0.7094	0.7133	0.7104
Lowercase+Lemmatization	0.7079	0.7119	0.7089
Lowercase+Remove Punctuation	0.7033	0.7094	0.7049
Lowercase+Remove HTML Tags	0.7094	0.7027	0.7102
Lowercase+Remove Stopwords	0.7094	0.7133	0.7104

2. Tuning parameter 'feature_num' (default=500)

討論: 在有 Preprocessing 的情況下，調整 feature_num 大小的數據如下表，可以看見 feature_num 越高，模型的表現越好，但同時花費的時間也提升不少。

	F1-Score	Precision	Recall
feature_num=100	0.6116	0.6206	0.6156
feature_num=300	0.6823	0.6864	0.6835
feature_num=500	0.7094	0.7133	0.7104
feature_num=700	0.7191	0.7238	0.7202
feature_num=900	0.7395	0.7435	0.7403
feature_num=1500	0.7637	0.7647	0.7639

3. Discuss about using Bigram and Sklearn GaussianNB:

I. Perplexity:

II. 在沒有 preprocess 的情況下，部分 corpus 的 perplexity 如下圖。

```
Perplexity of ngram: [2832.1228179467385, 2936.2994473619383, 2245.2861807423706, 6074.748386035511, 836508, 1839.6378053806302, 9487.195033801661, 986.2539523637664, 6508.375198516267, 3112.69697418, 8247434624, 5037.0244710471425, 5365.513701762085, 1902.4016102355909, 757.6795341288456, 4875.666, 4222713840961, 1676.5305172370336, 6065.52875695023, 5038.53123845883, 5664.506638359498, 3291.285, 8745939534, 7137.97625803662, 3107.778938064599, 3956.9246926270816, 4881.83089438037, 3729.933559, 14972196057, 3419.6223047210583, 5115.022334492303, 3012.5182933566452, 2176.616304165438, 2677.30, 63.95706926755, 5500.626895841658, 1886.1436723688787, 8377.511130525407, 3848.2533581721723, 6585, 3434.4855435598056, 3025.4313473399793, 5257.543804701354, 4699.374423902551, 2729.277139343587, 2, 3387.8174156059213, 3939.0310584003773, 5818.728704799422, 2615.604019484288, 3051.914512182716
```

而在有 preprocess 的情況下，部分 corpus 的 perplexity 如下圖，可發現沒有 preprocess 時的 perplexity 比有 preprocess 時的 perplexity 還大。

```
Perplexity of ngram: [1920.4022088115003, 2300.3537576648787, 1875.7350341693098, 4584.805146733251, 1429.69653043, 177212, 1536.8018423549386, 7594.587084331458, 819.9574352596497, 5046.432289525344, 2167.043830831896, 2414.67949, 0035925, 4216.790597921349, 3582.8802064631054, 1541.5602927632629, 613.2776934214475, 3739.6032695171375, 3175.86, 961027259194, 1095.8906262500486, 4449.769771406573, 3908.6596186052866, 4437.867177086313, 2294.4811941025173, 17, 086.3839229688308, 5347.065736631145, 2032.813073786227, 2640.785375660839, 3999.4127688245726, 3011.070638056011, , 1194.7432866293507, 2628.114804492444, 4156.4389014439585, 2297.4107728374547, 1734.6135070745893, 2074.66345421, 066947, 3786.333833783748, 4240.6894096601, 1501.9756615326137, 6462.194792662504, 3043.3828484178352, 4913.311062, 5514232, 2748.6582154020166, 2247.427975286093, 4008.462177134413, 3728.473206620272, 2274.69135341245, 5419.91720, 68745056, 2507.7782569531855, 3258.5829706752716, 4744.387081877004, 2003.9356898132246, 2533.073274885115, 3269.6, .2874452693532, 2477.6560317563603, 2508.855318057731, 2126.6022381251073, 3207.53660477323, 3375.3237715912555, 2, 980.9422043180234, 2856.0804607365862, 4101.647688475028, 767.7721463099997, 688.8146580624874, 6516.656988050756,
```

III. F1-Score, Precision, Recall 在 Preprocessing 與否的比較如下表(feature_num=500)

	F1-Score	Precision	Recall
Without Preprocessing	0.7057	0.7088	0.7065
With Preprocessing	0.7094	0.7133	0.7104

IV. Encounter zero division error:

運用 Laplace Smoothing 解決，避免計算時分母為零。

4. Bigram and DistilBert

I. Bigram:

運用 Markov 假設，當前這個詞僅僅與前一個詞相關，機率如下圖。

- Bigram Model :

$$P(w_i | w_{i-1}) = \frac{C(w_{i-1} w_i)}{C(w_{i-1})}$$

DistilBert:

較為簡易的 BERT，保留了 97% 的語言理解能力，但同時減少空間大小 40% 且加速 60%。它是一種深度雙向、無監督、且僅使用純文字進行 Pretrain 的模型，它疊加多層 Transformer Encoder，可以更好地理解語意，並在各 NLP 項目中被不斷使用。

- II. Bigram 為當前這個詞僅僅與前一個詞相關，但 BERT 則是利用 Self-Attention 機制，讓訊息長度過長時仍不易丟失訊息，在這方面上可看出 BERT 的表現較 Bigram 好。詳細表現數據如下圖，可發現 DistilBert 的表現確實比 Bigram 好。

	F1-Score	Precision	Recall
Bigram without peprocessing	0.7057	0.7088	0.7065
Bigram with preprocessing	0.7094	0.7133	0.7104
DistilBert without preprocessing	0.9329	0.9333	0.9329
DistilBert with preprocessing	0.9333	0.9331	0.9333

```
ddeng@LAPTOP-DLIGRF6B MINGW64 /c/thomas/NYCU_Courses/Second_Semester/AI/HW2
$ python main.py --model_type BERT --preprocess 0 --part 2
Some weights of the model checkpoint at distilbert-base-uncased were not used when initializing DistilBertModel: ['vocab_layer_norm.weight', 'vocab_projector.weight', 'vocab_transform.bias']
- This IS expected if you are initializing DistilBertModel from the checkpoint of a model trained on another task or with another architecture (e.g. initializing a BertForPreTraining model).
- This IS NOT expected if you are initializing DistilBertModel from the checkpoint of a model that you expect to be exactly identical (e.g. initializing a BertForSequenceClassification model).
second_part
100%|#####| 5000/5000 [11:59:29<00:00, 8.63s/it]
100%|#####| 1/1 [12:35:13<00:00, 45313.67s/it]it]
Epoch: 0, F1 score: 0.9329, Precision: 0.933, Recall: 0.9329, Loss: 0.2285
end from second_part
ddeng@LAPTOP-DLIGRF6B MINGW64 /c/thomas/NYCU_Courses/Second_Semester/AI/HW2
$ python main.py --model_type BERT --preprocess 1 --part 2
Some weights of the model checkpoint at distilbert-base-uncased were not used when initializing DistilBertModel: ['vocab_layer_norm.weight', 'vocab_projector.weight', 'vocab_transform.bias']
- This IS expected if you are initializing DistilBertModel from the checkpoint of a model trained on another task or with another architecture (e.g. initializing a BertForPreTraining model).
- This IS NOT expected if you are initializing DistilBertModel from the checkpoint of a model that you expect to be exactly identical (e.g. initializing a BertForSequenceClassification model).
second_part
100%|#####| 5000/5000 [11:43:58<00:00, 8.45s/it]
100%|#####| 1/1 [12:19:36<00:00, 44376.39s/it]it]
Epoch: 0, F1 score: 0.933, Precision: 0.9331, Recall: 0.933, Loss: 0.2297
end from second_part
```

- III. Preprocessing method 會增進 Bigram 模型的表現，因為這會萃取出一個句子中重要的詞彙、刪除多餘標點符號、詞型還原等等，這些都能對模型解語意有一定的幫助。

IV. Convert words as [UNK]:

- A. Convert all words that appeared less than 10 times as [UNK]: Perplexity 如下圖。

```
Perplexity of ngram: [603.5730264555942, 724.5841506276375, 586.3906653245771, 1278.48
6123, 473.89539514732826, 2177.7158748389575, 269.4816784103672, 1437.8424693736156, 6
26582806, 1252.2735717012981, 1113.7477146274666, 509.97111280749243, 204.897485145005
8830748541894, 358.9994027775824, 1383.9896608058611, 1168.589916637615, 1373.55781306
.127960298605, 1588.0555617296168, 633.7151286559109, 815.6158892444593, 1255.77839147
0796869303726, 775.8774682503217, 1262.5734904860597, 682.9732143785474, 563.367087531
67.6709674935425, 1307.6491963813403, 486.39696384664165, 2036.5249206679448, 925.9953
95, 831.7645716587524, 646.4355948335516, 1214.3774879634552, 1146.2771849004062, 734.
870077, 783.8783141487954, 969.6990514999237, 1402.5083900045004, 612.7917466502346, 7
4989449, 755.6748515462572, 800.7220976351971, 657.5478104284218, 1021.0970282787832,
00460952, 834.4849539423797, 1181.553034558805, 248.89809497925407, 232.69683788467106
24328618594, 830.8480815988617, 530.758211122687, 675.7015089274429, 654.2070677919247
858975560554, 898.0936759408656, 919.8431818773006, 736.9341290479138, 891.08523600318
002092345808, 228.50579648394412, 1377.3052340655042, 2181.9267416241537, 891.75672069
12.6358407214939, 680.9513436576775, 890.2526476259181, 1762.2827866212444, 823.604749
364.5335759907654, 697.7929476676472, 646.9593824805971, 1660.0727485812022, 567.70454
25, 585.9819168277257, 1211.0341386177413, 653.0659499670928, 971.147882311628, 993.75
```

B. Converts the words that appeared just once as [UNK]: Perplexity 如下圖。

```
Perplexity of ngram: [1128.0563056317278, 1357.8218646317814, 1113.114782880541, 2590.11812
870614, 906.1797847012732, 4365.501542953433, 487.1933495883005, 2874.8532571578526, 1278.3
1733079, 2451.9559381239806, 2116.012551990264, 930.0996456902316, 369.7817703641593, 2186.
6891762, 657.4448178329424, 2632.115718818669, 2283.080812339434, 2595.5566204097277, 1371.
744104453124, 3115.2148630346546, 1198.2454007795434, 1523.9007721079608, 2374.201728981155
711.3110919051834, 1524.4112200570714, 2412.4644543275904, 1331.9470081076117, 1040.7841921
636, 2225.8359131764896, 2476.6857503140263, 893.4224041898475, 3828.304917412661, 1787.002
59275, 1602.7054802173527, 1285.7604048193703, 2343.927096775015, 2221.0161194009897, 1363.
3046326072, 1473.1494670937957, 1888.2441212037418, 2751.672421130073, 1178.351648984093, 1
.9855896310936, 1453.5529749900672, 1499.1413523772337, 1246.7731844691953, 1909.6566572641
, 1161.52366482305, 1647.6618629839643, 2350.8467893986935, 459.60453001183384, 416.3569091
5148, 1832.33954472406, 1573.1982850274653, 993.5806339581119, 1257.119016528424, 1245.3062
422141, 1839.0956264772044, 1733.1599982757834, 1695.5567025514406, 1379.173516432013, 1716
.6504576420843, 2114.705266270336, 415.6199443436014, 2660.246346164391, 4269.050984892953,
121.857760254081, 1876.8567648337619, 1309.1579891973583, 1751.3039125510006, 3496.16394184
78, 1359.610609505472, 677.4621740641337, 1380.8798522462894, 1281.784008997053, 3298.67035
```

C. 討論:

在 A.情況中的 perplexity 減少許多，在情況 B.中的 perplexity 減少幅度較小，而表現結果如下表，可看見在情況 A.與 B.下，f1-score 反而較原本的低，但幅度過小因此也難以就此定論。

	F1-Score	Precision	Recall
Bigram without replacing as [UNK]	0.7094	0.7133	0.7104
Convert all words that appeared less than 10 times as [UNK]	0.7046	0.7085	0.7056
Converts the words that appeared just once as [UNK]	0.7069	0.7112	0.7080

- 這次作業中總共遇到了幾個問題，首先是發現未經過 preprocess 時的 f1-score 比經過 preprocess 時高，且經過 preprocess 時的 perplexity 比未經過 preprocess 時高，因此在試驗過後，留下對本次 Sentiment analysis 有幫助的 preprocess method，成功提升 f1-score。接下來遇到的問題是使用 Bigram 時無論有沒有經過 preprocess 都發現 f1-score=0.3333，因此推論 bug 出現在 train_sentiment 中，最後找到 bug 並改善。