

	precision	recall	f1-score	support
0	0.84	0.93	0.88	347
1	0.83	0.71	0.76	299
2	0.62	0.84	0.63	305
accuracy			0.83	951
macro avg	0.83	0.83	0.83	951
weighted avg	0.83	0.83	0.83	951

	2	0.87	0.82	0.84	301
	accuracy			0.83	951
	macro avg	0.83	0.83	0.83	951
	weighted avg	0.83	0.83	0.83	951

```
In [29]: _ = get_model(tfidf_2, True)
```

	precision	recall	f1-score	support
0	0.83	0.91	0.87	352
1	0.86	0.72	0.78	315
2	0.82	0.86	0.84	284
	accuracy		0.83	951
	macro avg	0.84	0.83	951
	weighted avg	0.84	0.83	951

```
In [30]: _ = get_model(tfidf_5, True)
```

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

	2	0.82	0.87	0.85	286
accuracy				0.84	951
macro avg	0.84	0.84	0.84		951
weighted avg	0.84	0.84	0.84		951

```
In [31]: _ = get_model(tridf_10, True)
```

	precision	recall	f1-score	support
0	0.85	0.87	0.86	345
1	0.78	0.77	0.78	296
2	0.84	0.83	0.83	310
accuracy			0.83	951
macro avg	0.82	0.82	0.82	951
weighted avg	0.83	0.83	0.83	951

Answer The model with tfidf with min_df 5 is the best model out of the generated models based on its macro average f1-score, where it obtained 0.84. However, it only performs marginally better than the other models. We chose to use the macro average as a benchmark as the performance amongst the classes is quite balanced. It seems to perform similarly for all cases.

of the other models, but generally, the precision, recall, f1-score and macro average etc. have little variation across the other three models.

[4 points] Question 6: Inspecting the best scoring features

- Train the scikit-learn classifier (Naive Bayes) model with the following settings (airline tweets 80% training and 20% test; Bag of words representation (airline_count), min_df=2)
 - [1 point] a. Generate the list of best scoring features per class (see function **important_features_per_class**) [1 point]
- [3 points] b. Look at the lists and consider the following issues:
 - [1 point] Which features did you expect for each separate class and why?
 - [1 point] Which features did you not expect and why?
 - [1 point] The list contains all kinds of words such as names of airlines, punctuation, numbers and content words (e.g., 'delay' and 'bad'). Which words would you remove or keep when trying to improve the model and why?

In [32]:

```
def important_features_per_class(vectorizer, classifier, n=30):
    class_labels = classifier.classes_
    feature_names = vectorizer.get_feature_names()
    topn_class0 = sorted(zip(classifier.feature_count_[0], feature_names), reverse=True)[:n]
    topn_class1 = sorted(zip(classifier.feature_count_[1], feature_names), reverse=True)[:n]
    topn_class2 = sorted(zip(classifier.feature_count_[2], feature_names), reverse=True)[:n]
```

```
print("-----")
print("Important words in neutral documents")
for coef, feat in topn_class2:
    print(class_labels[1], coef, feat)
print("-----")
print("Important words in positive documents")
for coef, feat in topn_class3:
    print(class_labels[2], coef, feat)

# example of how to call from notebook:
important_features_per_class(vectorizer_2, clf)
```

```
Important words in negative documents
0 1491.0 #
0 1182.0 united
0 1232.0 .
0 405.0 '
0 405.0 ?
0 394.0 flight
0 380.0 !
0 294.0 #
0 212.0 n't
0 159.0 's
0 127.0 's
```

Answer The model with tfidf with min_df 5 is the best model out of the generated models based on its macro average f1-score, where it obtained 0.84. However, it only performs marginally better than the other models. We chose to use the macro average as a benchmark as the performance amongst the classes is quite balanced. It seems to perform similarly for all cases.

The frequency threshold seems to affect the scores, but rather minimally, since the model tfidf with min_df 5 performs the best out of the other models, but generally, the precision, recall, f1-score and macro average etc. have little variation across the other tfidf models.

[4 points] Question 6: Inspecting the best scoring features

- Train the scikit-learn classifier (Naive Bayes) model with the following settings (airline tweets 80% training and 20% test: Bag of words representation (airline_count), min_df=2)
- [1 point] a. Generate the list of best scoring features per class (see function `important_features_per_class` below) [1 point]
- [3 points] b. Look at the lists and consider the following issues:
 - [1 point] Which features did you expect for each separate class and why?
 - [1 point] Which features did you not expect and why?
 - [1 point] The list contains all kinds of words such as names of airlines, punctuation, numbers and content words (e.g., 'delay' and 'bad'). Which words would you remove or keep when trying to improve the model and why?

In [32]:	<pre>def important_features_per_class(vectorizer, classifier, n=80): class_labels = classifier.classes_ feature_names = vectorizer.get_feature_names() topn_class2 = sorted(zip(classifier.feature_count_[0], feature_names), reverse=True)[n:] topn_class2 = sorted(zip(classifier.feature_count_[1], feature_names), reverse=True)[n:] topn_class3 = sorted(zip(classifier.feature_count_[2], feature_names), reverse=True)[n:] print("Important words in negative documents") for coef, feat in topn_class1: print(class_labels[0], coef, feat) print("-----") print("Important words in neutral documents") for coef, feat in topn_class2: print(class_labels[1], coef, feat) print("-----") print("Important words in positive documents") for coef, feat in topn_class3: print(class_labels[2], coef, feat) # example of how to call from notebook: important_features_per_class(vectorizer_2, clf)</pre>
	<pre>Important words in negative documents 0 1491.0 @ 0 1382.0 united 0 1325.0 0 405.0 '' 0 405.0 0 384.0 flight 0 380.0 ! 0 294.0 0 202.0 'n't 0 159.0 '' 0 127.0 's 0 115.0 0 102.0 virginamerica 0 99.0 get 0 88.0 service 0 93.0 cancelled 0 88.0 time 0 87.0 bag 0 85.0 delayed 0 85.0 customer 0 84.0 plane 0 79.0 ... 0 78.0 - 0 75.0 'n 0 74.0 http 0 74.0 ; 0 72.0 hours 0 68.0 & 0 66.0 gate 0 65.0 hour 0 63.0 still 0 61.0 late 0 60.0 help 0 60.0 airline 0 58.0 amp 0 58.0 2 0 57.0 would 0 54.0 worst 0 51.0 one 0 49.0 flights 0 48.0 flightled 0 47.0 newer 0 47.0 like 0 46.0 ca 0 45.0 waiting 0 45.0 delay 0 45.0 've 0 44.0 S 0 43.0 back 0 43.0 J 0 42.0 I 0 40.0 us 0 40.0 lost 0 39.0 luggage 0 39.0 ever 0 39.0 check 0 39.0) 0 38.0 due 0 37.0 really 0 37.0 day 0 36.0 wait 0 36.0 seat 0 36.0 people 0 36.0 fly 0 36.0 another 0 34.0 u 0 33.0 trying 0 33.0 problems 0 33.0 hold 0 32.0 ticket 0 32.0 got 0 32.0 days 0 31.0 thanks 0 31.0 going 0 31.0 bags 0 31.0 baggage 0 30.0 last 0 30.0 even 0 30.0 airport 0 30.0 4 0 30.0 4 ----- Important words in neutral documents 1 1426.0 @ 1 1336.0 ? 1 503.0 1 320.0 jetblue 1 276.0 1 261.0 united 1 259.0 # 1 258.0 southwestair 1 247.0 1 233.0 flight 1 204.0 americanair 1 178.0 http 1 176.0 ! 1 164.0 usairways 1 136.0 's 1 86.0 get 1 77.0 1 75.0 virginamerica 1 72.0 '' 1 71.0 flights 1 63.0 please 1 60.0) 1 58.0 help 1 53.0 need 1 52.0 n't 1 51.0 ... 1 51.0 (1 47.0 ; 1 45.0 dm 1 44.0 tomorrow 1 43.0 us 1 40.0 flying 1 40.0 & 1 37.0 would 1 36.0 way 1 36.0 thanks 1 35.0 'm 1 34.0 know 1 34.0 cancelled 1 32.0 " 1 32.0 hi 1 32.0 fleet 1 32.0 fleek 1 31.0 1 31.0 one 1 31.0 amp 1 30.0 change 1 30.0 like 1 28.0 new 1 28.0 fly 1 28.0 could 1 27.0 number 1 26.0 time 1 26.0 go 1 25.0 airport 1 24.0 travel 1 23.0 want 1 23.0 see 1 23.0 check 1 22.0 guys 1 21.0 trying 1 21.0 today 1 21.0 ticket 1 21.0 make 1 21.0 going 1 21.0 follow 1 21.0 back 1 20.0 use 1 20.0 tickets 1 20.0 first 1 20.0 chance 1 20.0 2 1 19.0 want 1 19.0 trip 1 19.0 start 1 19.0 rt 1 19.0 add 1 18.0 think 1 18.0 still 1 18.0 seat ----- Important words in positive documents 2 1003.0 @ 2 1008.0 ! 2 744.0 . 2 318.0 2 306.0 southwestair 2 286.0 thanks 2 285.0 jetblue 2 243.0 thank 2 241.0 united 2 241.0 2 175.0 flight 2 170.0 americanair 2 167.0 2 135.0 usairways 2 129.0 great 2 88.0 service 2 83.0 virginamerica 2 81.0) 2 78.0 love 2 71.0 http 2 65.0 much 2 63.0 's 2 62.0 customer 2 61.0 best 2 58.0 guys 2 58.0 ; 2 54.0 awesome 2 49.0 airline 2 48.0 - 2 46.0 got 2 46.0 good 2 43.0 amazing 2 42.0 & 2 40.0 n't 2 39.0 us 2 39.0 today 2 39.0 time 2 36.0 help 2 36.0 get 2 34.0 made 2 33.0 gate 2 33.0 flying 2 32.0 crew 2 31.0 home 2 30.0 . 2 27.0 new 2 27.0 back 2 26.0 see 2 26.0 appreciate 2 25.0 nice 2 25.0 ? 2 25.0 're 2 24.0 tonight 2 24.0 day 2 23.0 work 2 23.0 u 2 23.0 response 2 23.0 '' 2 22.0 would 2 22.0 well 2 22.0 team 2 22.0 like 2 21.0 yes 2 21.0 plane 2 21.0 know 2 21.0 flights 2 21.0 ever 2 21.0 always 2 20.0 staff 2 20.0 first 2 20.0 class 2 20.0 'll 2 19.0 please 2 19.0 'n 2 18.0 tha 2 18.0 southwest 2 18.0 make 2 18.0 helpul C:\Users\Gebruiker\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:87: FutureWarning: Function get_ feature_names is deprecated; get_feature_names is deprecated in 1.0 and will be removed in 1.2. Please use g et_feature_names_out instead. warnings.warn(msg, category=FutureWarning)</pre>

Answer

For the negatively labeled documents, features like *cancelled*, *delayed* and *and*, are expected words related to reviews with a negative sentiment (or experience). Regarding the positively labeled document, positively associated words like *great*, *awesome* and *amazing* can also be seen as part of the important features. Lastly, for the neutral labeled documents, words like *help* and *would*, which are often associated with questions rather than direct positive or negative feedback, are expected to be part of the important words regarding that class.

In general, interperction, nouns related to organizations and numbers, are not the type of words that one would expect to be quite decisive in class labeling. Leaving those type of words out of the models, would in our opinion improve the models since these type of words often do not have strong sentimental meaning/association. On the other hand, words that do have a strong sentimental meaning/association would need to remain in order to have those type of words be more decisive for the labeling

[Optional! (will not be graded)] Question 7

Train the model on airline tweets and test it on your own set of tweets

- Train the model with the following settings (airline tweets 80% training and 20% test: Bag of words representation ('airline_count', min_df=2)
- Apply the model on your own set of tweets and generate the classification report
- [1 point] a. Carry out a quantitative analysis
- [1 point] b. Carry out an error analysis on 10 correctly and 10 incorrectly classified tweets and discuss them
- [2 points] c. Compare the results (cf. classification report) with the results obtained by VADER on the same tweets and discuss the differences.

[Optional! (will not be graded)] Question 8: trying to improve the model

- [2 points] a. Think of some ways to improve the scikit-learn Naive Bayes model by playing with the settings or applying linguistic preprocessing (e.g., by filtering on part-of-speech, or removing punctuation). Do not change the classifier but continue using the Naive Bayes classifier. Explain what the effects might be of these other settings
- [1 point] b. Apply the model with at least one new setting (train on the airline tweets using 80% training, 20% test) and generate the scores
- [1 point] c. Discuss whether the model achieved what you expected.

End of this notebook

In []:	
---------	--