



```

In [23]: 1 """ Title: Sentiment Analysis using Python
2 Author: SasikalaAlaguvel
3 Date: 2019.11.20
4 Availability: https://www.kaggle.com/sasikalal1/sentiment-analysis-using-python """
5 import pandas as pd
6 import numpy as np
7 from matplotlib import pyplot as plt
8 import seaborn as sns
9 import plotly.graph_objs as go
10 import plotly.offline as py
11 color = sns.color_palette()
12 #py.init_notebook_mode(connected=True)
13 import plotly.tools as tlsdata
14 import nltk
15 from nltk.stem.porter import *
16 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
17 from sklearn.cluster import KMeans
18 from gensim.models import word2vec
19
20 from sklearn.manifold import TSNE
21 from sklearn import metrics
22 import sklearn
23 from sklearn.metrics import jaccard_similarity_score
24 cv = CountVectorizer()
25 from nltk.corpus import stopwords
26 from sklearn.metrics.pairwise import cosine_similarity
27 stop = set(stopwords.words("english"))
28
29 import warnings
30 warnings.filterwarnings('ignore')
31 import os
32 os.listdir("../input")
33
34 #input data files are available in the "../input" directory
35 #for example, running this will list the files in the input directory
36 data = pd.read_csv('../input/1429_1.csv', encoding="ISO-8859-1")
37 #keeping only the necessary columns
38 #print(data.head())
39
40 #any results you write to the current directory are saved as output.
41
42 #-----

```

```

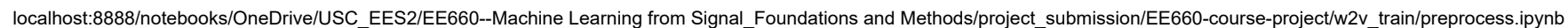
43
44 print(data.shape)
45 print(data.dtypes)
46 print(data.isnull().sum())
47 data = data.dropna(subset=['reviews.text'])
48
49 from wordcloud import WordCloud, STOPWORDS
50 stopwords = set(STOPWORDS)
51
52 def show_wordcloud(data, title = None):
53     wordcloud = WordCloud(
54         background_color='black',
55         stopwords=stopwords,
56         max_words=200,
57         max_font_size=40,
58         scale=3,
59         random_state=1
60     ).generate(str(data))
61
62     fig = plt.figure(1, figsize=(15,15))
63     plt.axis('off')
64     if title:
65         fig.suptitle(title, fontsize=20)
66         fig.subplots_adjust(top=2.3)
67
68     plt.imshow(wordcloud)
69     plt.show()
70
71
72 show_wordcloud(data['reviews.text'])

```

```
(34660, 21)
```

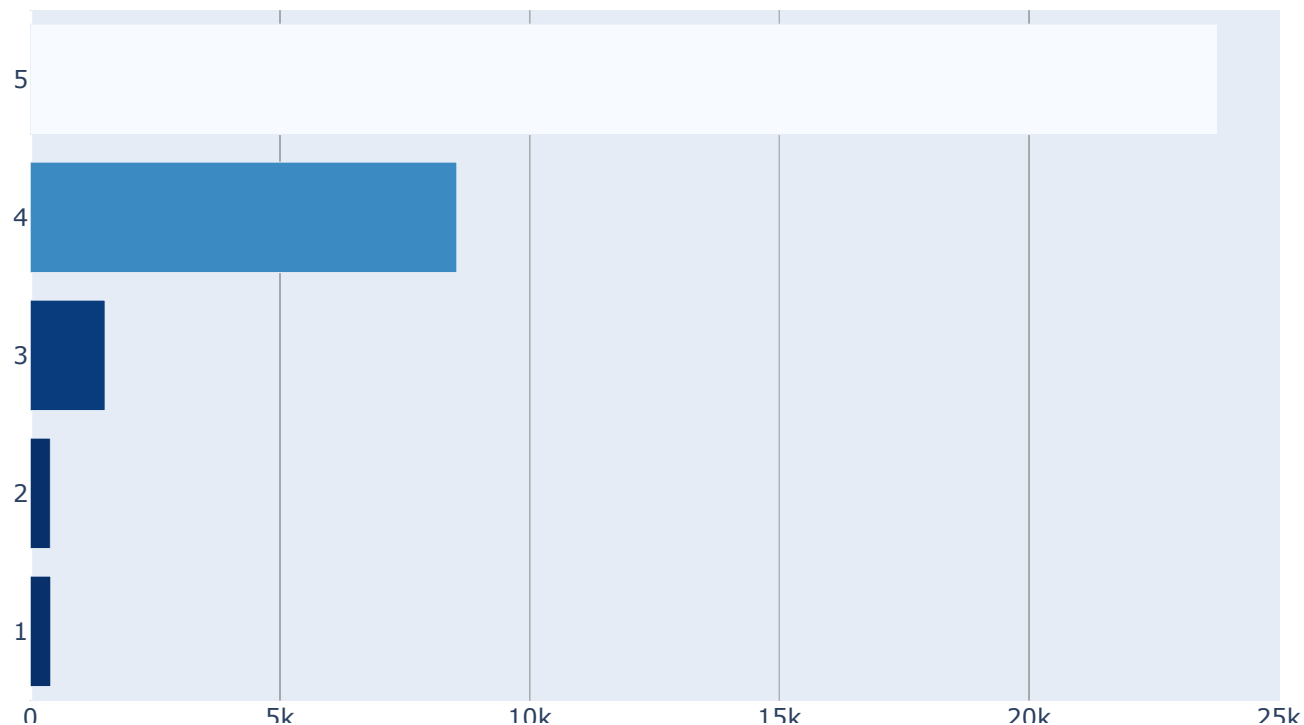
id	object
name	object
asins	object
brand	object
categories	object
keys	object
manufacturer	object
reviews.date	object
reviews.dateAdded	object
reviews.dateSeen	object

```
reviews.didPurchase    object
reviews.doRecommend    object
reviews.id              float64
reviews.numHelpful      float64
reviews.rating          float64
reviews.sourceURLs      object
reviews.text            object
reviews.title           object
reviews.userCity        float64
reviews.userProvince    float64
reviews.username        object
dtype: object
id                      0
name                   6760
asins                  2
brand                  0
categories             0
keys                   0
manufacturer           0
reviews.date           39
reviews.dateAdded      10621
reviews.dateSeen       0
reviews.didPurchase    34659
reviews.doRecommend    594
reviews.id             34659
reviews.numHelpful     529
reviews.rating         33
reviews.sourceURLs     0
reviews.text           1
reviews.title          5
reviews.userCity       34660
reviews.userProvince   34660
reviews.username       2
dtype: int64
```



```
In [2]: 1 cnt_srs = data['reviews.rating'].value_counts().head()
2 trace = go.Bar(
3     y=cnt_srs.index[::-1],
4     x=cnt_srs.values[::-1],
5     orientation='h',
6     marker=dict(
7         color=cnt_srs.values[::-1],
8         colorscale='Blues',
9         reversescale=True),
10 )
11 layout = dict(title='Ratings distribution')
12 data1 = [trace]
13 fig = go.Figure(data=data1, layout=layout)
14 py.iplot(fig, filename='Ratings')
```

## Ratings distribution





```

In [3]: 1 from subprocess import check_output
2 #print(check_output(["ls", "../input"]).decode("utf8"))
3 from mpl_toolkits.mplot3d import Axes3D
4 from mpl_toolkits.mplot3d import proj3d
5
6 from IPython.display import HTML
7 cat_hist = data.groupby('categories', as_index=False).count()
8 HTML(pd.DataFrame(cat_hist['categories']).to_html())
9 import nltk
10 from nltk import word_tokenize
11 from nltk.corpus import stopwords
12 import re
13 import string
14 from collections import Counter
15 from sklearn.model_selection import train_test_split
16 from sklearn.preprocessing import normalize
17 from sklearn.metrics import f1_score
18 from sklearn.naive_bayes import GaussianNB
19 def removePunctuation(x):
20     x = x.lower()
21     x = re.sub(r'[\x00-\x7f]', r' ', x) # replace the char that is not in the ASCII table
22     return re.sub('['+string.punctuation+']', " ", x)
23
24 stops = set(stopwords.words("english"))
25 def removeStopwords(x):
26     filtered_words = [word for word in x.split() if word not in stops]
27     return " ".join(filtered_words)
28 '''
29 When we deal with text problem in Natural Language Processing,
30 stop words removal process is a one of the important step to have a
31 better input for any models. Stop words means that it is a very
32 common words in a language (e.g. a, an, the in English. 的, 了 in
33 Chinese. え, も in Japanese). It does not help on most of NLP problem
34 such as semantic analysis, classification etc.
35 '''
36
37 def removeAmzString(x):
38     return re.sub(r'[0-9]+ people found this helpful\.. Was this review helpful to you Yes No', "", x)
39 # remove the amazon fixed sentence.
40

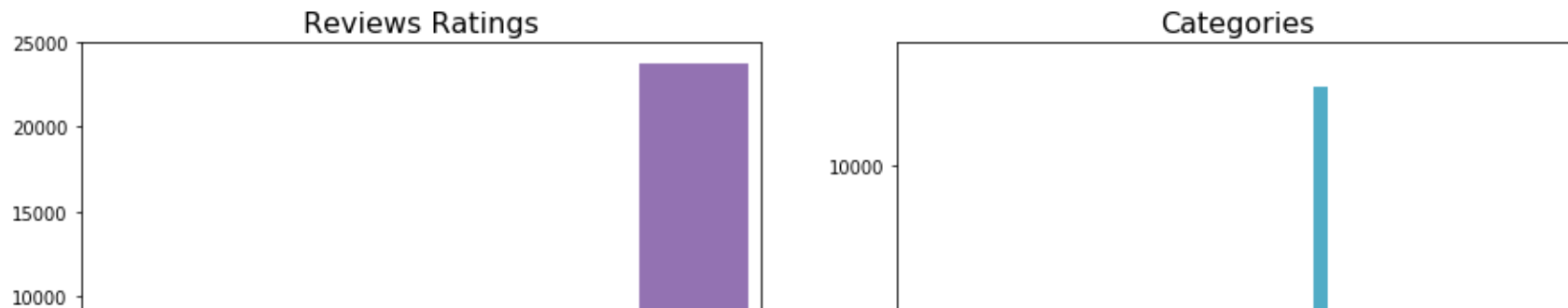
```

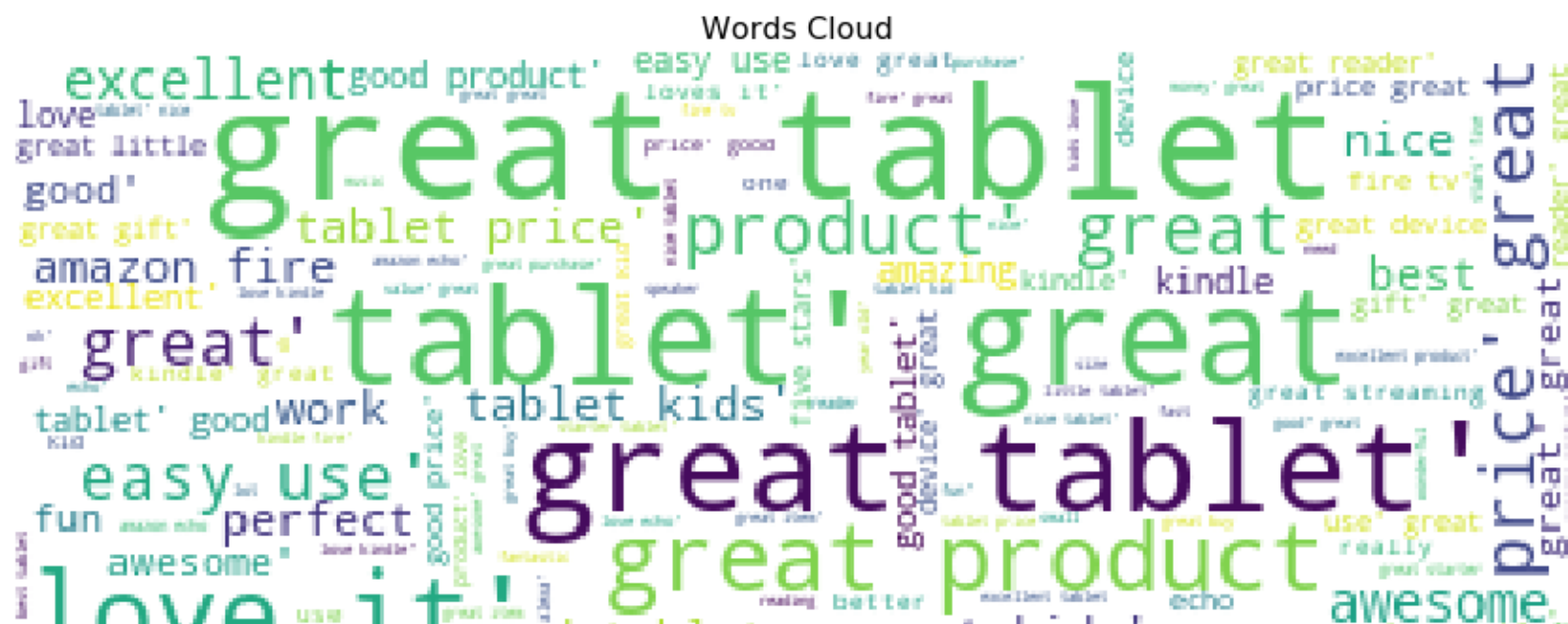
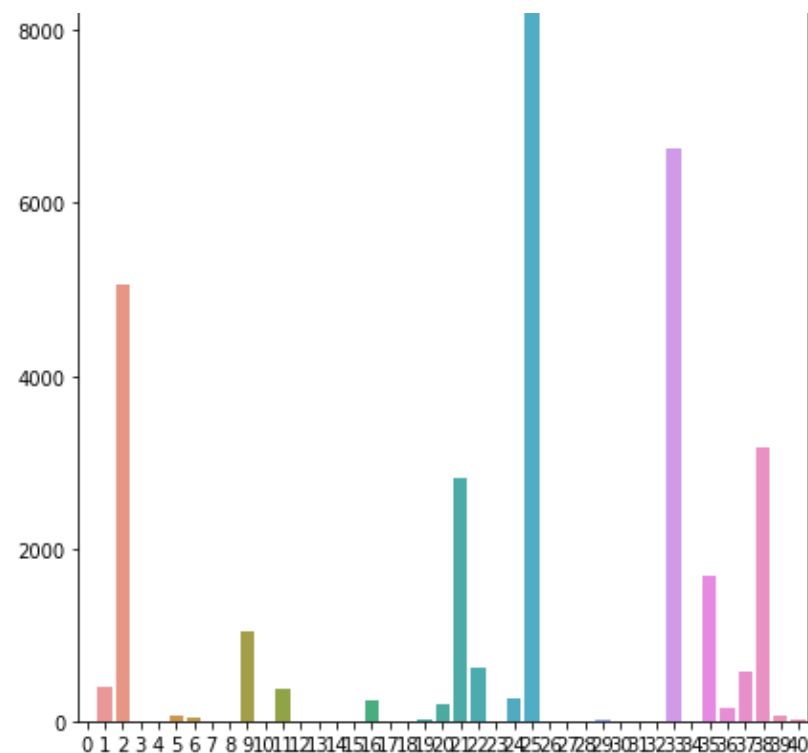


```

In [4]: 1 reviews = [sent if type(sent)==str else "" for sent in data['reviews.title'].values]
2 reviews = [removeAmzString(sent) for sent in reviews]
3
4 reviews = [removePunctuation(sent) for sent in reviews]
5
6 stopwords = set(STOPWORDS)
7 wordcloud = WordCloud(background_color='white', stopwords=stopwords, max_words=200,
8                       max_font_size=40, random_state=42).generate(str(reviews))
9 plt.figure(figsize=(15,20))
10 ax1 = plt.subplot2grid((4, 2), (0, 0))
11 ax2 = plt.subplot2grid((4, 2), (1, 0))
12 ax3 = plt.subplot2grid((4, 2), (0, 1), rowspan=2)
13 ax4 = plt.subplot2grid((4, 2), (2, 0), colspan=2, rowspan=2)
14
15 rat_hist = data.groupby('reviews.rating', as_index=False).count()
16 sns.barplot(x=rat_hist['reviews.rating'].values, y=rat_hist['id'].values, ax=ax1)
17
18 cat_hist = cat_hist.sort_values(by='id')
19 sns.barplot(x=cat_hist['categories'].index, y=cat_hist['id'].values, ax=ax3)
20
21 hf_hist = data.groupby('reviews.numHelpful', as_index=False).count()[0:30]
22 sns.barplot(x=hf_hist['reviews.numHelpful'].values.astype(int), y=hf_hist['id'].values, ax=ax2)
23
24 ax1.set_title("Reviews Ratings", fontsize=16)
25 ax3.set_title("Categories", fontsize=16)
26 ax2.set_title("Helpful Feedback", fontsize=16)
27 ax4.set_title("Words Cloud", fontsize=16)
28 ax4.imshow(wordcloud)
29 ax4.axis('off')
30 plt.show()
31

```







```

In [36]: 1 import nltk.stem as ns
2 from spellchecker import SpellChecker
3 def cleaning(s):
4     s = str(s)
5     s = s.lower()
6     s = re.sub('\s\W', ' ', s)
7     s = re.sub('\W\s', ' ', s)
8     s = re.sub(r'[\w]', ' ', s)
9     s = re.sub("\d+", "", s)
10    s = re.sub('\s+', ' ', s)
11    s = re.sub('[!@#$_]', '', s)
12    s = s.replace("co", "")
13    s = s.replace("https", "")
14    s = s.replace(", ", "")
15    s = s.replace("[\w*", " ")
16    s = s.replace(' s ', ' is ')
17    s = s.replace('don t', 'do not')
18    s = s.replace('doesn t', 'does not')
19    s = s.replace('can t', 'cannot')
20    s = s.replace('isn t', 'is not')
21    s = s.replace('uldn t', 'uld not')
22    s = s.replace('aren t', 'are not')
23    s = s.replace('wasn t', 'was not')
24    s = s.replace('weren t', 'were not')
25    s = s.replace('haven t', 'have not')
26    s = s.replace('hasn t', 'has not')
27    s = s.replace(' ve ', ' have ')
28    s = s.replace(' wa ', ' was ')
29    s = s.replace(' s ', '')
30    s = s.replace(' t ', '')
31    s = s.replace(' ntact', ' contact')
32    s = s.replace(' nnect', ' connect')
33    s = s.replace(' wasn t', ' was not ')
34    words = s.split()
35    lemmatizer = ns.WordNetLemmatizer()
36    for i in range(len(words)):
37        words[i] = lemmatizer.lemmatize(words[i], 'n')
38    s = " ".join(words)
39    return s
40    data['reviews.text'] = [cleaning(s) for s in data['reviews.text']]
41    data['revies.title'] = [cleaning(s) for s in data['reviews.title']]
42

```

```
In [37]: 1 from pathlib import Path
2 data = data.reset_index(drop=True)
3 outfile = Path('reviews.txt')
4 with outfile.open('w', encoding='utf-8') as w:
5     for k in range(data['reviews.text'].shape[0]):
6         w.write(data['reviews.text'][k]+'\\n')


In [38]: 1 from sklearn.model_selection import train_test_split
2 from collections import Counter
3 review = data['reviews.text'].to_numpy()
4 rate = data['reviews.rating'].to_numpy()
5 idx = np.argwhere(np.isnan(rate))
6 review_miss_rate = review[idx]
7 rate_miss = rate[idx]
8 review1 = np.delete(review, idx.T, axis=0)
9 rate1 = np.delete(rate, idx.T, axis=0)
10 review_train, review_test, rate_train, rate_test = train_test_split(review1,
11     rate1, test_size=0.2, random_state=10, stratify=rate1)
12 pd.DataFrame({'review':review_train}).to_csv('train_review.csv')
13 pd.DataFrame({'review':review_test}).to_csv('test_review.csv')
14 pd.DataFrame({'rate':rate_train}).to_csv('train_rate.csv')
15 pd.DataFrame({'rate':rate_test}).to_csv('test_rate.csv')
16 pd.DataFrame({'review':review_miss_rate.reshape(33)}).to_csv('miss_review_rate.csv')
17 pd.DataFrame({'rate':rate_miss.reshape(33)}).to_csv('miss_rate.csv')
18


In [39]: 1 outfile = Path('train_review.txt')
2 with outfile.open('w', encoding='utf-8') as w:
3     for k in range(review_train.shape[0]):
4         w.write(review_train[k]+'\\n')
5
6 outfile = Path('test_review.txt')
7 with outfile.open('w', encoding='utf-8') as w:
8     for k in range(review_test.shape[0]):
9         w.write(review_test[k]+'\\n')
```

```
In [ ]: 1 from sklearn.model_selection import train_test_split
2 df_train = pd.read_csv('review_train.csv')
3 df_train = df_train.dropna(axis=0, how='any')
4 df_train = df_train.reset_index(drop=True)
5 x = df_train['review'].to_numpy()
6 y = pd.read_csv('rate_train.csv')
7 y = y['rate'].to_numpy()
8 y = np.delete(y, 12675, axis=0) # we don't have comment on idx = 12675
9 # whole training set
10 #train, val, y_train, y_val = train_test_split(x,
11 #                                             y, test_size=0.2, random_state=10, stratify=y)
```

```
In [ ]: 1 from pathlib import Path
2 outfile = Path('../input/train.txt')
3 with outfile.open('w', encoding='utf-8') as w:
4     for k in range(train.shape[0]):
5         w.write(train[k]+'\\n')
6 pd.DataFrame({'review':train}).to_csv('../input/train.csv')
```

```
In [60]: 1 from pathlib import Path
2 outfile = Path('../input/val.txt')
3 with outfile.open('w', encoding='utf-8') as w:
4     for k in range(val.shape[0]):
5         w.write(val[k]+'\\n')
6 pd.DataFrame({'review':val}).to_csv('../input/val.csv')
```

```
In [ ]: 1 pd.DataFrame({'rate':y_train}).to_csv('../input/y_train.csv')
2 pd.DataFrame({'rate':y_val}).to_csv('../input/y_val.csv')
```



```

In [ ]: 1 """ Title: Word2Vec Model, Software Framework for Topic Modelling with Large Corpora
2 Author: Radim et al. and Petr Sojka
3 Date: 2010.5.22
4 Availability: https://radimrehurek.com/gensim/auto_examples/tutorials/run_word2vec.html#sphx-glr-download-auto-examples-
5 #generate a doc vector by averaging the word2vec vector
6 #use gensim.utils.simple_preprocess function
7 import pandas as pd
8 import numpy as np
9 df_train = pd.read_csv('../input/train.csv')
10 train = df_train['review'].to_numpy()
11 df_y_train = pd.read_csv('../input/y_train.csv')
12 y_train = df_y_train['rate'].to_numpy()
13 df_val = pd.read_csv('../input/val.csv')
14 val = df_val['review'].to_numpy()
15 df_y_val = pd.read_csv('../input/y_val.csv')
16 y_val = df_y_val['rate'].to_numpy()
17
18
19
20 # train word2vec model
21 import logging
22 logging.basicConfig(format='%(asctime)s : %(levelname)s : %(message)s', level=logging.INFO)
23 from gensim.test.utils import datapath
24 from gensim import utils
25 from pathlib import Path
26
27 class MyCorpus(object):
28     """An iterator that yields sentences (lists of str)."""
29
30     def __iter__(self):
31         corpus_path = Path('../input/train.txt')
32         for line in open(corpus_path):
33             # assume there's one document per line, tokens separated by whitespace
34             yield utils.simple_preprocess(line)
35
36 import gensim.models
37 sentences = MyCorpus()
38 model = gensim.models.Word2Vec(sentences=sentences, size=50, iter=5)
39 model.save('../embedding/word2vec.model')
40 model = gensim.models.Word2Vec.load('../embedding/word2vec.model')
41
42 # generate sentence vector for each sentence

```



```
43 word_vectors = model.wv
44 len(word_vectors.vocab)
45 for k, seq in enumerate(train):
46     tokens = gensim.utils.simple_preprocess(seq)
47     l = 0.
48     for i in tokens:
49         if i in word_vectors.vocab:
50             ave_vec[k] += model.wv[i]
51             l += 1
52     ave_vec[k] /= l
53 print(ave_vec.shape)
```

```
In [15]: 1 import pandas as pd
2 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
3 from matplotlib import pyplot as plt
4 from sklearn.metrics import f1_score
5 from sklearn import linear_model
6 from sklearn.metrics import confusion_matrix
7 import numpy as np
8 from sklearn.externals import joblib
9 #import warnings
10 #warnings.filterwarnings("ignore")
11 df_train = pd.read_csv('../input/train.csv')
12 df_y_train = pd.read_csv('../input/y_train.csv')
13 df_val = pd.read_csv('../input/val.csv')
14 df_y_val = pd.read_csv('../input/y_val.csv')
15 train = df_train['review'].to_numpy()
16 y_train = df_y_train['rate'].to_numpy()
17 val = df_val['review'].to_numpy()
18 y_val = df_y_val['rate'].to_numpy()
19
20 # the data set is merge here for crossvalidation in the Bayesian Inference analysis.
21 '''
22 train_total = pd.DataFrame({'review':train.tolist()+val.tolist()})
23 train_y_total = pd.DataFrame({'rate':y_train.tolist()+y_val.tolist()})
24 train_total.to_csv('../input/train_total.csv')
25 train_y_total.to_csv('../input/train_y_total.csv')
26 '''
27
28 vectorizer = TfidfVectorizer(stop_words='english', use_idf=True)
29 model_tr = vectorizer.fit_transform(train)
30 model_val = vectorizer.transform(val)
```

```
In [14]: 1 from sklearn.decomposition import PCA
2         #from sklearn.preprocessing import normalize
3         #model_tr_normalize = normalize(model_tr, norm='l2', axis=0, copy=True, return_norm=False)
4         pca = PCA(n_components = 0.95)
5         pca.fit(model_tr.toarray())
6         reduced = pca.transform(model_tr.toarray())
7         print(reduced.shape)
```

(22159, 3192)

```
In [16]: 1 from sklearn.decomposition import PCA
2         model_tr1 = model_tr.copy()
3         model_val1 = model_val.copy()
4         #pca = joblib.load('../input/model/pca.pkl')
5         model_val = pca.transform(model_val.toarray())
6         model_tr = pca.transform(model_tr.toarray())
7
```

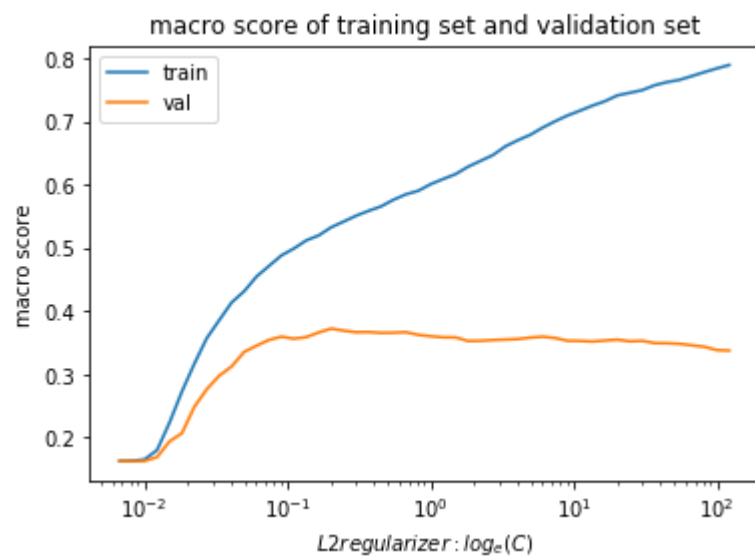
```

In [4]: 1 # Logistic regression, L2 regularization, weighted loss, hyperparameter selection use validation set, macro-weighted score
2 # reason: we assume that all the variables are necessary for classification so we use a L2 regularizer.
3 # and we give more weight to the minor class to remedy for the class imbalance.
4 import warnings
5 from matplotlib import pyplot as plt
6 from sklearn.metrics import f1_score
7 from sklearn.metrics import confusion_matrix
8 warnings.filterwarnings("ignore")
9 from sklearn.linear_model import LogisticRegression
10 n_alphas = 76
11 #alphas = np.exp(np.linspace(-5, 10, n_alphas))
12 C = np.exp(np.linspace(-5, 10, n_alphas))
13 train_score = []
14 val_score = []
15 for i in C:
16     m = LogisticRegression(C=i, class_weight='balanced').fit(model_tr, y_train)
17     s = f1_score(y_val, m.predict(model_val), average='macro')
18     train_score.append(f1_score(y_train, m.predict(model_tr), average='macro'))
19     val_score.append(f1_score(y_val, m.predict(model_val), average='macro'))
20     print('Val macro score is', s)
21 #plt.plot(np.log10(C), train_score, label='train')
22 #plt.plot(np.log10(C), val_score, label='val')
23 #plt.title('macro score of training set and validation set')
24 #plt.legend()
25 #plt.xlabel('log10(C)')
26 #plt.ylabel('macro score')
27 #plt.show()
28 #-----

```

...

```
In [5]: 1 from sklearn.metrics import confusion_matrix
2 ax = plt.gca()
3 # randomly choose 50 coefficients in the 9676 variables
4 # to plot figure of coefs vs regularizer.
5 ax.plot(C[:len(train_score)], train_score, label='train')
6 ax.plot(C[:len(train_score)], val_score, label='val')
7 ax.set_xscale('log')
8 #ax.set_xlim(ax.get_xlim()[::-1]) # reverse axis
9 plt.xlabel('$L2$ regularizer: log_e(C)$')
10 plt.ylabel('macro score')
11 plt.legend()
12 plt.title('macro score of training set and validation set')
13 plt.axis('tight')
14 plt.show()
15 C_opt = C[np.argmax(val_score)]
16 m_opt = LogisticRegression(C=C_opt, class_weight='balanced').fit(model_tr, y_train)
17 s = f1_score(y_val, m_opt.predict(model_val), average='macro')
18 mat = confusion_matrix(y_val, m_opt.predict(model_val))
19 print('Best Val macro score is', s)
20 print('confusion matrix is: \n', mat)
21 print('recall of each class:\n', [round(mat[0,0]/mat[0,:].sum(), 4), round(mat[1,1]/mat[1,:].sum(), 4),
22                                   round(mat[2,2]/mat[2,:].sum(), 4), round(mat[3,3]/mat[3,:].sum(), 4),
23                                   round(mat[4,4]/mat[4,:].sum(), 4)])
24 joblib.dump(m_opt, '../input/model/logistic_L2.pkl')
```



Best Val macro score is 0.3724644780023236

confusion matrix is:

```
[[ 31  12   8   2  12]
 [ 18  12  13   6  15]
 [ 22  23  73  39  83]
 [ 25  25 121 299 897]
 [ 33  37  78 262 3394]]
```

recall of each class:

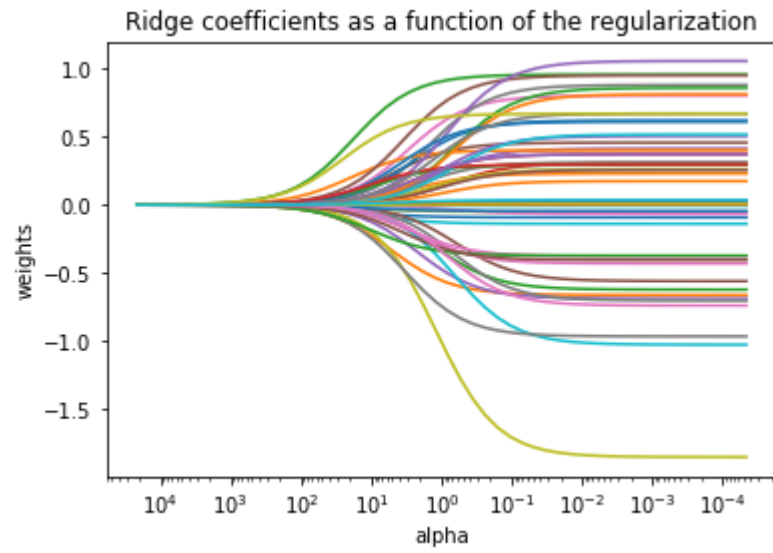
```
[0.4769, 0.1875, 0.3042, 0.2187, 0.8922]
```

Out[5]: ['../input/model/logistic\_L2.pkl']

```

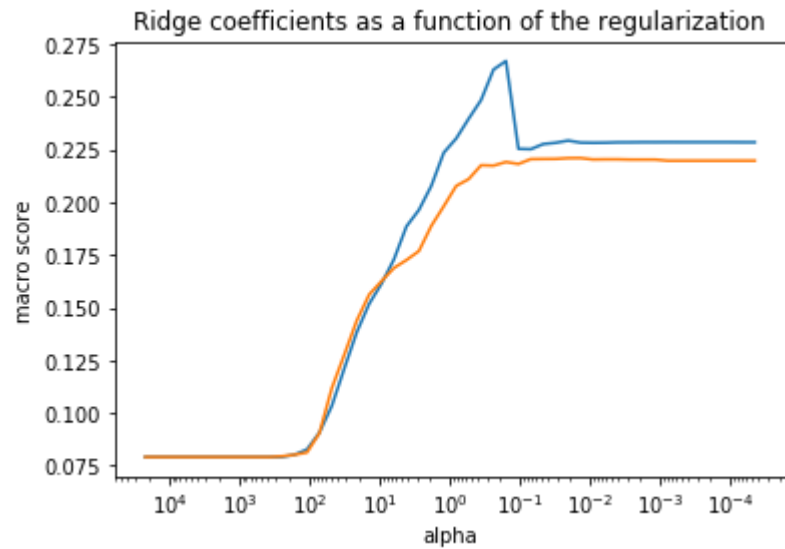
In [66]: 1 # ridge regression (L2 regularizer)
          2 #
          3 from sklearn import linear_model
          4 n_alphas = 50
          5 alphas = np.exp(np.linspace(-10, 10, n_alphas))
          6 coefs = []
          7 train_score = []
          8 val_score = []
          9 for a in alphas:
10     ridge = linear_model.Ridge(alpha=a, fit_intercept=True)
11     ridge.fit(model_tr, y_train)
12     train_pre = ridge.predict(model_tr).astype('int')
13     train_pre[train_pre > 5] = np.int(5)
14     train_pre[train_pre < 0] = np.int(1)
15     val_pre = ridge.predict(model_val).astype('int')
16     val_pre[val_pre > 5] = np.int(5)
17     val_pre[val_pre < 1] = np.int(1)
18     train_score.append(f1_score(y_train, train_pre, average='macro'))
19     val_score.append(f1_score(y_val, val_pre, average='macro'))
20     coefs.append(ridge.coef_)
21
22 cc = np.array(coefs)
23
24 ax = plt.gca()
25 # randomly choose 50 coefficients in the 9676 variables
26 # to plot figure of coefs vs regularizer.
27 ax.plot(alphas, cc[:, np.random.randint(0, 3192, 50)])
28 ax.set_xscale('log')
29 ax.set_xlim(ax.get_xlim()[::-1]) # reverse axis
30 plt.xlabel('alpha')
31 plt.ylabel('weights')
32 plt.title('Ridge coefficients as a function of the regularization')
33 plt.axis('tight')
34 plt.show()

```





```
In [67]: 1 train_score = np.array(train_score)
2 val_score = np.array(val_score)
3 ax1 = plt.gca()
4 # randomly choose 50 coefficients in the 9676 variables
5 # to plot figure of coefs vs regularizer.
6 ax1.plot(alphas, train_score)
7 ax1.plot(alphas, val_score)
8 ax1.set_xscale('log')
9 ax1.set_xlim(ax1.get_xlim()[::-1]) # reverse axis
10 plt.xlabel('alpha')
11 plt.ylabel('macro score')
12 plt.title('Ridge coefficients as a function of the regularization')
13 plt.axis('tight')
14 plt.show()
15
16 alphas_best = alphas[np.argmax(val_score)]
17 print('the best macro score on validation set:', np.max(val_score))
18 ridge = linear_model.Ridge(alpha=alphas_best, fit_intercept=True).fit(model_tr, y_train)
19 val_pre = ridge.predict(model_val).astype('int')
20 val_pre[val_pre > 5] = np.int(5)
21 val_pre[val_pre < 1] = np.int(1)
22 mat = confusion_matrix(y_val, val_pre)
23 print('confusion matrix is: \n', mat)
24 print('f1 score of each class:\n', f1_score(y_val, val_pre, average='macro'))
25 #print('the number of wi equal to 0 =', len(np.where(np.abs(m.coef_.flatten())==0)[0]))
26
```



the best macro score on validation set: 0.2209726787766575

confusion matrix is:

```
[[ 0  13  31  19   2]
 [ 0  12  29  23   0]
 [ 1  17  88 130   4]
 [ 0  16 187 1067  97]
 [ 0   5 123 3027 649]]
```

f1 score of each class:

0.2209726787766575

```
In [68]: 1 joblib.dump(ridge, '../input/model/tf-idf_ridge.pkl')
```

```
Out[68]: ['../input/model/tf-idf_ridge.pkl']
```

## Decision Tree

There are 7 hyperparameters in Decision Tree classifier. Since the dataset in this project has a large size and a large number of variables. Here I use some prior knowledge of this task to determine some hyperparameters and the order hyperparameter tuning. In the tuning process, I will use greedy search to find out the optimal hyperparameter.

1."criterion": It defines the criterion whether a node should be split and the two options are gini-index and entropy. When using entropy criterion, or saying the information gain between parent node and child node, the degree of impurity in a node can be better enlarged than using gini-index. But since our task has too many variables, using entropy criterion will result in more overfitting. So gini-index will be chose

to be the criterion.

2."splitter": we will choose "best" here so that our algorithm will choose the most important feature to split each time. Although we can use "random" and then merge the leave using some criterion. But because our task is too complecated, splitting the best node each can provide better stability.

3.parameters that I use val score to choose:

The parameters are ordered in their degree of importance to the model and they are choose according to validation set score.

3.1.max\_depth: max\_depth is the most important hyperparameter in decision, because it mainly decides the ability of the tree to fit the train set. Here I tune max\_depth from 3 and fixed minimum of samples in leafs to 5, minmum of samples to split to 10. And give no constraint to maximum of feature number and minimum impurity decrease.

3.2.min\_samples\_leaf: min\_samples\_leaf can be use to smooth the boundary and with similar validation score a model with bigger min\_samples\_leaf can give a better boundary. I use the optimal max\_length and the search for the best min\_sample\_leaf.

3.3.min\_samples\_split: This hyperparameter also can be use to make the boundary smoother.

```

In [3]: 1 from sklearn.tree import DecisionTreeClassifier
2 import warnings
3 warnings.filterwarnings("ignore")
4 clf = DecisionTreeClassifier(random_state=0, criterion='gini',
5                             splitter='best' )
6
7 depth = np.linspace(3, 183, 61).astype('int')
8 train_score1 = []
9 val_score1 = []
10 for i in depth:
11     m = DecisionTreeClassifier(random_state=0,
12                               criterion='gini',
13                               splitter='best',
14                               max_depth=i,
15                               min_samples_leaf=5,
16                               min_samples_split=10,
17                               class_weight='balanced'
18                               )
19     m.fit(model_tr, y_train)
20     s = f1_score(y_val, m.predict(model_val), average='macro')
21     train_score1.append(f1_score(y_train, m.predict(model_tr), average='macro'))
22     val_score1.append(f1_score(y_val, m.predict(model_val), average='macro'))
23     #print ('Val macro-weighted score is', s)
24
25 train_score1 = np.array(train_score1)
26 val_score1 = np.array(val_score1)
27
28 min_leaf = np.linspace(1, 101, 51).astype('int')
29 train_score2 = []
30 val_score2 = []
31 for i in min_leaf:
32     m = DecisionTreeClassifier(random_state=0,
33                               criterion='gini',
34                               splitter='best',
35                               max_depth=depth[np.argmax(val_score1)],
36                               min_samples_leaf=i,
37                               min_samples_split=10,
38                               class_weight='balanced'
39                               )
40     m.fit(model_tr, y_train)
41     s = f1_score(y_val, m.predict(model_val), average='macro')
42     train_score2.append(f1_score(y_train, m.predict(model_tr), average='macro'))

```

```

43     val_score2.append(f1_score(y_val, m.predict(model_val), average='macro'))
44
45 min_samples_spl = np.linspace(5, 50, 17).astype('int')
46 train_score3 = []
47 val_score3 = []
48 for i in min_samples_spl:
49     m = DecisionTreeClassifier(random_state=0,
50                               criterion='gini',
51                               splitter='best',
52                               max_depth=depth[np.argmax(val_score1)],
53                               min_samples_leaf=min_leaf[np.argmax(val_score2)],
54                               min_samples_split=i,
55                               class_weight='balanced'
56                               )
57     m.fit(model_tr, y_train)
58     s = f1_score(y_val, m.predict(model_val), average='macro')
59     train_score3.append(f1_score(y_train, m.predict(model_tr), average='macro'))
60     val_score3.append(f1_score(y_val, m.predict(model_val), average='macro'))
61
62 m_opt = DecisionTreeClassifier(random_state=0,
63                               criterion='gini',
64                               splitter='best',
65                               max_depth=depth[np.argmax(val_score1)],
66                               min_samples_leaf=min_leaf[np.argmax(val_score2)],
67                               min_samples_split=min_samples_spl[np.argmax(val_score3)],
68                               class_weight='balanced'
69                               )
70 m_opt.fit(model_tr, y_train)
71 mat = confusion_matrix(y_val, m_opt.predict(model_val))
72 print('optimal score', f1_score(y_val, m_opt.predict(model_val), average='macro'))
73 print('confusion matrix is: \n', mat)
74 print('accuracy of each class:\n', [round(mat[0,0]/mat[0,:].sum(), 4), round(mat[1,1]/mat[1,:].sum(), 4),
75                                     round(mat[2,2]/mat[2,:].sum(), 4), round(mat[3,3]/mat[3,:].sum(), 4),
76                                     round(mat[4,4]/mat[4,:].sum(), 4)])
77
78

```

optimal score 0.23399102471368582

confusion matrix is:

```

[[ 4  5  8 19 29]
 [ 3  7  2 25 27]
 [ 5 16 38 73 108]

```

```
[ 29  36 124 433 745]
[ 75  74 283 1067 2305]]
accuracy of each class:
[0.0615, 0.1094, 0.1583, 0.3168, 0.6059]
```

```
In [4]: 1 plt.figure()
2 plt.plot(depth, train_score1, label='train')
3 plt.plot(depth, val_score1, label='val')
4 plt.legend()
5 plt.title('depth vs f1-macro')
6 print(depth[np.argmax(val_score1)])
7 print(np.max(val_score1))
8 pass
9 plt.figure()
10 plt.plot(min_leaf, train_score2, label='train')
11 plt.plot(min_leaf, val_score2, label='val')
12 plt.legend()
13 plt.title('min_leaf vs f1-macro')
14 print(min_leaf[np.argmax(val_score2)])
15 print(np.max(val_score2))
16 pass
17 plt.figure()
18 plt.plot(min_samples_spl, train_score3, label='train')
19 plt.plot(min_samples_spl, val_score3, label='val')
20 plt.legend()
21 plt.title('min_samples_spl vs f1-macro')
22 print(min_samples_spl[np.argmax(val_score3)])
23 print(np.max(val_score3))
24 pass
```

27

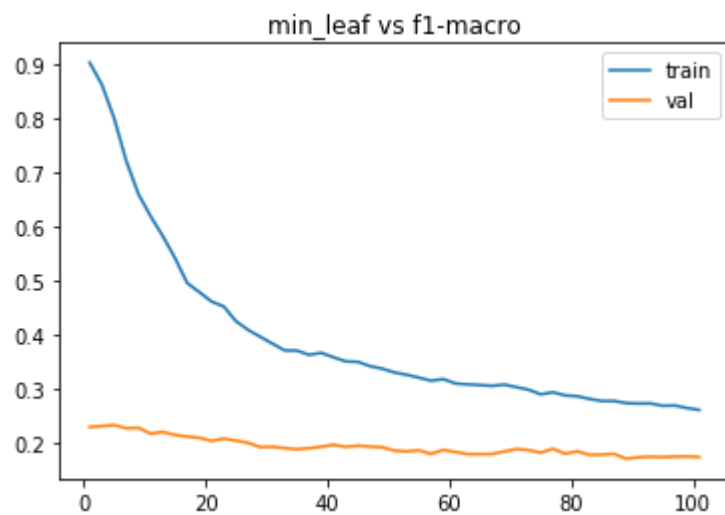
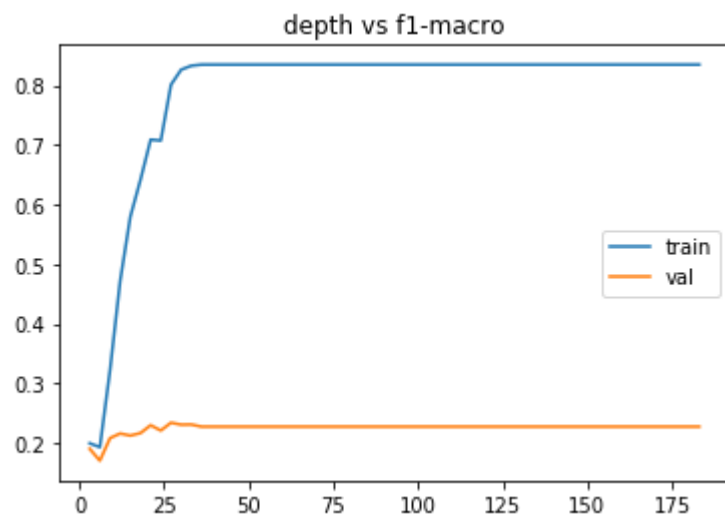
0.23370315113393864

5

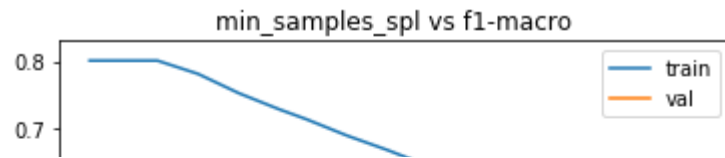
0.23370315113393864

13

0.23399102471368582







```
In [5]: 1 joblib.dump(m_opt, '../input/decision_tree.pkl')
```

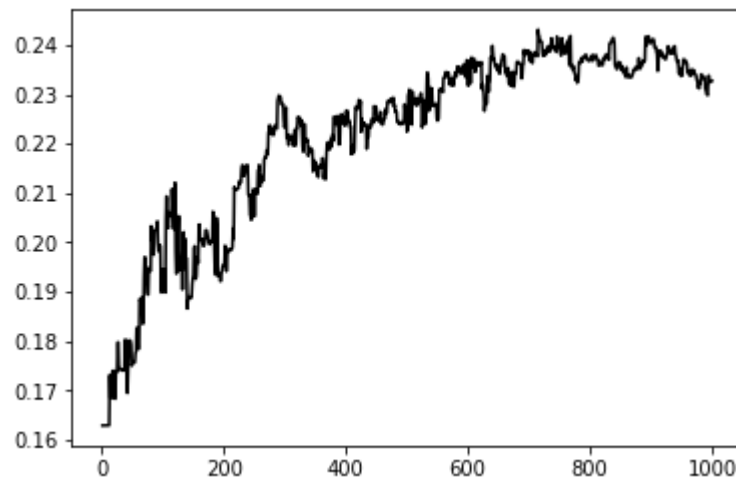
```
Out[5]: ['../input/decision_tree.pkl']
```

```

In [6]: 1 # Adaboost with tree method
        2 #import warnings
        3 #warnings.filterwarnings("ignore")
        4 from sklearn.ensemble import AdaBoostClassifier
        5 from sklearn.metrics import f1_score
        6 from sklearn.metrics import confusion_matrix
        7 clf = AdaBoostClassifier(n_estimators=1000, random_state=1)
        8 clf.fit(model_tr, y_train)
        9 real_val_macro = [0]
       10 for real_test_predict in clf.staged_predict(model_val):
       11     if f1_score(y_val, real_test_predict, average='macro') > np.max(real_val_macro) :
       12         pred_opt = real_test_predict
       13     real_val_macro.append(
       14         f1_score(y_val, real_test_predict, average='macro'))
       15
       16
       17
       18 n_trees_real = len(clf)
       19 plt.figure()
       20 plt.plot(range(1, n_trees_real + 1),
       21         real_val_macro[1:], c='black',
       22         label='SAMME. R')
       23
       24

```

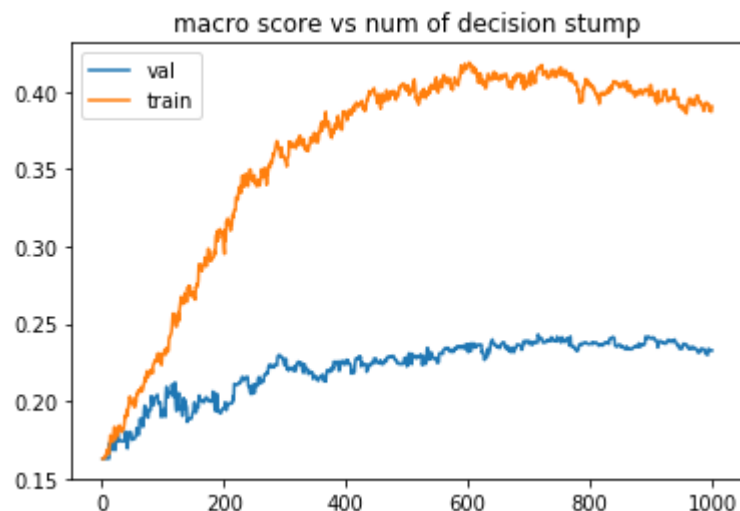
Out[6]: [<matplotlib.lines.Line2D at 0x7f6d9887b6d8>]



```

In [9]: 1 real_val_macro = [0]
        2 train_macro = []
        3 for real_test_predict in clf.staged_predict(model_val):
        4     if f1_score(y_val, real_test_predict, average='macro') > np.max(real_val_macro) :
        5         pred_opt = f1_score(y_val, real_test_predict, average='macro')
        6         real_val_macro.append(
        7             f1_score(y_val, real_test_predict, average='macro'))
        8 for real_train_predict in clf.staged_predict(model_tr):
        9     train_macro.append(f1_score(y_train, real_train_predict, average='macro'))
        10
        11
        12 plt.figure()
        13 plt.plot(range(1, n_trees_real + 1),
        14         real_val_macro[1:],
        15         label='val')
        16 plt.plot(range(1, n_trees_real + 1),
        17         train_macro,
        18         label='train')
        19 plt.legend()
        20 plt.title('macro score vs num of decision stump')
        21 adaboost_opt = AdaBoostClassifier(n_estimators=np.argmax(real_val_macro[1:])+1,
        22                                 random_state=1)

```



```
In [32]: 1 np.max(real_val_macro)
```

```
Out[32]: 0.24310635913853074
```

```
In [29]: 1 clf_ada_pca = AdaBoostClassifier(n_estimators=715, random_state=1).fit(model_tr,y_train)
```

```
In [64]: 1 joblib.dump(clf_ada_pca,'../input/model/adaboost1.pkl')
```

```
Out[64]: ['../input/model/adaboost1.pkl']
```

```
In [31]: 1 w2v_adaboost = clf_ada_pca
2 print('TF-IDF feature: Adaboost, best F1-macro score on validation set:',f1_score(y_val, w2v_adaboost.predict(model_val))
3 mat = confusion_matrix(y_val, w2v_adaboost.predict(model_val))
4 print('confusion matrix on validation set is: \n',mat)
5 print('F1 score of each class:\n',f1_score(y_val, w2v_adaboost.predict(model_val), average=None))
6 print('accuracy on validation set:', w2v_adaboost.score(model_val,y_val),'\n\n')
```

TF-IDF feature: Adaboost, best F1-macro score on validation set: 0.24310635913853074

confusion matrix on validation set is:

```
[[ 3  0  5 15 42]
 [ 1  0  9 17 37]
 [ 0  0 21 67 152]
 [ 2  0 49 300 1016]
 [ 3 11 61 574 3155]]
```

F1 score of each class:

```
[0.08108108 0.          0.10909091 0.25641026 0.76894955]
```

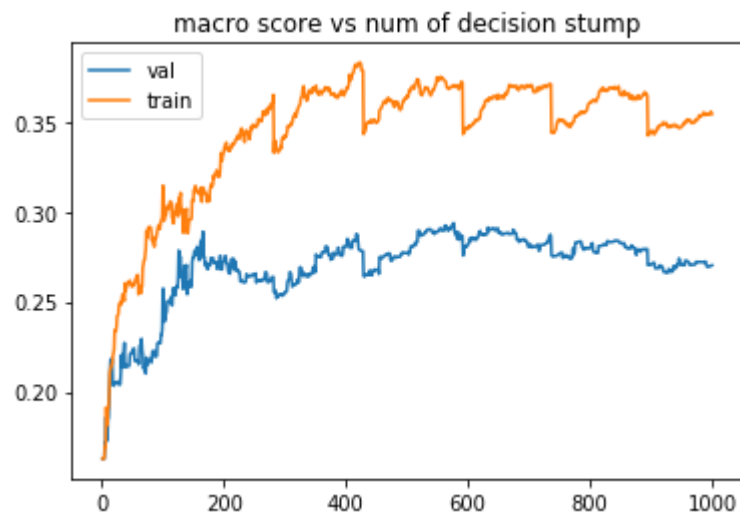
accuracy on validation set: 0.6279783393501805

```
In [20]: 1 # Adaboost with tree method with out PCA
2 #import warnings
3 #warnings.filterwarnings("ignore")
4 from imblearn.over_sampling import SMOTE # doctest: +NORMALIZE_WHITESPACE
5 from collections import Counter
6
7 from sklearn.ensemble import AdaBoostClassifier
8 from sklearn.metrics import f1_score
9 from sklearn.metrics import confusion_matrix
10 clf_no_pca = AdaBoostClassifier(n_estimators=1000, random_state=1)
11 clf_no_pca.fit(model_tr1, y_train)
12 real_val_macro1 = [0]
13 train_macro1 = []
14 for real_test_predict in clf_no_pca.staged_predict(model_val1):
15     if f1_score(y_val, real_test_predict, average='macro') > np.max(real_val_macro1) :
16         pred_opt1 = real_test_predict
17     real_val_macro1.append(
18         f1_score(y_val, real_test_predict, average='macro'))
19 for real_train_predict in clf_no_pca.staged_predict(model_tr1):
20     train_macro1.append(f1_score(y_train, real_train_predict, average='macro'))
```

```

In [22]: 1 plt.figure()
          2 plt.plot(range(1, n_trees_real + 1),
          3               real_val_macro1[1:],
          4               label='val')
          5 plt.plot(range(1, n_trees_real + 1),
          6               train_macro1,
          7               label='train')
          8 plt.legend()
          9 plt.title('macro score vs num of decision stump')
         10 adaboost_opt = AdaBoostClassifier(n_estimators=np.argmax(real_val_macro[1:])+1,
         11                                   random_state=1)
         12

```



```

In [39]: 1 np.argmax(real_val_macro1)

```

Out[39]: 577

```

In [40]: 1 clf_no_pca_opt = AdaBoostClassifier(n_estimators=577, random_state=1)
          2 clf_no_pca_opt.fit(model_tr1, y_train)

```

Out[40]: AdaBoostClassifier(algorithm='SAMME.R', base\_estimator=None, learning\_rate=1.0, n\_estimators=577, random\_state=1)

```
In [42]: 1 w2v_adaboost = clf_no_pca_opt
2 print('TF-IDF feature: Adaboost, best F1-macro score on validation set:', f1_score(y_val, w2v_adaboost.predict(model_val1)))
3 mat = confusion_matrix(y_val, w2v_adaboost.predict(model_val1))
4 print('confusion matrix on validation set is: \n', mat)
5 print('F1 score of each class:\n', f1_score(y_val, w2v_adaboost.predict(model_val1), average=None))
6 print('accuracy on validation set:', w2v_adaboost.score(model_val1, y_val), '\n\n')
```

TF-IDF feature: Adaboost, best F1-macro score on validation set: 0.29365991902671673

confusion matrix on validation set is:

```
[[ 6  4  8 16 31]
 [ 2  8  4 20 30]
 [ 7  8 28 70 127]
 [12  7 22 251 1075]
 [28 11 39 329 3397]]
```

F1 score of each class:

```
[0.1      0.15686275 0.16422287 0.24452021 0.80269376]
```

accuracy on validation set: 0.6660649819494585

```
In [25]: 1 joblib.dump(adaboost_opt, '../input/model/adaboost_nopca.pkl')
```

Out[25]: ['../input/model/adaboost\_nopca.pkl']

```
In [45]: 1 from sklearn.linear_model import LogisticRegression
2 n_alphas = 41
3 #alphas = np.exp(np.linspace(-5, 10, n_alphas))
4 C = np.exp(np.linspace(-5, 10, n_alphas))
5 a = np.linspace(0.2, 0.7, 6)
6 train_score = [], [], [], [], [], []
7 val_score = [], [], [], [], [], []
8 for n, u in enumerate(a):
9     for j, i in enumerate(C):
10         m = LogisticRegression(penalty='elasticnet', C=i, class_weight='balanced', n_jobs=-1, l1_ratio=u, solver='saga').fit(model_tr, y_train)
11         train_pre = m.predict(model_tr)
12         val_pre = m.predict(model_val)
13         train_score[n].append(f1_score(y_train, train_pre, average='macro'))
14         val_score[n].append(f1_score(y_val, val_pre, average='macro'))
15
```

```
In [60]: 1 elastic_opt = LogisticRegression(penalty='elasticnet',
2                                     C=C[(np.argmax(val_score))%41],
3                                     class_weight='balanced', n_jobs=-1,
4                                     l1_ratio=a[(np.argmax(val_score))//41],
5                                     solver='saga')
6 elastic_opt.fit(model_tr, y_train)
7 print('the best macro score on validation set:', f1_score(y_val, elastic_opt.predict(model_val), average='macro'))
8 mat = confusion_matrix(y_val, elastic_opt.predict(model_val))
9 print('confusion matrix is: \n', mat)
10 print('accuracy of each class:\n', f1_score(y_val, elastic_opt.predict(model_val), average=None))
```

the best macro score on validation set: 0.31207363766143204

confusion matrix is:

```
[[ 30  14   4   5  12]
 [ 20  18   6  11   9]
 [ 32  34  39  92  43]
 [ 50  76  99 566 576]
 [102 113  95 1019 2475]]
```

accuracy of each class:

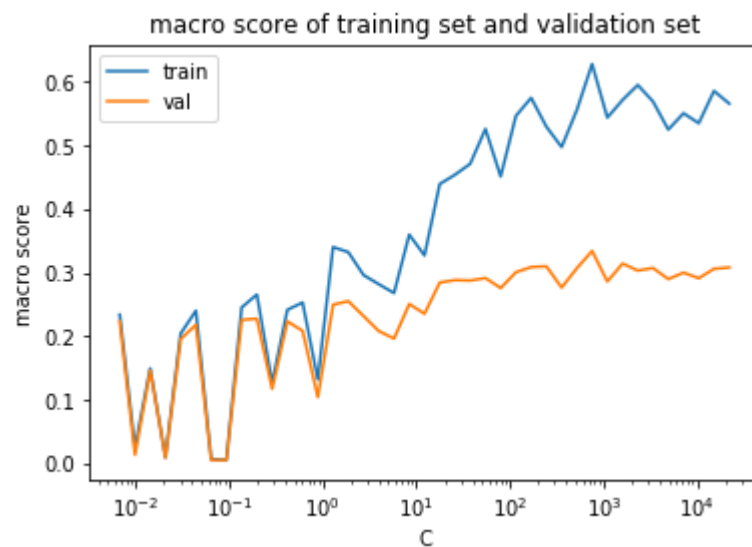
```
[0.2006689  0.11285266 0.16149068 0.36993464 0.7154213 ]
```



```
In [62]: 1 joblib.dump(elastic_opt, '../input/model/logistic_elsticnet.pkl')
```

```
Out[62]: ['../input/model/logistic_elsticnet.pkl']
```

```
In [63]: 1 train_score = np.array(train_score)
2 val_score = np.array(val_score)
3 ax1 = plt.gca()
4 # randomly choose 50 coefficients in the 9676 variables
5 # to plot figure of coefs vs regularizer.
6 ax1.plot(C, train_score[(np.argmax(val_score))//41], label='train')
7 ax1.plot(C, val_score[(np.argmax(val_score))//41], label='val')
8 ax1.set_xscale('log')
9 #ax1.set_xlim(ax1.get_xlim()[::-1]) # reverse axis
10 plt.xlabel('C')
11 plt.ylabel('macro score')
12 plt.title('macro score of training set and validation set')
13 plt.axis('tight')
14 plt.legend()
15 plt.show()
```





```

In [1]: 1 import pandas as pd
2 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
3 from matplotlib import pyplot as plt
4 from sklearn.metrics import f1_score
5 from sklearn import linear_model
6 from sklearn.metrics import confusion_matrix
7 import numpy as np
8 from sklearn.externals import joblib
9 #import warnings
10 #warnings.filterwarnings("ignore")
11 df_train = pd.read_csv('../input/train.csv')
12 df_y_train = pd.read_csv('../input/y_train.csv')
13 df_val = pd.read_csv('../input/val.csv')
14 df_y_val = pd.read_csv('../input/y_val.csv')
15 train = df_train['review'].to_numpy()
16 y_train = df_y_train['rate'].to_numpy()
17 val = df_val['review'].to_numpy()
18 y_val = df_y_val['rate'].to_numpy()
19
20 # the data set is merge here for crossvalidation in the Bayesian Inference analysis.
21 '''
22 train_total = pd.DataFrame({'review':train.tolist()+val.tolist()})
23 train_y_total = pd.DataFrame({'rate':y_train.tolist()+y_val.tolist()})
24 train_total.to_csv('../input/train_total.csv')
25 train_y_total.to_csv('../input/train_y_total.csv')
26 '''
27
28 vectorizer = TfidfVectorizer(stop_words='english', use_idf=True)
29 model_tr = vectorizer.fit_transform(train)
30 model_val = vectorizer.transform(val)

```

/home/zixi/anaconda3/lib/python3.7/site-packages/sklearn/externals/joblib/\_\_init\_\_.py:15: DeprecationWarning: sklearn.externals.joblib is deprecated in 0.21 and will be removed in 0.23. Please import this functionality directly from joblib, which can be installed with: pip install joblib. If this warning is raised when loading pickled models, you may need to re-serialize those models with scikit-learn 0.21+.

```
warnings.warn(msg, category=DeprecationWarning)
```

```
In [14]: 1 from sklearn.decomposition import PCA
2         #from sklearn.preprocessing import normalize
3         #model_tr_normalize = normalize(model_tr, norm='l2', axis=0, copy=True, return_norm=False)
4         pca = PCA(n_components = 0.95)
5         pca.fit(model_tr.toarray())
6         reduced = pca.transform(model_tr.toarray())
7         print(reduced.shape)
```

(22159, 3192)

```
In [2]: 1 from sklearn.decomposition import PCA
2         model_tr1 = model_tr
3         model_val1 = model_val
4         #pca = joblib.load('../input/model/pca.pkl')
5         model_val = pca.transform(model_val.toarray())
6         model_tr = pca.transform(model_tr.toarray())
7
```

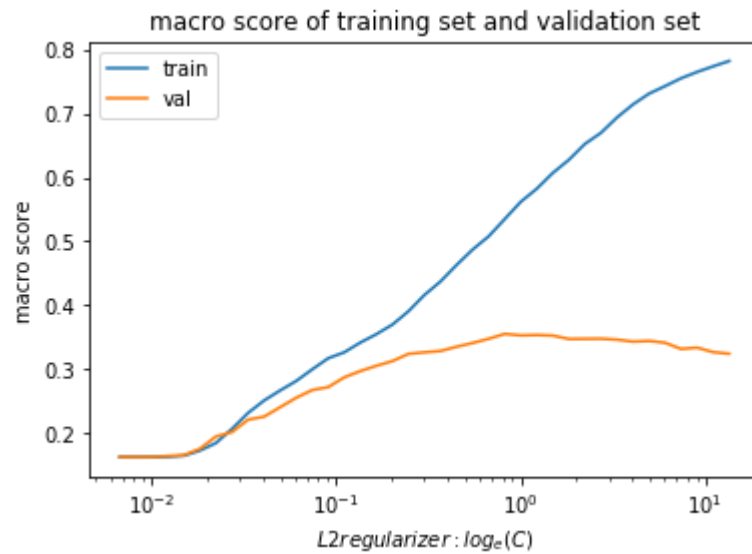
```

In [5]: 1 # Logistic regression, L1 regularization, weighted loss, hyperparameter selection use validation set, macro-weighted score
2 # reason: we assume that some of the variable which is corresponding to the word contribute little to emotion(rating)
3 # and we give more weight to the minor class to remedy for the class imbalance.
4 import warnings
5 from matplotlib import pyplot as plt
6 from sklearn.metrics import f1_score
7 from sklearn.metrics import confusion_matrix
8 from sklearn.linear_model import LogisticRegression
9 warnings.filterwarnings("ignore")
10 n_alphas = 76
11 #alphas = np.exp(np.linspace(-5, 10, n_alphas))
12 C = np.exp(np.linspace(-5, 10, n_alphas))
13 train_score = []
14 val_score = []
15 for i in C:
16     m = LogisticRegression(penalty='l1', C=i, class_weight='balanced').fit(model_tr, y_train)
17     s = f1_score(y_val, m.predict(model_val), average='macro')
18     train_score.append(f1_score(y_train, m.predict(model_tr), average='macro'))
19     val_score.append(f1_score(y_val, m.predict(model_val), average='macro'))
20     print('Val macro score is', s)
21 #plt.plot(np.log10(C), train_score, label='train')
22 #plt.plot(np.log10(C), val_score, label='val')
23 #plt.title('macro score of training set and validation set')
24 #plt.legend()
25 #plt.xlabel('log10(C)')
26 #plt.ylabel('macro score')
27 #plt.show()
28 #-----

```

...

```
In [6]: 1 from sklearn.metrics import confusion_matrix
2 ax = plt.gca()
3 # randomly choose 50 coefficients in the 9676 variables
4 # to plot figure of coefs vs regularizer.
5 ax.plot(C[:len(train_score)], train_score, label='train')
6 ax.plot(C[:len(train_score)], val_score, label='val')
7 ax.set_xscale('log')
8 #ax.set_xlim(ax.get_xlim()[::-1]) # reverse axis
9 plt.xlabel('$L2$ regularizer: log_e(C)$')
10 plt.ylabel('macro score')
11 plt.legend()
12 plt.title('macro score of training set and validation set')
13 plt.axis('tight')
14 plt.show()
15 C_opt = C[np.argmax(val_score)]
16 m_opt = LogisticRegression(C=C_opt, class_weight='balanced').fit(model_tr, y_train)
17 s = f1_score(y_val, m_opt.predict(model_val), average='macro')
18 mat = confusion_matrix(y_val, m_opt.predict(model_val))
19 print('Best Val macro score is', s)
20 print('confusion matrix is: \n', mat)
21 print('recall of each class:\n', [round(mat[0,0]/mat[0,:].sum(), 4), round(mat[1,1]/mat[1,:].sum(), 4),
22                                     round(mat[2,2]/mat[2,:].sum(), 4), round(mat[3,3]/mat[3,:].sum(), 4),
23                                     round(mat[4,4]/mat[4,:].sum(), 4)])
24 joblib.dump(m_opt, '../input/model/logistic_L1.pkl')
```



Best Val macro score is 0.3626043565870517

confusion matrix is:

```
[[ 26   8  16   3  12]
 [ 16  14  18   8   8]
 [ 22  22  73  54  69]
 [ 23  29 159 369 787]
 [ 37  42 120 400 3205]]
```

recall of each class:

```
[0.4, 0.2188, 0.3042, 0.2699, 0.8425]
```

Out[6]: ['../input/model/logistic\_L2.pkl']

```
In [60]: 1 import pandas as pd
2 import numpy as np
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import f1_score
5
6 df_train = pd.read_csv('../input/train.csv')
7 train = df_train['review'].to_numpy()
8 df_y_train = pd.read_csv('../input/y_train.csv')
9 y_train = df_y_train['rate'].to_numpy()
10 df_val = pd.read_csv('../input/val.csv')
11 val = df_val['review'].to_numpy()
12 df_y_val = pd.read_csv('../input/y_val.csv')
13 y_val = df_y_val['rate'].to_numpy()
14
15 import logging
16 logging.basicConfig(format='%(asctime)s : %(levelname)s : %(message)s', level=logging.INFO)
17 from gensim.test.utils import datapath
18 from gensim import utils
19 from pathlib import Path
20 import gensim
21 model = gensim.models.Word2Vec.load('../embedding/word2vec.model')
22 ave_vec = np.zeros((train.shape[0], 50), dtype='float')
23 word_vectors = model.wv
24 len(word_vectors.vocab)
25 for k, seq in enumerate(train):
26     tokens = gensim.utils.simple_preprocess(seq)
27     for i in tokens:
28         if i in word_vectors.vocab:
29             ave_vec[k] += model.wv[i]
30     ave_vec[k] /= len(tokens)
31 print(ave_vec.shape)
```

(22159, 50)



```
In [61]: 1 word_vectors = model.wv
          2 len(word_vectors.vocab)
          3 ave_vec_val = np.zeros((val.shape[0], 50), dtype='float')
          4 for k, seq in enumerate(val):
          5     tokens = gensim.utils.simple_preprocess(seq)
          6     for i in tokens:
          7         if i in word_vectors.vocab:
          8             ave_vec_val[k] += model.wv[i]
          9     ave_vec_val[k] /= np.array(len(tokens))
         10 print(ave_vec_val.shape)
```

(5540, 50)

```

In [97]: 1 # ridge regression (L2 regularizer)
2 #
3 print(Counter(y_train))
4 print(Counter(y_val))
5 stra = {5.0:15214,
6         4.0:10000,
7         3.0:2000,
8         2.0:1000,
9         1.0:1000}
10 sm = SMOTE(random_state=42, sampling_strategy=stra)
11 X_res, y_res = sm.fit_resample(ave_vec, y_train)
12 print(Counter(y_res))
13 from sklearn import linear_model
14 n_alphas = 50
15 alphas = np.exp(np.linspace(-5, 13, n_alphas))
16 coefs = []
17 train_score = []
18 val_score = []
19 for a in alphas:
20     ridge = linear_model.Ridge(alpha=a, fit_intercept=True)
21     ridge.fit(X_res, y_res)
22     train_pre = ridge.predict(ave_vec).astype('int')
23     train_pre[train_pre > 5] = np.int(5)
24     train_pre[train_pre < 0] = np.int(1)
25     val_pre = ridge.predict(ave_vec_val).astype('int')
26     val_pre[val_pre > 5] = np.int(5)
27     val_pre[val_pre < 1] = np.int(1)
28     train_score.append(f1_score(y_train, train_pre, average='macro'))
29     val_score.append(f1_score(y_val, val_pre, average='macro'))
30     coefs.append(ridge.coef_)
31
32 cc = np.array(coefs)
33
34 ax = plt.gca()
35 # randomly choose 50 coefficients in the 9676 variables
36 # to plot figure of coefs vs regularizer.
37 ax.plot(alphas, cc)
38 ax.set_xscale('log')
39 ax.set_xlim(ax.get_xlim()[::-1]) # reverse axis
40 plt.xlabel('alpha')
41 plt.ylabel('weights')
42 plt.title('Ridge coefficients as a function of the regularization')

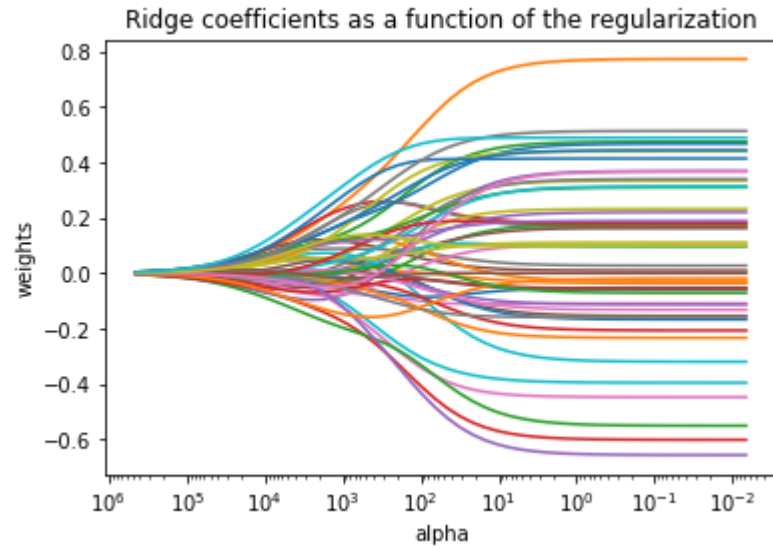
```

```
43 plt.axis('tight')  
44 plt.show()
```

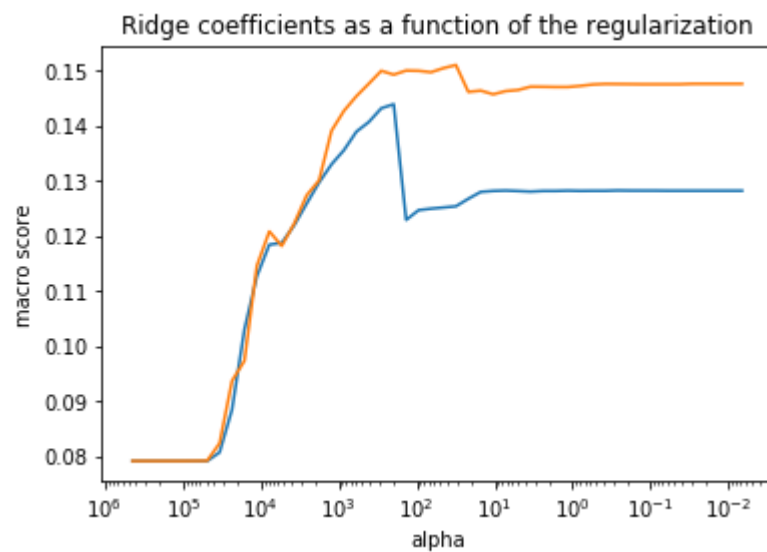
```
Counter({5.0: 15214, 4.0: 5466, 3.0: 959, 1.0: 263, 2.0: 257})
```

```
Counter({5.0: 3804, 4.0: 1367, 3.0: 240, 1.0: 65, 2.0: 64})
```

```
Counter({5.0: 15214, 4.0: 10000, 3.0: 2000, 2.0: 1000, 1.0: 1000})
```



```
In [98]: 1 train_score = np.array(train_score)
2 val_score = np.array(val_score)
3 ax1 = plt.gca()
4 # randomly choose 50 coefficients in the 9676 variables
5 # to plot figure of coefs vs regularizer.
6 ax1.plot(alphas, train_score)
7 ax1.plot(alphas, val_score)
8 ax1.set_xscale('log')
9 ax1.set_xlim(ax1.get_xlim()[::-1]) # reverse axis
10 plt.xlabel('alpha')
11 plt.ylabel('macro score')
12 plt.title('Ridge coefficients as a function of the regularization')
13 plt.axis('tight')
14 plt.show()
15
16 alphas_best = alphas[np.argmax(val_score)]
17 print('the best macro score on validation set:', np.max(val_score))
18 ridge = linear_model.Ridge(alpha=alphas_best, fit_intercept=True).fit(X_res, y_res)
19 val_pre = ridge.predict(ave_vec_val).astype('int')
20 val_pre[val_pre > 5] = np.int(5)
21 val_pre[val_pre < 1] = np.int(1)
22 mat = confusion_matrix(y_val, val_pre)
23 print('confusion matrix is: \n', mat)
24 print('accuracy of each class:\n', [round(mat[0,0]/mat[0,:].sum(), 4), round(mat[1,1]/mat[1,:].sum(), 4),
25                                     round(mat[2,2]/mat[2,:].sum(), 4), round(mat[3,3]/mat[3,:].sum(), 4),
26                                     round(mat[4,4]/mat[4,:].sum(), 4)])
27 print('the number of wi equal to 0 =', len(np.where(np.abs(ridge.coef_.flatten())==0)[0]))
28
```



the best macro score on validation set: 0.15098002820153628

confusion matrix is:

```
[[ 0  1  49  15  0]
 [ 0  1  54  9  0]
 [ 0  6 157  74  3]
 [ 0  4 424 878 61]
 [ 0  6 497 2927 374]]
```

accuracy of each class:

```
[0.0, 0.0156, 0.6542, 0.6423, 0.0983]
```

the number of wi equal to 0 = 0

```
In [99]: 1 jolib.dump(ridge, '../input/model_w2v/w2v_ridge.pkl')
```

```
Out[99]: ['../input/model_w2v/w2v_ridge.pkl']
```

```

In [80]: 1 from sklearn.tree import DecisionTreeClassifier
          2 import warnings
          3 warnings.filterwarnings("ignore")
          4 clf = DecisionTreeClassifier(random_state=0, criterion='entropy',
          5                             splitter='best' )
          6
          7 depth = np.linspace(3, 183, 61).astype('int')
          8 train_score1 = []
          9 val_score1 = []
         10 for i in depth:
         11     m = DecisionTreeClassifier(random_state=0,
         12                             criterion='entropy',
         13                             splitter='best',
         14                             max_depth=i,
         15                             min_samples_leaf=5,
         16                             min_samples_split=10,
         17                             class_weight='balanced'
         18                             )
         19     m.fit(X_res, y_res)
         20     s = f1_score(y_val, m.predict(ave_vec_val), average='macro')
         21     train_score1.append(f1_score(y_train, m.predict(ave_vec), average='macro'))
         22     val_score1.append(f1_score(y_val, m.predict(ave_vec_val), average='macro'))
         23     #print ('Val macro-weighted score is', s)
         24
         25 train_score1 = np.array(train_score1)
         26 val_score1 = np.array(val_score1)
         27
         28 min_leaf = np.linspace(1, 101, 51).astype('int')
         29 train_score2 = []
         30 val_score2 = []
         31 for i in min_leaf:
         32     m = DecisionTreeClassifier(random_state=0,
         33                             criterion='entropy',
         34                             splitter='best',
         35                             max_depth=depth[np.argmax(val_score1)],
         36                             min_samples_leaf=i,
         37                             min_samples_split=10,
         38                             class_weight='balanced'
         39                             )
         40     m.fit(X_res, y_res)
         41     s = f1_score(y_val, m.predict(ave_vec_val), average='macro')
         42     train_score2.append(f1_score(y_train, m.predict(ave_vec), average='macro'))

```

```

43     val_score2.append(f1_score(y_val, m.predict(ave_vec_val), average='macro'))
44
45 min_samples_spl = np.linspace(5, 50, 17).astype('int')
46 train_score3 = []
47 val_score3 = []
48 for i in min_samples_spl:
49     m = DecisionTreeClassifier(random_state=0,
50                               criterion='entropy',
51                               splitter='best',
52                               max_depth=depth[np.argmax(val_score1)],
53                               min_samples_leaf=min_leaf[np.argmax(val_score2)],
54                               min_samples_split=i,
55                               class_weight='balanced'
56                               )
57     m.fit(X_res, y_res)
58     s = f1_score(y_val, m.predict(ave_vec_val), average='macro')
59     train_score3.append(f1_score(y_train, m.predict(ave_vec), average='macro'))
60     val_score3.append(f1_score(y_val, m.predict(ave_vec_val), average='macro'))
61
62 m_opt = DecisionTreeClassifier(random_state=0,
63                               criterion='entropy',
64                               splitter='best',
65                               max_depth=depth[np.argmax(val_score1)],
66                               min_samples_leaf=min_leaf[np.argmax(val_score2)],
67                               min_samples_split=min_samples_spl[np.argmax(val_score3)],
68                               class_weight='balanced'
69                               )
70 m_opt.fit(X_res, y_res)
71 mat = confusion_matrix(y_val, m_opt.predict(ave_vec_val))
72 print('optimal score', f1_score(y_val, m_opt.predict(ave_vec_val), average='macro'))
73 print('confusion matrix is: \n', mat)
74 print('accuracy of each class:\n', [round(mat[0,0]/mat[0,:].sum(), 4), round(mat[1,1]/mat[1,:].sum(), 4),
75                                     round(mat[2,2]/mat[2,:].sum(), 4), round(mat[3,3]/mat[3,:].sum(), 4),
76                                     round(mat[4,4]/mat[4,:].sum(), 4)])
77
78

```

optimal score 0.23430500637813695

confusion matrix is:

```

[[ 4  3  8 21 29]
 [ 7  4 21 13 19]
 [11 14 39 94 82]
 [48 48 133 477 661]

```

```
[ 73  66 237 1118 2310]]  
accuracy of each class:  
[0.0615, 0.0625, 0.1625, 0.3489, 0.6073]
```

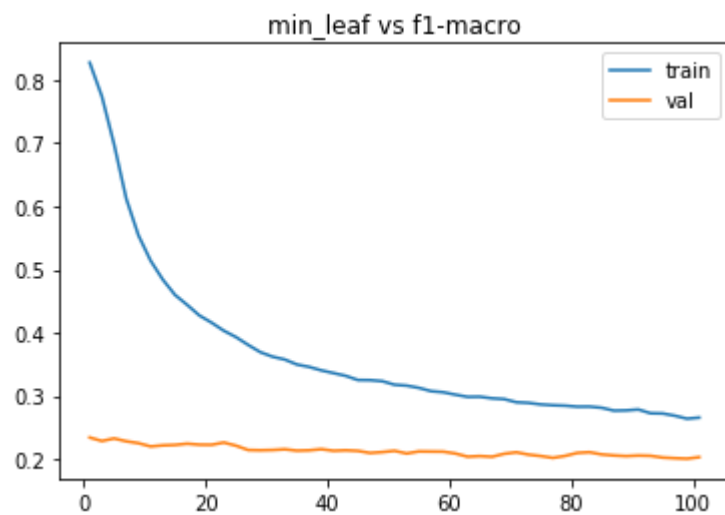
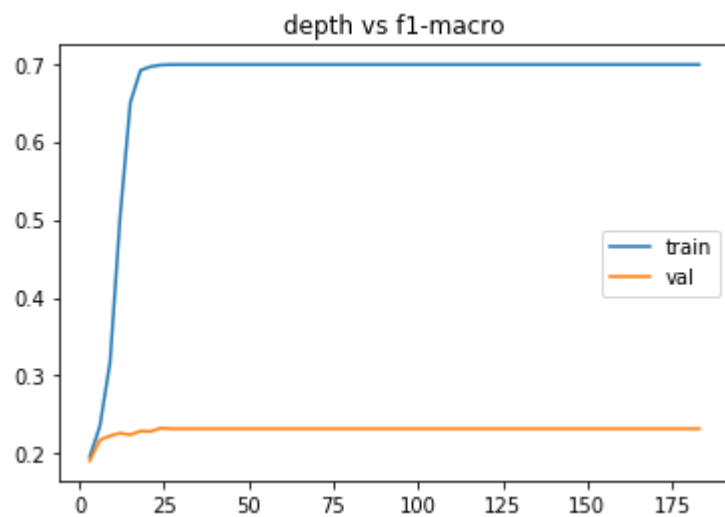
```
In [81]: 1 mat = confusion_matrix(y_val, m_opt.predict(ave_vec_val))  
2 print('optimal score', f1_score(y_val, m_opt.predict(ave_vec_val), average='macro'))  
3 print('confusion matrix is: \n', mat)  
4 print('f1 of each class:\n', f1_score(y_val, m_opt.predict(ave_vec_val), average=None))  
5
```

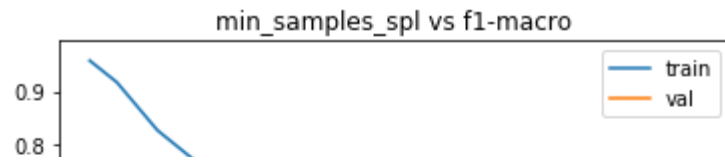
```
optimal score 0.23430500637813695  
confusion matrix is:  
[[ 4  3  8 21 29]  
 [ 7  4 21 13 19]  
 [11 14 39 94 82]  
 [48 48 133 477 661]  
 [73 66 237 1118 2310]]  
f1 of each class:  
[0.03846154 0.04020101 0.11504425 0.30873786 0.66908038]
```



```
In [82]: 1 plt.figure()
2         plt.plot(depth, train_score1, label='train')
3         plt.plot(depth, val_score1, label='val')
4         plt.legend()
5         plt.title('depth vs f1-macro')
6         print(depth[np.argmax(val_score1)])
7         print(np.max(val_score1))
8         pass
9         plt.figure()
10        plt.plot(min_leaf, train_score2, label='train')
11        plt.plot(min_leaf, val_score2, label='val')
12        plt.legend()
13        plt.title('min_leaf vs f1-macro')
14        print(min_leaf[np.argmax(val_score2)])
15        print(np.max(val_score2))
16        pass
17        plt.figure()
18        plt.plot(min_samples_spl, train_score3, label='train')
19        plt.plot(min_samples_spl, val_score3, label='val')
20        plt.legend()
21        plt.title('min_samples_spl vs f1-macro')
22        print(min_samples_spl[np.argmax(val_score3)])
23        print(np.max(val_score3))
24        pass
```

```
24
0.23279484750733576
1
0.23430500637813695
10
0.23430500637813695
```





```
In [83]: 1 joblib.dump(m_opt, '../input/model_w2v/decision_tree_w2v.pkl')
```

```
Out[83]: ['../input/model_w2v/decision_tree_w2v.pkl']
```

```

In [71]: 1 # Adaboost with tree method
2 #import warnings
3 #warnings.filterwarnings("ignore")
4 from sklearn.ensemble import AdaBoostClassifier
5 from sklearn.metrics import f1_score
6 from sklearn.metrics import confusion_matrix
7 print(Counter(y_train))
8 print(Counter(y_val))
9 stra = {5.0:15214,
10         4.0:10000,
11         3.0:2000,
12         2.0:1000,
13         1.0:1000}
14 sm = SMOTE(random_state=42, sampling_strategy=stra)
15 X_res, y_res = sm.fit_resample(ave_vec, y_train)
16 print(Counter(y_res))
17 clf = AdaBoostClassifier(n_estimators=1000, random_state=1)
18 clf.fit(X_res, y_res)
19 real_val_macro = [0]
20 train_macro = []
21 for real_test_predict in clf.staged_predict(ave_vec_val):
22     if f1_score(y_val, real_test_predict, average='macro') > np.max(real_val_macro) :
23         pred_opt = real_test_predict
24         real_val_macro.append(
25             f1_score(y_val, real_test_predict, average='macro'))
26 for real_train_predict in clf.staged_predict(ave_vec):
27     train_macro.append(f1_score(y_train, real_train_predict, average='macro'))
28
29
30 n_trees_real = len(clf)
31 plt.figure()
32 plt.plot(range(1, n_trees_real + 1),
33         real_val_macro[1:],
34         label='val')
35 plt.plot(range(1, n_trees_real + 1),
36         train_macro,
37         label='train')

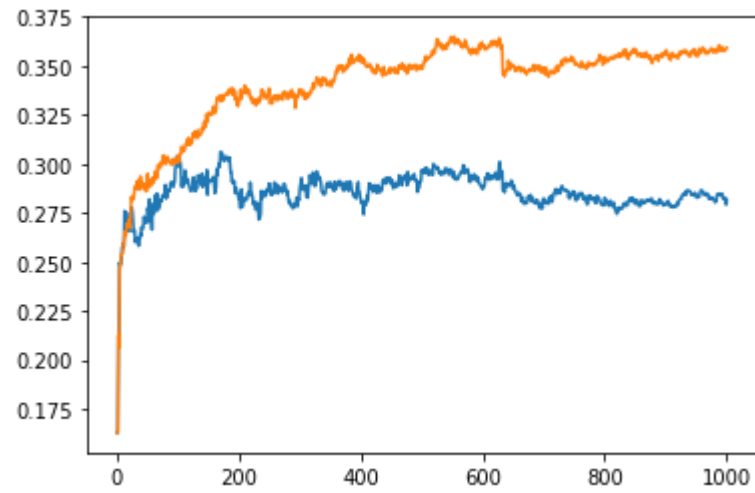
```

Counter({5.0: 15214, 4.0: 5466, 3.0: 959, 1.0: 263, 2.0: 257})

Counter({5.0: 3804, 4.0: 1367, 3.0: 240, 1.0: 65, 2.0: 64})

Counter({5.0: 15214, 4.0: 10000, 3.0: 2000, 2.0: 1000, 1.0: 1000})

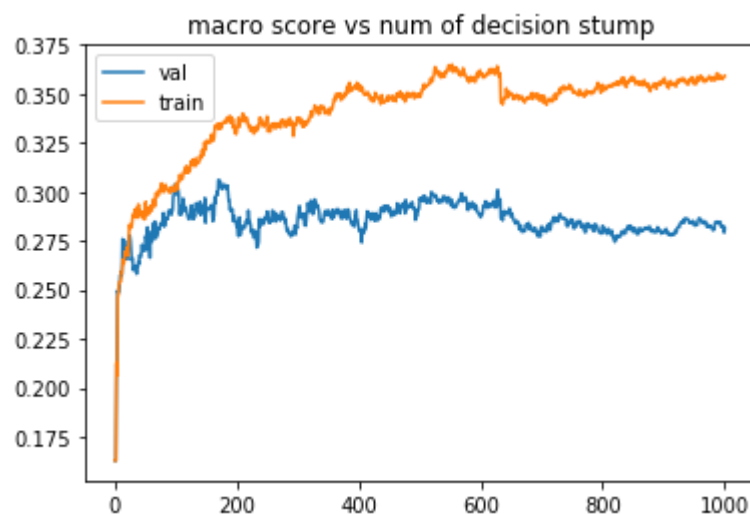
Out[71]: [<matplotlib.lines.Line2D at 0x7fcea9b37ac8>]



```

In [72]: 1 plt.figure()
          2 plt.plot(range(1, n_trees_real + 1),
          3             real_val_macro[1:],
          4             label='val')
          5 plt.plot(range(1, n_trees_real + 1),
          6             train_macro,
          7             label='train')
          8 plt.legend()
          9 plt.title('macro score vs num of decision stump')
         10 adaboost_opt = AdaBoostClassifier(n_estimators=np.argmax(real_val_macro[1:])+1,
         11                                   random_state=1)

```



```

In [73]: 1 pred_opt

```

```

Out[73]: 0.30599352334303764

```

```

In [76]: 1 adaboost_opt.fit(X_res, y_res)

```

```

Out[76]: AdaBoostClassifier(algorithm='SAMME.R', base_estimator=None, learning_rate=1.0,
                             n_estimators=170, random_state=1)

```

```
In [84]: 1 f1_score(y_val, adaboost_opt.predict(ave_vec_val), average='macro')
```

```
Out[84]: 0.30599352334303764
```

```
In [89]: 1 joblib.dump(adaboost_opt, '../input/model_w2v/w2v_adaboost.pkl')
```

```
Out[89]: ['../input/model_w2v/w2v_adaboost.pkl']
```

```
In [86]: 1 np.argmax(real_val_macro[1:])+1
```

```
Out[86]: 170
```

```
In [88]: 1 w2v_adaboost = adaboost_opt
2 print('Word2Vec feature: Adaboost, best F1-macro score on validation set:', f1_score(y_val, w2v_adaboost.predict(ave_vec_val)))
3 mat = confusion_matrix(y_val, w2v_adaboost.predict(ave_vec_val))
4 print('confusion matrix on validation set is: \n', mat)
5 print('F1 score of each class:\n', f1_score(y_val, w2v_adaboost.predict(ave_vec_val), average=None))
6 print('accuracy on validation set:', w2v_adaboost.score(ave_vec_val, y_val), '\n\n')
```

Word2Vec feature: Adaboost, best F1-macro score on validation set: 0.30599352334303764

confusion matrix on validation set is:

```
[[ 18   7   4  15  21]
 [ 15   6   2  19  22]
 [ 18  14  21 111  76]
 [ 28  32  28 429 850]
 [ 46  30  26 580 3122]]
```

F1 score of each class:

```
[0.18947368 0.07843137 0.13084112 0.34034113 0.7908803 ]
```

accuracy on validation set: 0.6490974729241877

```
In [2]: 1 import pandas as pd
2 import numpy as np
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import f1_score
5
6 df_train = pd.read_csv('../input/train.csv')
7 train = df_train['review'].to_numpy()
8 df_y_train = pd.read_csv('../input/y_train.csv')
9 y_train = df_y_train['rate'].to_numpy()
10 df_val = pd.read_csv('../input/val.csv')
11 val = df_val['review'].to_numpy()
12 df_y_val = pd.read_csv('../input/y_val.csv')
13 y_val = df_y_val['rate'].to_numpy()
14
15 import logging
16 logging.basicConfig(format='%(asctime)s : %(levelname)s : %(message)s', level=logging.INFO)
17 from gensim.test.utils import datapath
18 from gensim import utils
19 from pathlib import Path
20 import gensim
21 model = gensim.models.Word2Vec.load('../embedding/word2vec.model')
22 ave_vec = np.zeros((train.shape[0], 50), dtype='float')
23 word_vectors = model.wv
24 len(word_vectors.vocab)
25 for k, seq in enumerate(train):
26     tokens = gensim.utils.simple_preprocess(seq)
27     for i in tokens:
28         if i in word_vectors.vocab:
29             ave_vec[k] += model.wv[i]
30     ave_vec[k] /= len(tokens)
31 print(ave_vec.shape)
```

...



```
In [3]: 1 word_vectors = model.wv
2 len(word_vectors.vocab)
3 ave_vec_val = np.zeros((val.shape[0], 50), dtype='float')
4 for k, seq in enumerate(val):
5     tokens = gensim.utils.simple_preprocess(seq)
6     for i in tokens:
7         if i in word_vectors.vocab:
8             ave_vec_val[k] += model.wv[i]
9     ave_vec_val[k] /= np.array(len(tokens))
10 print(ave_vec_val.shape)
```

(5540, 50)

```
In [4]: 1 from imblearn.over_sampling import SMOTE # doctest: +NORMALIZE_WHITESPACE
2 from collections import Counter
3 n_alphas = 36
4 #alphas = np.exp(np.linspace(-5, 10, n_alphas))
5 C = np.exp(np.linspace(-5, 10, n_alphas))
6 print(Counter(y_train))
7 print(Counter(y_val))
8 stra = {5.0:15214,
9         4.0:10000,
10        3.0:2000,
11        2.0:1000,
12        1.0:1000}
13 sm = SMOTE(random_state=42, sampling_strategy=stra)
14 X_res, y_res = sm.fit_resample(ave_vec, y_train)
15 print(Counter(y_res))
16
17 train_score = []
18 val_score = []
19 for j, i in enumerate(C):
20     m = LogisticRegression(penalty='l2', C=i, class_weight='balanced', n_jobs=-1).fit(X_res, y_res)
21     train_pre = m.predict(ave_vec)
22     val_pre = m.predict(ave_vec_val)
23     train_score.append(f1_score(y_train, train_pre, average='macro'))
24     val_score.append(f1_score(y_val, val_pre, average='macro'))
25
```

...

```
In [7]: 1 from matplotlib import pyplot as plt
2 from sklearn.metrics import confusion_matrix
3 from sklearn.externals import joblib
4 ax = plt.gca()
5 ax.plot(C, train_score, label='train')
6 ax.plot(C, val_score, label='val')
7 ax.set_xscale('log')
8 ax.set_xlim(ax.get_xlim()) # reverse axis
9 plt.xlabel('C')
10 plt.ylabel('weights')
11 plt.title('logistic regression with L2 penalty(Word2Vec)')
12 plt.axis('tight')
13 plt.legend()
14 plt.show()
15 C_opt = C[np.argmax(val_score)]
16 m_opt = LogisticRegression(penalty='l2', C=C_opt, class_weight='balanced', n_jobs=-1).fit(X_res, y_res)
17 print(f1_score(y_val, m_opt.predict(ave_vec_val), average='macro'))
18 mat = confusion_matrix(y_val, m_opt.predict(ave_vec_val))
19 print('confusion matrix is: \n', mat)
20 print('f1 of each class:\n', f1_score(y_val, m_opt.predict(ave_vec_val), average=None))
21 joblib.dump(m_opt, '../input/model_w2v/logistic_L2.pkl')
```

...

```
In [9]: 1 from imblearn.over_sampling import SMOTE # doctest: +NORMALIZE_WHITESPACE
2 from collections import Counter
3 n_alphas = 36
4 #alphas = np.exp(np.linspace(-5, 10, n_alphas))
5 C = np.exp(np.linspace(-5, 10, n_alphas))
6 print(Counter(y_train))
7 print(Counter(y_val))
8 stra = {5.0:15214,
9         4.0:10000,
10        3.0:2000,
11        2.0:1000,
12        1.0:1000}
13 sm = SMOTE(random_state=42, sampling_strategy=stra)
14 X_res, y_res = sm.fit_resample(ave_vec, y_train)
15 print(Counter(y_res))
16
17 train_score = []
18 val_score = []
19 for j, i in enumerate(C):
20     m = LogisticRegression(penalty='l1', C=i, class_weight='balanced', n_jobs=-1, solver='saga', tol=0.001).fit(X_res, y_
21     train_pre = m.predict(ave_vec)
22     val_pre = m.predict(ave_vec_val)
23     train_score.append(f1_score(y_train, train_pre, average='macro'))
24     val_score.append(f1_score(y_val, val_pre, average='macro'))
```

...

```
In [11]: 1 # SMOTE with L1
2 from matplotlib import pyplot as plt
3 from sklearn.metrics import confusion_matrix
4 from sklearn.externals import joblib
5 ax = plt.gca()
6 ax.plot(C, train_score, label='train')
7 ax.plot(C, val_score, label='val')
8 ax.set_xscale('log')
9 ax.set_xlim(ax.get_xlim()) # reverse axis
10 plt.xlabel('C')
11 plt.ylabel('F1-macro')
12 plt.title('macro score of training set and validation set')
13 plt.axis('tight')
14 plt.legend()
15 plt.show()
16 C_opt = C[np.argmax(val_score)]
17 m_opt = LogisticRegression(penalty='l1', C=C_opt, class_weight='balanced', n_jobs=-1, solver='saga', tol=0.001).fit(X_res,
18 print(f1_score(y_val, m_opt.predict(ave_vec_val), average='macro'))
19 mat = confusion_matrix(y_val, m_opt.predict(ave_vec_val))
20 print('confusion matrix is: \n', mat)
21 print('F1 score of each class:\n', f1_score(y_val, m_opt.predict(ave_vec_val), average=None))
22 joblib.dump(m_opt, '../input/model_w2v/logistic_L1.pkl')
```

...

```

In [19]: 1 # elastic net
2 from imblearn.over_sampling import SMOTE # doctest: +NORMALIZE_WHITESPACE
3 from collections import Counter
4 n_alphas = 36
5 #alphas = np. exp(np. linspace(-5, 10, n_alphas))
6 C = np. exp(np. linspace(-5, 10, n_alphas))
7 print(Counter(y_train))
8 print(Counter(y_val))
9 stra = {5.0:15214,
10         4.0:10000,
11         3.0:2000,
12         2.0:1000,
13         1.0:1000}
14 sm = SMOTE(random_state=42, sampling_strategy=stra)
15 X_res, y_res = sm.fit_resample(ave_vec, y_train)
16 print(Counter(y_res))
17 a = np. linspace(0.2,0.7,6)
18 train_score = [[], [], [], [], [], []]
19 val_score = [[], [], [], [], [], []]
20 for n, u in enumerate(a):
21     for j, i in enumerate(C):
22         m = LogisticRegression(penalty='elasticnet', C=i, class_weight='balanced', n_jobs=-1, l1_ratio=u, solver='saga').fit(X_res, y_res)
23         train_pre = m.predict(ave_vec)
24         val_pre = m.predict(ave_vec_val)
25         train_score[n].append(f1_score(y_train, train_pre, average='macro'))
26         val_score[n].append(f1_score(y_val, val_pre, average='macro'))
27

```

...

```

In [29]: 1 elastic_opt = LogisticRegression(penalty='elasticnet', C=C[(np.argmax(val_score))%36], class_weight='balanced', n_jobs=-1,

```

```
In [30]: 1 elastic_opt.fit(X_res, y_res)
          2 print('the best macro score on validation set:', f1_score(y_val, elastic_opt.predict(ave_vec_val), average='macro'))
          3 mat = confusion_matrix(y_val, elastic_opt.predict(ave_vec_val))
          4 print('confusion matrix is: \n', mat)
          5 print('accuracy of each class:\n', f1_score(y_val, elastic_opt.predict(ave_vec_val), average=None))
```

/home/zixi/anaconda3/lib/python3.7/site-packages/sklearn/linear\_model/logistic.py:469: FutureWarning: Default multi\_class will be changed to 'auto' in 0.22. Specify the multi\_class option to silence this warning.

"this warning.", FutureWarning)

the best macro score on validation set: 0.2848942737772833

confusion matrix is:

```
[[ 28  20   9   2   6]
 [ 19  34   6   3   2]
 [ 32  65  76  43  24]
 [ 88 161 263 374 481]
 [242 230 286 667 2379]]
```

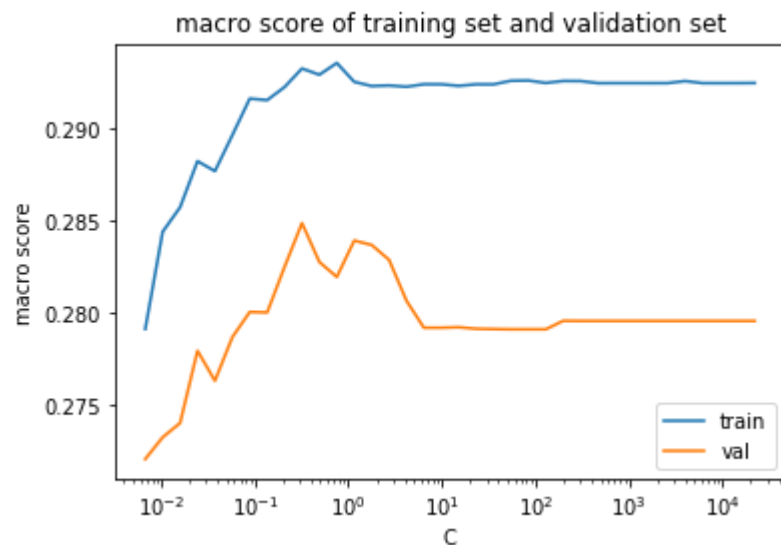
accuracy of each class:

```
[0.11814346 0.1184669 0.17272727 0.30456026 0.71057348]
```

```

In [33]: 1 train_score = np.array(train_score)
          2 val_score = np.array(val_score)
          3 ax1 = plt.gca()
          4 # randomly choose 50 coefficients in the 9676 variables
          5 # to plot figure of coefs vs regularizer.
          6 ax1.plot(C, train_score[(np.argmax(val_score))//36], label='train')
          7 ax1.plot(C, val_score[(np.argmax(val_score))//36], label='val')
          8 ax1.set_xscale('log')
          9 #ax1.set_xlim(ax1.get_xlim()[::-1]) # reverse axis
         10 plt.xlabel('C')
         11 plt.ylabel('macro score')
         12 plt.title('macro score of training set and validation set')
         13 plt.axis('tight')
         14 plt.legend()
         15 plt.show()

```



```

In [27]: 1 val_score[(np.argmax(val_score))//36][(np.argmax(val_score))%36]

```

Out[27]: 0.2848942737772833

```
In [34]: 1 joblib.dump(elastic_opt, '../input/model_w2v/logistic_elastic.pkl')
```

```
Out[34]: ['../input/model_w2v/logistic_elastic.pkl']
```

```
In [ ]: 1
```

```
In [ ]: 1
```