```
In [2]: import numpy as np
         import json
         import pandas as pd
         from sklearn import metrics
         from sklearn.model selection import KFold
         from sklearn.naive bayes import GaussianNB
         from sklearn.naive bayes import BernoulliNB
         from sklearn.linear model import LogisticRegression
 In [3]: # loading data from training data
         f = open('./whats-cooking/train.json','r')
         train data = json.loads(f.read())
         f.close()
In [57]: # loading data from testing data
         f = open('./whats-cooking/test.json','r')
         test data = json.loads(f.read())
         f.close()
In [17]: #(b) count
         dish count = len(train data)
         cuisine = []
         uniq ingredients = []
         for d in train data:
             cuisine.append(d['cuisine'])
             uniq ingredients += d['ingredients']
         uniq ingredients = np.unique(uniq ingredients)
         cuisine = np.unique(cuisine)
         cuisine count = len(cuisine)
         uniq ingredients count = len(uniq ingredients)
         print("dish count: {}".format(dish_count))
         print("unique cuisines: {}".format(cuisine count))
         print("unique ingredients: {}".format(uniq ingredients count))
         dish count: 39774
         unique cuisines: 20
         unique ingredients: 6714
```

## Out[38]:

	id	cuisine	( oz.) tomato sauce	( oz.) tomato paste	(10 oz.) frozen chopped spinach	frozen chopped spinach, thawed and squeezed dry	(14 oz.) sweetened condensed milk	(14.5 oz.) diced tomatoes	(15 oz.) refried beans	
0	10259	greek	0	0	0	0	0	0	0	
1	25693	southern_us	0	0	0	0	0	0	0	
2	20130	filipino	0	0	0	0	0	0	0	
3	22213	indian	0	0	0	0	0	0	0	
4	13162	indian	0	0	0	0	0	0	0	
39769	29109	irish	0	0	0	0	0	0	0	
39770	11462	italian	0	0	0	0	0	0	0	
39771	2238	irish	0	0	0	0	0	0	0	
39772	41882	chinese	0	0	0	0	0	0	0	
39773	2362	mexican	0	0	0	0	0	0	0	

(10 oz.)

39774 rows × 6716 columns

 $[0, 0, 0, \ldots, 0, 0, 0]]$ 

```
In [40]: #print(binary train data.shape)
         (39774, 6714)
In [46]: | #(d) Using Naïve Bayes Classifier to perform 3 fold cross-validation on
          the training set
         # Report your average classification accuracy.
         # Try both Gaussian distribution prior assumption and Bernoulli distribu
         tion prior assumption.
         gaussian = GaussianNB()
         bernoulli = BernoulliNB()
         Gaussian accuracy = []
         Bernoulli accuracy = []
         kf = KFold(n splits=3)
         for train index, test index in kf.split(binary train data):
             trainx = binary train data[train index]
             testx = binary train data[test index]
             trainy = train labels[train index]
             testy = train labels[test index]
             gaussian.fit(trainx, trainy)
             score gaussian = gaussian.score(testx, testy)
             Gaussian accuracy.append(score gaussian)
             avg_score_g = np.average(Gaussian accuracy)
             bernoulli.fit(trainx, trainy)
             score bernoulli = bernoulli.score(testx, testy)
             Bernoulli accuracy.append(score bernoulli)
             avg score b = np.average(Bernoulli accuracy)
             print("Average accuracy of Gaussian:", avg score g)
             print("Average accuracy of Bernoulli:", avg score b)
         Average accuracy of Gaussian: 0.37901644290239855
```

(e) For Gaussian prior and Bernoulli prior which performs better in terms of cross-validation accuracy? Why? Please give specific arguments.

Average accuracy of Bernoulli: 0.684190677326897 Average accuracy of Gaussian: 0.3809775230049781 Average accuracy of Bernoulli: 0.6818524664353598 Average accuracy of Gaussian: 0.3798461306381053 Average accuracy of Bernoulli: 0.6835369839593705

Average accuracy of Gaussian prior is aroung 0.38 and average accuracy of Bernoulli is around 0.68. Bernoulli performs better in terms of cross-validation accuracy. The ingredients have two features which are represented by 0 or 1 (i.e. existing or non-existing), which is better described by a Bernoulli model. Thus, Bernoulli distribution prior assumption is better.

```
In [48]: # (f) Using Logistic Regression Model to perform 3 fold cross-validation
         on the training set
         # Report your average classification accuracy
         # import warnings filter
         import warnings
         # ignore all future warnings
         warnings.simplefilter('ignore')
         logreg = LogisticRegression()
         logreg accuracy = []
         for train index, test index in kf.split(binary train data):
             trainx = binary train data[train index]
             testx = binary train data[test index]
             trainy = train labels[train index]
             testy = train labels[test index]
             logreq.fit(trainx, trainy)
             score logreg = logreg.score(testx, testy)
             logreg accuracy.append(score logreg)
             avg score 1 = np.average(logreg accuracy)
             print("Average accuracy of LogisticRegression:", avg score 1)
```

Average accuracy of LogisticRegression: 0.7758334590435964 Average accuracy of LogisticRegression: 0.773985518177704 Average accuracy of LogisticRegression: 0.7755568964650275

```
In [67]: # (g) Train your best-performed classifier with all of the training dat
    a, and generate test labels on
# test set. Submit your results to Kaggle and report the accuracy.

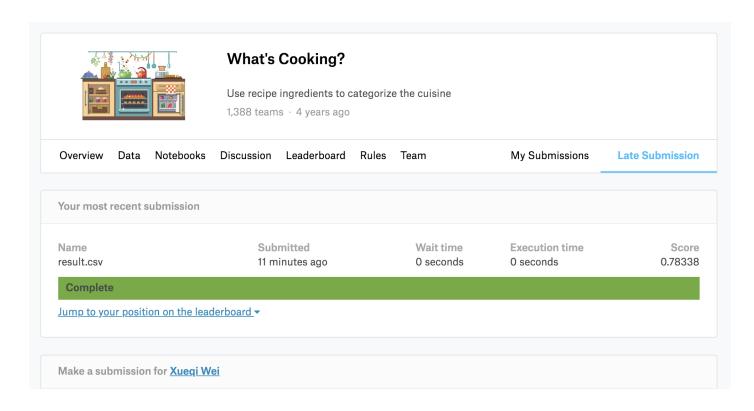
# Average accuracy of LogisticRegression is about 0.77, which is the bes
t-performed classifier
# train LogisticRegression with all of the training data

# prepare test data
binary_test_data = np.zeros([len(test_data), uniq_ingredients_count])

for i, dish in enumerate(test_data):
    for j in dish['ingredients']:
        if j in uniq_ingredients:
            binary_test_data[i][np.where(uniq_ingredients == j)] = 1
```

```
In [74]: #train model
    logreg.fit(binary_train_data, train_labels)
    pred = logreg.predict(binary_test_data)
```

```
In [75]: #output
    df = pd.DataFrame(data = {"id" : test_id, "cuisine" : pred})
    df.to_csv(path_or_buf="result.csv", index=False)
```



The accuracy on Kaggle after submission is 0.78338.