

Machine Learning Engineer Nanodegree

Capstone Project Report

Sai Marasani
October 3, 2018

Project Overview

Picture yourself strolling through your local, open-air market... What do you see? What do you smell? What will you make for dinner tonight?

If you're in Northern California, you'll be walking past the inevitable bushels of leafy greens, spiked with dark purple kale and the bright pinks and yellows of chard. Across the world in South Korea, mounds of bright red kimchi greet you, while the smell of the sea draws your attention to squids squirming nearby. India's market is perhaps the most colorful, awash in the rich hues and aromas of dozens of spices: turmeric, star anise, poppy seeds, and garam masala as far as the eye can see. Some of our strongest geographic and cultural associations are tied to a region's local foods.

In this project, I have developed supervised machine learning model which will predict the cuisine type provided with ingredients as input.

Data set link: <https://www.kaggle.com/c/whats-cooking/data>

Problem Statement:

This is Multi Class Classification problem where there are 20 possible responses from the features we pass as input. Final goal of this project is to predict cuisine type from the ingredients list provided. For out of sample performances I have considered Cross-validation scores. The ultimate goal is to accurately predict the cuisine type provided with ingredients as input. We will use accuracy measurement as a model performance metrics which gives us a percentage based on all correct and incorrect classifications of our model on the test data.

Tasks involved are:

1. Download and preprocess the JSON format data
2. Build and train the ML model to classify cuisine types
3. Our trained model will be used to predict the cuisine type provided ingredients as input.
4. Based on accuracy scores on test data, improve the model

After reading data, I will apply feature engineering techniques to create new features like number of ingredients, mean ingredient length. Ingredients are in the form of list of strings so I will convert them into a string. Then using CountVectorizer we will convert text into a matrix of token counts, Use Pipeline for proper cross validation and finally using GridSearchCV search (can be exhaustive search) and tune parameters and finally as mentioned above we will test our model on test data.

Metrics:

As the problem comes under multi classification problem and we have little imbalance in class distribution as shown in figure 2, f1_score will be my metric.

The F1 score can be interpreted as a weighted average of the precision and recall, where an F1 score reaches its best value at 1 and worst score at 0. The relative contribution of precision and recall to the F1 score are equal. The formula for the F1 score is:

$$F1 = 2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$$

Where precision is the ratio $tp / (tp + fp)$ where tp is the number of true positives and fp the number of false positives. The precision is intuitively the ability of the classifier not to label as positive a sample that is negative and recall is the ratio $tp / (tp + fn)$ where tp is the number of true positives and fn the number of false negatives. The recall is intuitively the ability of the classifier to find all the positive samples.

Analysis

Data Exploration:

“Yummly” has provided the data set which consists of 39774 samples which consists of list of ingredients which are going to be our features and cuisine column which is going to be our response variable. Data looks like below picture (Figure 1)

	cuisine	id	ingredients
0	greek	10259	[romaine lettuce, black olives, grape tomatoes...
1	southern_us	25693	[plain flour, ground pepper, salt, tomatoes, g...
2	filipino	20130	[eggs, pepper, salt, mayonaise, cooking oil, g...
3	indian	22213	[water, vegetable oil, wheat, salt]
4	indian	13162	[black pepper, shallots, cornflour, cayenne pe...

Figure 1

We have 20 different kinds of cuisines in the data set with respective value counts as shown below:

italian	7838
mexican	6438
southern_us	4320
indian	3003
chinese	2673
french	2646
cajun_creole	1546
thai	1539
japanese	1423
greek	1175
spanish	989
korean	830
vietnamese	825
moroccan	821
british	804
filipino	755
irish	667
jamaican	526
russian	489
brazilian	467

Exploratory Visualization:

Using value_counts function I have obtained the count of sample data we have for each cuisine in our training data and following plot (Figure2) shows the distribution of cuisines

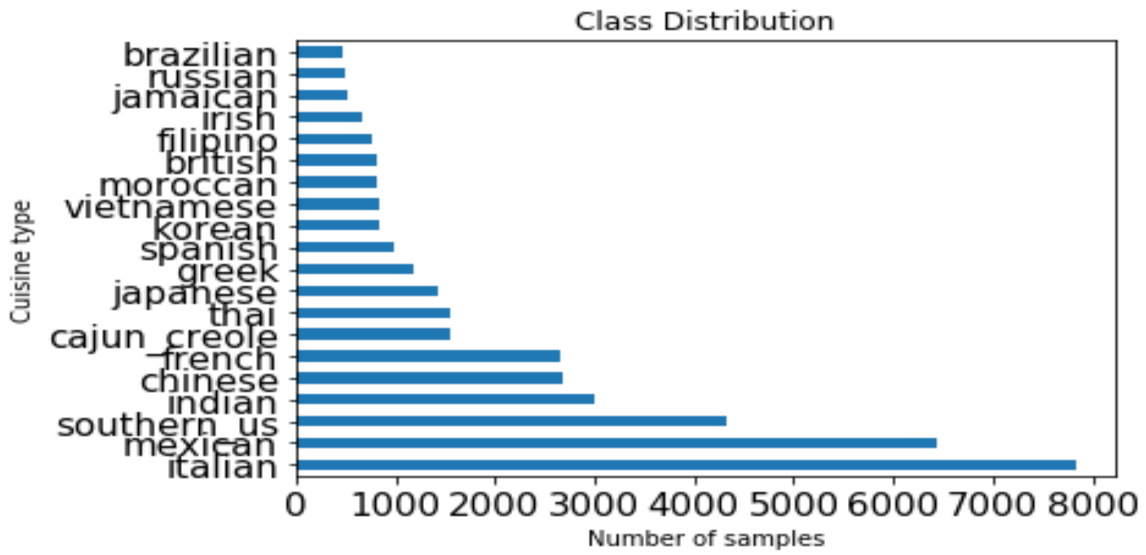


Figure 2

In order to check the variation between classes let's describe the number of ingredients by displaying count, mean, std, min, max.

```
train.groupby('cuisine').num_ingredients.describe()
```

Out[14]:

	count	mean	std	min	25%	50%	75%	max
cuisine								
brazilian	467.0	9.520343	5.555139	2.0	5.0	9.0	13.0	59.0
british	804.0	9.708955	4.165011	2.0	7.0	9.0	12.0	30.0
cajun_creole	1546.0	12.617076	4.611601	2.0	9.0	12.0	16.0	31.0
chinese	2673.0	11.982791	4.042125	2.0	9.0	12.0	14.0	38.0
filipino	755.0	10.000000	3.855135	2.0	7.0	10.0	12.0	38.0
french	2646.0	9.817838	4.144744	1.0	7.0	9.0	12.0	31.0
greek	1175.0	10.182128	3.729461	1.0	7.0	10.0	12.0	27.0
indian	3003.0	12.705961	5.016806	1.0	9.0	12.0	16.0	49.0
irish	667.0	9.299850	3.700505	2.0	7.0	9.0	12.0	27.0
italian	7838.0	9.909033	3.806708	1.0	7.0	10.0	12.0	65.0
jamaican	526.0	12.214829	4.763897	2.0	9.0	12.0	15.0	35.0
japanese	1423.0	9.735067	4.245882	1.0	7.0	9.0	12.0	34.0
korean	830.0	11.284337	3.878880	2.0	9.0	11.0	14.0	29.0
mexican	6438.0	10.877446	4.660183	1.0	7.0	10.0	14.0	52.0
moroccan	821.0	12.909866	4.799813	2.0	9.0	13.0	16.0	31.0
russian	489.0	10.224949	4.051223	2.0	7.0	10.0	13.0	25.0
southern_us	4320.0	9.634954	3.869404	1.0	7.0	9.0	12.0	40.0
spanish	989.0	10.423660	4.160919	1.0	7.0	10.0	13.0	35.0
thai	1539.0	12.545809	4.411794	1.0	9.0	12.0	15.0	40.0
vietnamese	825.0	12.675152	5.256173	1.0	9.0	12.0	16.0	31.0

Figure 3

Algorithms and Techniques:

CountVectorizer:

Text Analysis is a major application field for machine learning algorithms. However the raw data, a sequence of symbols cannot be fed directly to the algorithms themselves as most of them expect numerical feature vectors with a fixed size rather than the raw text documents with variable length.

We will use CountVectorizer to convert text into matrix of token counts where each individual token occurrence frequency will be treated as a feature. Vectorization the general process of turning a

collection of text documents into numerical feature vectors. This specific strategy (tokenization, counting and normalization) is called the Bag of Words or "Bag of n-grams" representation.

Logistic Regression:

Logistic regression classifier is more like a linear classifier which uses the calculated logits (score) to predict the target class.

Definition from Wikipedia: "Logistic regression measures the relationship between the categorical dependent variable and one or more independent variables by estimating probabilities using a logistic function"

The dependent variable is the target class variable we are going to predict. However, the independent variables are the features or attributes we are going to use to predict the target class.

Following picture is for a model which predicts how likely that penguin will be happy based on its daily activities.

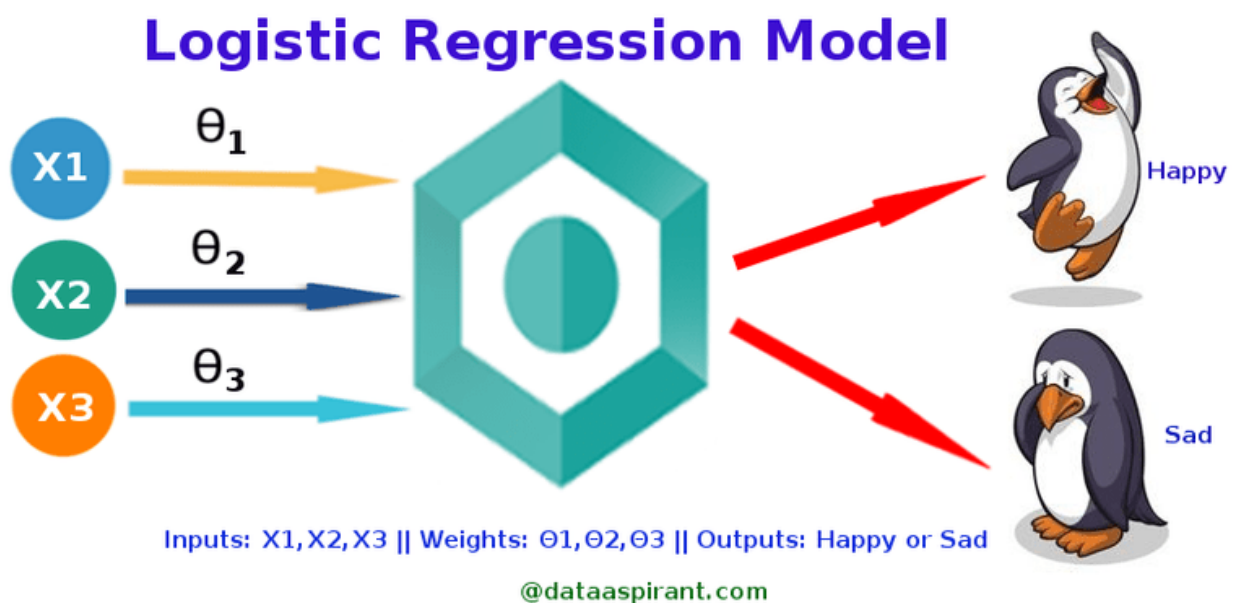


Figure 4

The "logit" model solves these problems:

$$\ln[p/(1-p)] = \alpha + \mathbf{B}X + e \text{ or}$$

$$[p/(1-p)] = \exp(\alpha + \mathbf{B}X + e)$$

where:

- \ln is the natural logarithm, \log_{\exp} , where $\exp=2.71828...$
- p is the probability that the event Y occurs, $p(Y=1)$
- $p/(1-p)$ is the "odds ratio"
- $\ln[p/(1-p)]$ is the log odds ratio, or "logit"
- all other components of the model are the same.

The logistic regression model is simply a non-linear transformation of the linear regression. The "logistic" distribution is an S-shaped distribution function which is similar to the standard-normal distribution (which results in a probit regression model) but easier to work with in most applications (the probabilities are easier to calculate). The logit distribution constrains the estimated probabilities to lie between 0 and 1.

For instance, the estimated probability is:

$$p = 1/[1 + \exp(-\alpha - \mathbf{B}X)]$$

With this functional form:

- if you let $\alpha + \mathbf{B}X = 0$, then $p = .50$
- as $\alpha + \mathbf{B}X$ gets really big, p approaches 1
- as $\alpha + \mathbf{B}X$ gets really small, p approaches 0.

In our model we will import Logistic Regression from sklearn and we will tune hyperparameters for better accuracy.

Real-World Applications: Credit Risk Analysis, measuring the success rate of the campaign, Sales (buying or not buying)

Strengths of the model: Logistic Regression is most useful in understanding the influence of several independent variables (features) on a single outcome(response) With this model we can also output the predicted probability whether the individual is willing to donate and as we are very particular about precision, predicted probability feature will help us. Adding new data is also easier using gradient descent.

Weakness of the model: Logistic Regression requires observations to be independent of each other. It may overfit the model as well as it attempts to predict the outcomes based on a set of independent variables.

GridSearchCV:

GridSearchCV implements a “fit” and a “score” method. It also implements “predict”, “predict_proba”, “decision_function”, “transform” and “inverse_transform” if they are implemented in the estimator used. The parameters of the estimator used to apply these methods are optimized by cross-validated grid-search over a parameter grid.

Pipeline:

Sequentially apply a list of transforms and a final estimator. Intermediate steps of the pipeline must be ‘transforms’, that is, they must implement fit and transform methods. The final estimator only needs to implement fit. The transformers in the pipeline can be cached using memory argument.

The purpose of the pipeline is to assemble several steps that can be cross-validated together while setting different parameters.

FeatureUnion:

This estimator applies a list of transformer objects in parallel to the input data, then concatenates the results. This is useful to combine several feature extraction mechanisms into a single transformer.

Parameters of the transformers may be set using its name and the parameter name separated by a ‘__’. A transformer may be replaced entirely by setting the parameter with its name to another transformer, or removed by setting to ‘drop’ or None.

Benchmark:

As initial baseline model, I have started with Null model as the benchmark model; null model always predicts the most frequent class from the training data.

```
y_test.value_counts().head(1) / y_test.shape
italian      0.199216
```

So, our model accuracy should be way greater than the null model accuracy which is 0.199216.

After the review comments from my proposal submission I am going to use Naïve Bayes as my benchmark model. Naïve Bayes model has produced an accuracy score of 0.67115 on the test data set.

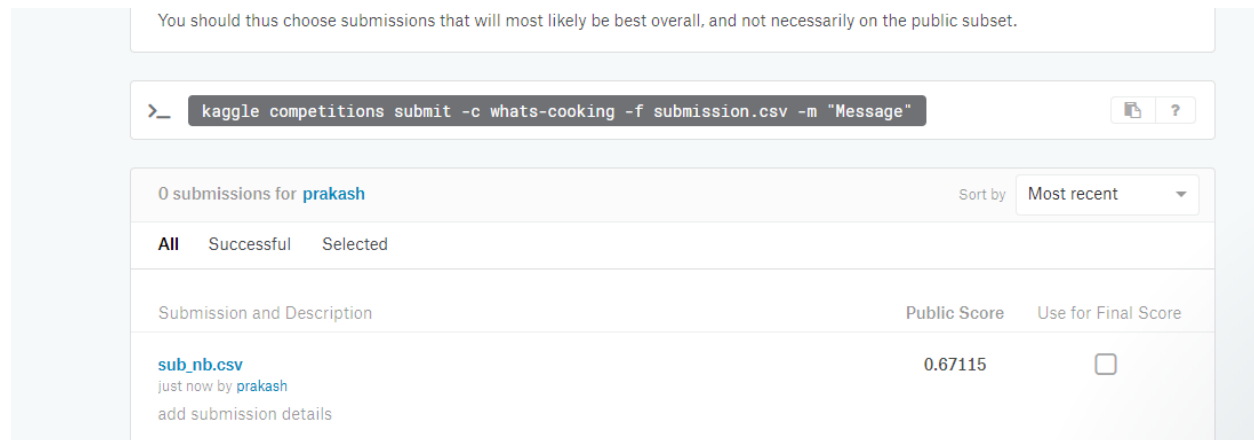


Figure 5

Methodology:

Data Preprocessing:

Preprocessing of data is carried out before the model is built and training process is executed:

Following are the steps carried out during preprocessing:

1. Dataset is divided into a training set and testing dataset.
2. Will apply feature engineering techniques to create new features like number of ingredients, mean ingredient length Ingredients are in the form of list of strings so I will convert them into a string
3. All the new features will be created for testing data set as well before we predict.

Implementation:

The first step in this process will be training a classifier with the following sub steps:

1. Reading the training data and look for null values, if we find minimal amount of null values then delete them from dataset.
2. Adding new features through Feature Engineering process
3. Adding new features to existing ones through FeatureUnion so that it will be compatible with sklearn terms once we vectorize.
4. Then use CountVectorizer to convert text into a matrix of token counts.
5. Use Pipeline for proper cross validation.
6. Using GridSearchCV search(can be exhaustive search) and tune parameters
7. Final step would be predicting on test data and calculate accuracy score, if accuracy is not high then return to step 5 to tune model or for change of hyper parameters.

Refinement:

Initially I have started with Support Vector Machines and the accuracy achieved was not satisfactory (68%) so I have tried with Logistic Regression model which resulted in 70% accuracy without using feature union and using the following parameters for grid search

```
param_grid = {}  
  
param_grid['countvectorizer__token_pattern'] = [r"\b\w\w+\b", r"([a-z ]+)"]  
  
param_grid['countvectorizer__max_df'] = [0.05, 0.06, 0.04, 0.07, 0.1]  
  
param_grid['countvectorizer__ngram_range'] = [(1,1), (1,2)]  
  
param_grid['logisticregression__C'] = [0.01, 0.1, 0.4, 0.6, 0.5, 1]
```

. Then I have concentrated my energy more on Feature Engineering through which I came to know about Feature union Process. I have started with Random Forests with different hyper parameters which resulted in accuracy of 75% and then I went back to Logistic Regression which resulted in accuracy of 78.73%.The parameter combinations I have used with feature union are:

```
param_grid = {}  
  
param_grid['featureunion__pipeline__countvectorizer__token_pattern'] = [r"\b\w\w+\b", r"([a-z ]+)"]  
  
param_grid['featureunion__pipeline__countvectorizer__max_df'] = [0.05, 0.06, 0.04, 0.07, 0.1]  
  
param_grid['featureunion__pipeline__countvectorizer__ngram_range'] = [(1,1), (1,2)]  
  
param_grid['logisticregression__C'] = [0.01, 0.1, 0.4, 0.6, 0.5, 1]
```

The best combination I have got with GridSearch are:

```
{'featureunion__pipeline__countvectorizer__token_pattern': '\\b\\w\\w+\\b'  
, 'featureunion__pipeline__countvectorizer__max_df': 0.1, 'featureunion__p  
ipeline__countvectorizer__ngram_range': (1, 2), 'logisticregression__C': 0  
.4}
```

Results

Model Evaluation and Validation:

The final hyper parameters were chosen because they performed best among the tried combinations. During model creation and development, a validation set was used to evaluate the model.

Complete Description of Final Model:

1. After Reading the training data and deep diving into the ingredients feature provided, I have come to a conclusion to add more features.
2. Adding new features through Feature Engineering process
3. Now we have new features created through step2 and in order to train a model on both (newly created and given), these features have to be combined to generate one feature matrix
4. Using FeatureUnion adding new features to the document term matrix was accomplished. Refer to the below picture on how it was achieved:

```
In [153]: from sklearn.preprocessing import FunctionTransformer
```

```
In [191]: new_features_ft = FunctionTransformer(new_features, validate=False)
```

5. Then use CountVectorizer to convert text into a matrix of token counts and then using Pipeline and grid search I have tuned and found best hyper parameters.

To verify the robustness of the model, I have used ingredients of food I prepare and few from internet and these are shown in free form visualization and the model performed pretty good.

Justification

Finally, the model is tried on the Kaggle's test data set which resulted in public score of 78.73% which is pretty higher than our null model and baseline model (Naïve Bayes). The public leaderboard score is 83.21% So our model has performed decently and this can be used for predictions.

Conclusion

Free-Form Visualization:

I have tried with some recipes' which I use to cook and then few from google to validate the model.

Prediction 1:

```
In [73]: # test_data = pd.read_json('../data/test.json')
test_data = number_of_ingredients(test_data)
test_data = mean_ingredient_length(test_data)
test_data = ingredients_as_string(test_data)

In [72]: test_data = pd.DataFrame({'ingredients': ('balsamic vinegar, pasta, basil, mozzarella')}, index=[0])

In [74]: test_data.head()

Out[74]:
```

	ingredients	num_of_ingredients	ingredient_length	ingredients_string
0	balsamic vinegar, pasta, basil, mozzarella	42	1.0	balsamic vinegar, pasta, basil, mozzarella

click to scroll output; double click to hide

```
In [32]: grid.best_estimator_

Out[32]: Pipeline(memory=None,
               steps=[('featureunion', FeatureUnion(n_jobs=1,
               transformer_list=[('pipeline', Pipeline(memory=None,
               steps=[('functiontransformer', FunctionTransformer(accept_sparse=False,
               func=<function get_ingredients at 0x00000000C2DCDD8>,
               inv_kw_args=None, inverse_func=None, kw_a...ty='l2', random_state=None, solver='liblinear', tol=0.0001,
               verbose=0, warm_start=False))])])])

In [75]: test_data_pred = grid.predict(test_data)
test_data_pred

<class 'pandas.core.series.Series'>
<class 'pandas.core.frame.DataFrame'>

Out[75]: array([u'italian'], dtype=object)
```

Figure 6

Prediction 2:

```

In [69]: # test_data = pd.read_json('../data/test.json')
test_data = number_of_ingredients(test_data)
test_data = mean_ingredient_length(test_data)
test_data = ingredients_as_string(test_data)

In [68]: test_data = pd.DataFrame({'ingredients': ('basmati rice, chicken, garam masala, egg')}, index=[0])

In [70]: test_data.head()

Out[70]:

```

	ingredients	num_of_ingredients	ingredient_length	ingredients_string
0	basmati rice, chicken, garam masala, egg	40	1.0	basmati rice, chicken, garam masala, egg

```

In [32]: grid.best_estimator_

Out[32]: Pipeline(memory=None,
steps=[('featureunion', FeatureUnion(n_jobs=1,
transformer_list=[('pipeline', Pipeline(memory=None,
steps=[('functiontransformer', FunctionTransformer(accept_sparse=False,
func=<function get_ingredients at 0x00000000C2DCDD8>,
inv_kw_args=None, inverse_func=None, kw_a...ty='l2', random_state=None, solver='liblinear', tol=0.0001,
verbose=0, warm_start=False))])),

In [71]: test_data_pred = grid.predict(test_data)
test_data_pred

<class 'pandas.core.series.Series'>
<class 'pandas.core.frame.DataFrame'>

Out[71]: array([u'indian'], dtype=object)

```

Figure 7

Prediction 3:

```

In [65]: # test_data = pd.read_json('../data/test.json')
test_data = number_of_ingredients(test_data)
test_data = mean_ingredient_length(test_data)
test_data = ingredients_as_string(test_data)

In [64]: test_data = pd.DataFrame({'ingredients': ('basil, tomatoes, oregano, garlic, parmesan cheese')}, index=[0])

In [67]: test_data.head()

Out[67]:

```

	ingredients	num_of_ingredients	ingredient_length	ingredients_string
0	basil, tomatoes, oregano, garlic, parmesan cheese	49	1.0	basil, tomatoes, oregano, garlic, parmesan cheese

```

In [32]: grid.best_estimator_

Out[32]: Pipeline(memory=None,
steps=[('featureunion', FeatureUnion(n_jobs=1,
transformer_list=[('pipeline', Pipeline(memory=None,
steps=[('functiontransformer', FunctionTransformer(accept_sparse=False,
func=<function get_ingredients at 0x00000000C2DCDD8>,
inv_kw_args=None, inverse_func=None, kw_a...ty='l2', random_state=None, solver='liblinear', tol=0.0001,
verbose=0, warm_start=False))])),

In [66]: test_data_pred = grid.predict(test_data)
test_data_pred

<class 'pandas.core.series.Series'>
<class 'pandas.core.frame.DataFrame'>

Out[66]: array([u'italian'], dtype=object)

```

Figure 8

Prediction 4:

```

In [77]: # test_data = pd.read_json('../data/test.json')
test_data = number_of_ingredients(test_data)
test_data = mean_ingredient_length(test_data)
test_data = ingredients_as_string(test_data)

In [76]: test_data = pd.DataFrame({'ingredients': ('avocadoes, beans, chipotles, tomatos')}, index=[0])

In [78]: test_data.head()

Out[78]:

```

	ingredients	num_of_ingredients	ingredient_length	ingredients_string
0	avocadoes, beans, chipotles, tomatos	36	1.0	avocadoes, beans, chipotles, tomatos

```

In [32]: grid.best_estimator_

Out[32]: Pipeline(memory=None,
steps=[('featureunion', FeatureUnion(n_jobs=1,
transformer_list=[('pipeline', Pipeline(memory=None,
steps=[('functiontransformer', FunctionTransformer(accept_sparse=False,
func=<function get_ingredients at 0x00000000C2DCDD8>,
inv_kw_args=None, inverse_func=None, kw_a...ty='l2', random_state=None, solver='liblinear', tol=0.0001,
verbose=0, warm_start=False))])),

```

```

In [79]: test_data_pred = grid.predict(test_data)
test_data_pred

<class 'pandas.core.series.Series'>
<class 'pandas.core.frame.DataFrame'>

Out[79]: array([u'filipino'], dtype=object)

```

Figure 9

Predictions were accurate for the first three ones but for the last one, model has predicted the cuisine as “Filipino” for ingredients which fall under “Mexican” cuisine.

The confusion matrix looks like

Reflection:

The process used for this project can be summarized using the following steps:

1. An initial problem of my interest and relevant, public datasets were found in Kaggle.
2. The data was downloaded and preprocessed(creating new features, FeatureUnion)
3. Benchmark model was created (Naïve Bayes model)
4. The classifier was trained using the data multiple times by tuning the parameters until the good set of hyper parameters were found.
5. The classifier was tested using the test data set and compared with public leader board.

After creating new features, I have struggled to find feature importance of the features to be used. There were instances where basic model used to perform well than complicated models(Ensembling) with more features. Then with permutation importance I was able to get more insights into feature importance and then as initial step I have appended new features I have created to the ingredients string and then did vectorization but I have found that weight of new features has become low may be due to the reasons of adding to the rest of the strings. After some research I have found that sklearn FeatureUnion can be used to accomplish this.

This is the most interesting supervised learning project I have ever worked and the way FeatureUnion was very interesting (bit complex).

Improvement:

To achieve more accuracy I think we may need bigger dataset with more data regarding Russian, Brazilian, Jamaican cuisines and few more. Other area of improvement can be use of more robust models and may be ensembling two or three models. May be we can more cuisine types rather than restricting us to 20 types to serve better purpose.

My Submissions to Kaggle:

12 submissions for prakash		Sort by	Most recent
All Successful Selected			
Submission and Description		Public Score	Use for Final Score
sub_fu_log_reg_tfidf.csv 2 hours ago by prakash add submission details		0.77916	<input type="checkbox"/>
sub_fu_log_reg_4.csv a day ago by prakash add submission details		0.78730	<input type="checkbox"/>

Figure 10

Public leaderboard:

Public LeaderboardPrivate Leaderboard

This leaderboard is calculated with all of the test data.

Raw DataRefresh





#	Δ1w	Team Name	Kernel	Team Members	Score ?	Entries	Last
1	▲2	yilihome			0.83216	31	3y
2	▲3	pavelgonchar			0.82853	49	3y
3	▼2	Kemal Atatürk			0.82461	163	3y
4	▲52	KKC			0.82441	11	

Figure 11

References:

<http://scikit-learn.org/stable/modules/generated/sklearn.pipeline.FeatureUnion.html>

<http://michelleful.github.io/code-blog/2015/06/20/pipelines/>

<http://zacstewart.com/2014/08/05/pipelines-of-featureunions-of-pipelines.html>

<https://codesachin.wordpress.com/2015/08/16/logistic-regression-for-dummies/>

<http://dataaspirant.com/2017/03/02/how-logistic-regression-model-works/>

<https://machinelearningmastery.com/logistic-regression-for-machine-learning/>

https://en.wikipedia.org/wiki/Logistic_function

http://scikit-learn.org/stable/modules/feature_extraction.html#text-feature-extraction

<http://scikit-learn.org/stable/modules/generated/sklearn.pipeline.Pipeline.html>

<http://www.appstate.edu/~whiteheadjc/service/logit/intro.htm>