

WHAT'S COOKING?

Jincai Ma

Xi'an Shiyou University, China

Introduction

The Data Source The data comes from Kaggle <https://www.kaggle.com/c/whats-cooking>.

Project Purpose Use recipe ingredients to categorize the cuisine.
Given the name of the condiment, predict the cuisine to which the dish belongs.

Related Field Name Interpretation In the dataset, including the recipe ID, the dish, and the list of ingredients for each recipe (variable length).The data is stored in JSON format.

- 1.train.json- A training set that contains the recipe ID, dish type, and ingredient list
- 2.test.json- A test set containing a recipe ID and a list of ingredients
- 3.sample_submission.csv- Properly formatted sample submission document

Feature extraction

- We convert the ingredients of the dish into a numerical feature vector.Consider that most dishes include salt, water, sugar, butter, etc,We will consider weighting the seasonings according to the occurrence times of the seasonings, that is, the more the occurrence times of the condiments, the lower the discriminability of the condiments.The feature we adopt is TF-IDF.
- We can get the top 5 characteristics:['greek','southern_us','filipino','indian','indian']
- The first five data in train_tfidf: $\begin{bmatrix} 0. & 0. & 0. & \dots & 0. & 0. & 0. \\ 0. & 0. & 0. & \dots & 0. & 0. & 0. \\ 0. & 0. & 0. & \dots & 0. & 0. & 0. \\ 0. & 0. & 0. & \dots & 0. & 0. & 0. \\ 0. & 0. & 0. & \dots & 0. & 0. & 0. \end{bmatrix}$

Data Import And Introduction

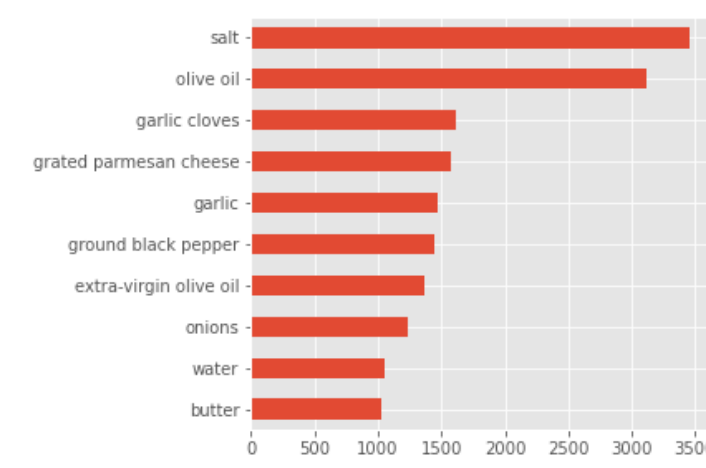
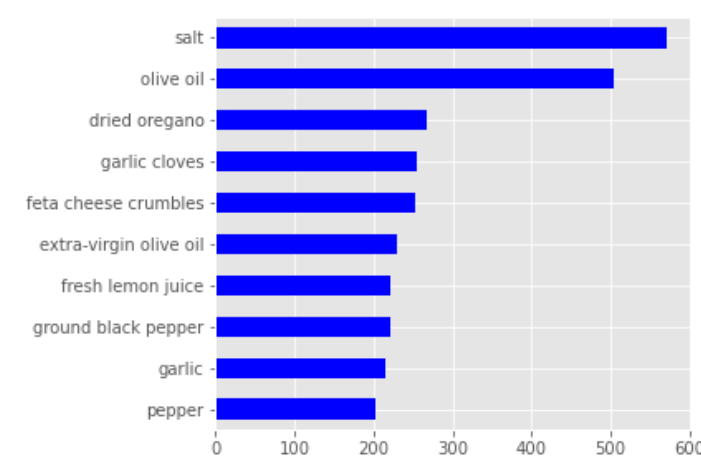
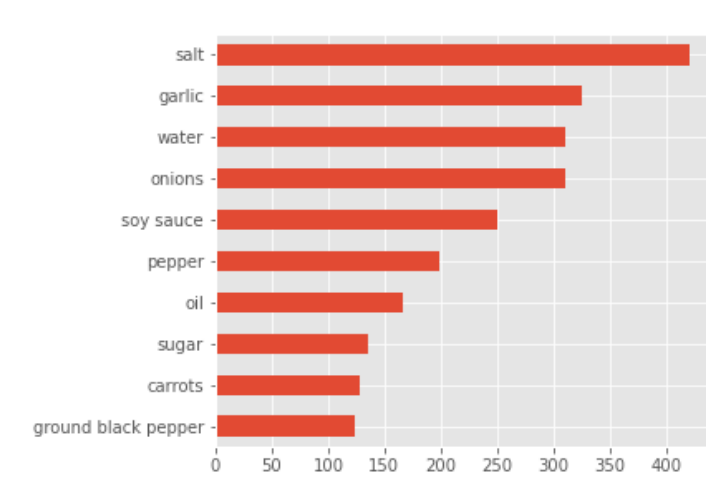
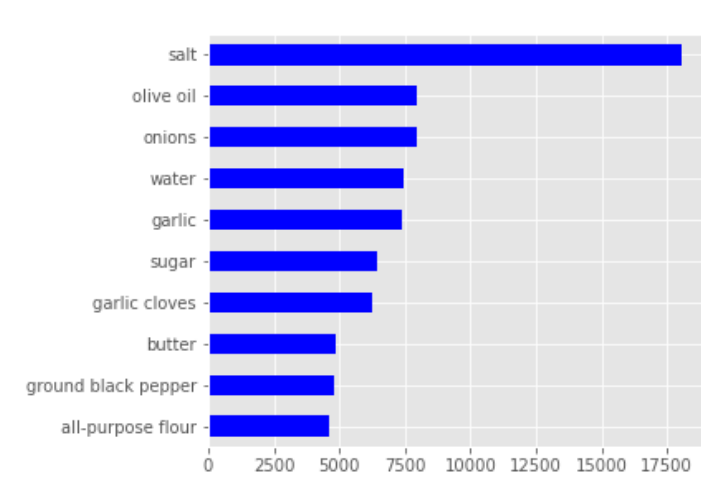
- Import the JSON file with Pandas: We can get the data set of dish names, including 39774 training data and 9944 test samples.
- To see the distribution of our data set and the total variety of dishes, we printed out some of the data samples.

	id	cuisine	ingredients
0	10259	greek	[romaine lettuce, black olives, grape tomatoes, garlic, pepper, purple onion, seasoning, garbanzo beans, feta cheese crumbles]
1	25693	southern_us	[plain flour, ground pepper, salt, tomatoes, ground black pepper, thyme, eggs, green tomatoes, yellow corn meal, milk, vegetable oil]
2	20130	filipino	[eggs, pepper, salt, mayonaise, cooking oil, green chilies, grilled chicken breasts, garlic powder, yellow onion, soy sauce, butter, chicken livers]
3	22213	indian	[water, vegetable oil, wheat, salt]
4	13162	indian	[black pepper, shallots, cornflour, cayenne pepper, onions, garlic paste, milk, butter, salt, lemon juice, water, chili powder, passata, oil, grou...
5	6602	jamaican	[plain flour, sugar, butter, eggs, fresh ginger root, salt, ground cinnamon, milk, vanilla extract, ground ginger, powdered sugar, baking powder]

- Total dish classification
There are 20 dishes in total, which are: ['brazilian' 'british' 'cajun_creole' 'chinese' 'filipino' 'french' 'greek' 'indian' 'irish' 'italian' 'jamaican' 'japanese' 'korean' 'mexican' 'moroccan' 'russian' 'southern_us' 'spanish' 'thai' 'vietnamese']

Data Analysis

- The data set is divided into Features and Target Variables.
- Features:'ingredients', we were given the names of the ingredients contained in each dish. Target variable:'cuisine', is the classification of cuisines that we want to predict.
- What are the top 10 most frequently used ingredients?
- What are the 10 most common ingredients in filipino,greek and Italian cuisine?



Validation set partitioning

- The training set is divided into a new training set and a validation set by calling train_test_split function, which is convenient for the subsequent accuracy observation of the model.
- 1.Use train_tfidf and train_targets as input variables for train_test_split
2.The test_size is set to 0.2, 20% of the validation set is divided, and 80% of the data is reserved for the new training set.
3.Set the random_state random seed to ensure that the same partition results are obtained every time you run it.

raining Model

- Invoke the logistic regression model in sklearn.
- 1.GridSearchCV is imported from sklearn.model_selection, and the parameters are automatically searched. As long as the parameters are typed in, the best results and parameters can be given.
2.Define the parameters variable: Create a dictionary for the C parameters, whose values are an array from 1 to 10;
3.Define the classifier variable: Create a classification function using the imported LogisticRegression;
4.Define Grid variables: Create a grid search object using the imported GridSearchCV;Pass the variables 'classifier', 'parameters' as arguments to the object constructor;
- After the model training, we calculated the prediction results of the model on the validation set X_VALID, and calculated the prediction accuracy of the model .
- The score on the validation set is: 0.7958516656191075.

Predictive test set

- Test set test_tfidf is predicted by the model grid, and then the predicted results are viewed.
The predicted number of test sets is 9944

id	ingredients	cuisine
0	[baking powder, eggs, all-purpose flour, raisins, milk, white sugar]	british
1	[sugar, egg yolks, corn starch, cream of tartar, bananas, vanilla wafers, milk, vanilla extract, toasted pecans, egg ...]	southern_us
2	[sausage links, fennel bulb, fronds, olive oil, cuban peppers, onions]	italian
3	[meat cuts, file powder, smoked sausage, okra, shrimp, andouille sausage, water, paprika, hot sauce, garlic cloves ...]	cajun_creole
4	[ground black pepper, salt, sausage casings, leeks, parmigiano reggiano cheese, cornmeal, water, extra-virgin olive ...]	italian
5	[baking powder, all-purpose flour, peach slices, corn starch, heavy cream, lemon juice, unsalted butter, salt, white ...]	southern_us
6	[grape juice, orange, white zinfandel]	french
7	[ground ginger, white pepper, green onions, orange juice, sugar, shiracha, vegetable oil, orange zest, chicken broth ...]	chinese
8	[diced onions, taco seasoning mix, all-purpose flour, chopped cilantro fresh, ground cumin, ground cinnamon, vegetab ...]	mexican
9	[eggs, cherries, dates, dark muscovado sugar, ground cinnamon, mixed spice, cake, vanilla extract, self-raising flour ...]	indian

Acknowledgement
• International Cooperation Project (Y7Z0511101)
of IIE, Chinese Academy of Sciences