

Applied Machine Learning HW2

1 Face Recognition

1. Read train data and test data. Display one training sample.

We first load the data with the help of the code in homework 2. Several libraries are also imported for future use. One training example is displayed.



Figure 1:

2. Average Face

We calculated the mean of the whole training set by using the mean function with axis equals 0. The following is the plot of average face.

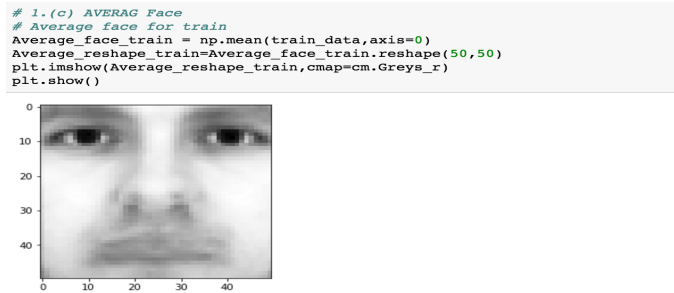


Figure 2:

3. Mean Subtraction

We subtract the mean from every column in X. And also subtract the training set mean from the test set. The following two graphs are the 10th sample from both training set and test set.

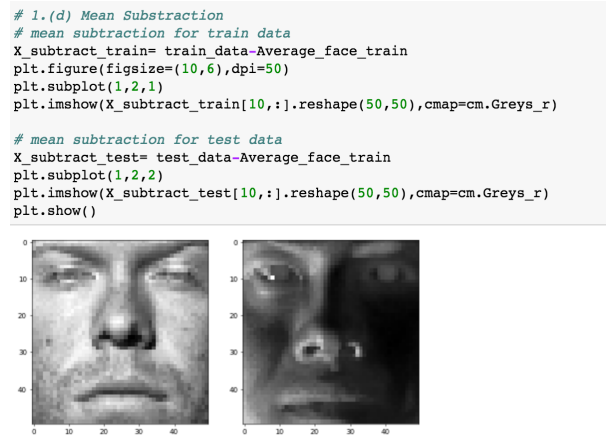


Figure 3:

4. Eigenface

We apply singular value decomposition by using function *np.linalg.svd*. The following are the 10 eigenfaces of mean subtracted training set.



Figure 4:

5. Low Rank Approximation

We defined function called rankapprox calculate $U[:, :r] * \sum[:, r, :r] * V_T[:, :r]$. Function approx error is used to calculate the rank-r approximation error of absolute value of $X - \hat{X}_r$. After we run r from 1,2,3....200, we plot the approximation error with r. In the following graph, the x axis is the value of r and the y axis is the approximation error.

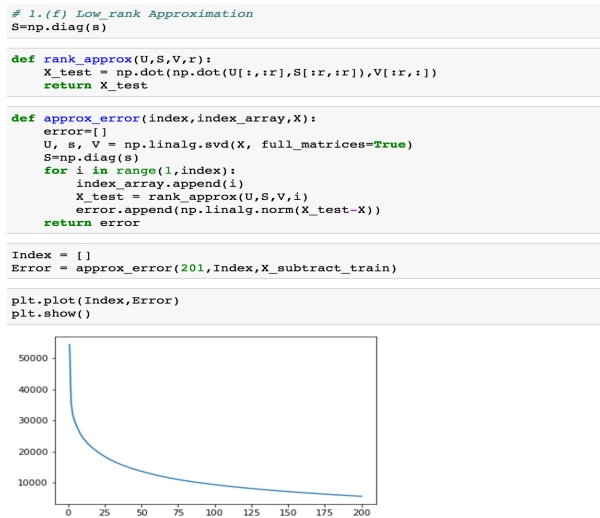


Figure 5:

From the graph, we could observe that with the increase of r , the approximation error drops with a decreasing speed. With smaller r , the approximation error drops faster.

6. Eigenface Feature

We write the function F to generate the r -dimensional feature matrix F and F_{test} for both training images X and test images X_{test} . Both F and F_{test} have the same number of rows as X and X_{test} and r columns. (we test $r=10$ in this following code)

```
# 1.(g) Eigenface Feature
def F(X,r):
    U, s, V = np.linalg.svd(X_subtract_train, full_matrices=True)
    f=np.dot(X,np.transpose(V[:r,:]))
    return f

F_train = F(X_subtract_train,10)
print(F_train.shape,X_subtract_train.shape)
F_test =F(X_subtract_test,10)
print(F_test.shape,X_subtract_test.shape)

(540, 10) (540, 2500)
(100, 10) (100, 2500)
```

Figure 6:

7. Face recognition

In this section, we extract training and test features for $r=10$ and train the logistic regression model using F and F_{test} . The classification accuracy on the test set is stable around 0.8. The following plot is the classification accuracy on the test set as a function of r when $r=1,2,3,\dots,200$. We found that after $r=25$, we already have a pretty low classification error. The approximation accuracy is stable around 0.9 after $r = 25$. This means that we could reduce the 2500 dimension feature into 25 dimensional feature with a acceptable classification error.

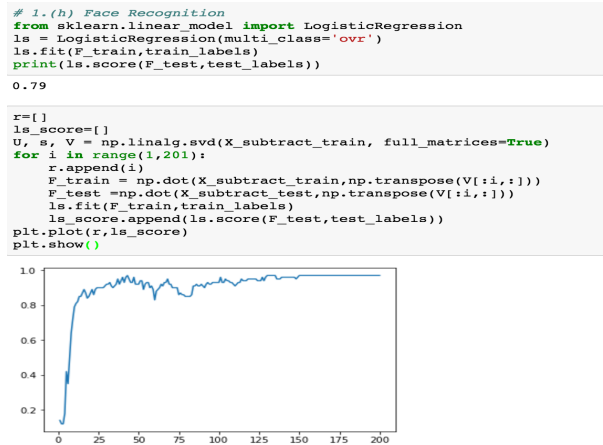


Figure 7:

2 Cooking

1. Load the data

```
# 1.(a) load the data
import json
import pandas as pd
with open('/Users/yingzhu/Desktop/train.json') as json_train:
    train_data = json.load(json_train)
with open('/Users/yingzhu/Desktop/test.json') as json_test:
    test_data = json.load(json_test)
```

```
import numpy as np
train = pd.DataFrame(train_data)
test = pd.DataFrame(test_data)
print(train.shape)
print(test.shape)
(39774, 3)
(9944, 2)
```

Figure 8:

2. Description of the data.

There are 39774 samples in the training set. There are total 20 categories. As what we displayed in the (c) part, there are 6714 unique ingredients. And list of unique ingredients called ingredientslist is created and listed in the (c) part.

```
# 1.(b) How many categories (types of cuisine)?
type_cuisine = train['cuisine'].value_counts()
print(type_cuisine)
print(len(type_cuisine))

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
Name: cuisine, dtype: int64
20
```

Figure 9:

3. Represent each dish by a binary ingredients feature vector. Use $n \times d$ feature matrix to represent all the dishes in training set and test set, where n is the number of dishes and d is the total number of unique ingredients.¹

The following is the display of the training data binary ingredients feature vectors. d equals 6714 in this problem. We also fit the test data into the $n \times d$ binary ingredients feature vector using for loop.

```
# 1.(c) Represent each dish by a binary ingredient feature vector
ingredient_matrix = train['ingredients'].str.join(sep='*').str.get_dummies(sep='*')

# citation for ingredient_matrix:
# https://stackoverflow.com/questions/18889588/create-dummies-from-column-with-multiple-values-in-pandas

ingredient_list = list(ingredient_matrix.columns.values)
print(ingredient_matrix.shape)
ingredient_matrix.head(n=10)

(39774, 6714)
```

	(oz.) tomato sauce	(oz.) tomato paste	(10 oz.) frozen chopped spinach	(10 oz.) frozen chopped spinach thawed and squeezed dry	(14 oz.) sweetened condensed milk	(14.5 oz.) diced tomatoes	(15 oz.) frozen refried beans	1% low- fat buttermilk	1% low- fat chocolate milk	1% low-fat cottage cheese	...	yukon gold potatoes	yuzu juice	yuzu za'atar	zest	zesty italian dressing	zinf
0	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0

10 rows x 6714 columns

Figure 10:

¹<https://stackoverflow.com/questions/18889588/create-dummies-from-column-with-multiple-values-in-pandas>

```

ingredient_matrix_test = np.zeros((len(test),len(ingredient_list)),dtype=np.int)

for i in range(0,len(test)):
    for j in range(0,len(test.values[i][1])):
        if test.values[i][1][j] in ingredient_list:
            pos = ingredient_list.index(test.values[i][1][j])
            ingredient_matrix_test[i][pos]=1

print(ingredient_matrix_test.shape)
ingredient_matrix_test[:10]

(9944, 6714)

array([[0, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0],
       ...,
       [0, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0]])

```

Figure 11:

4. Using Naive Bayes Classification to perform 3 fold cross-validation on the training set and report your average classification accuracy. Try both Gaussian distribution prior assumption and Bernoulli distribution prior assumption.

First, we tried the Gaussian distribution prior assumption. The accuracy by 3 fold cross validation is around 0.38.

```

# 1.(d) Using Naive Bayes Classifier to perform 3 fold cross-validation on the training set
# report average classification accuracy
from sklearn import cross_validation
from sklearn.naive_bayes import GaussianNB, BernoulliNB
clf_g = GaussianNB()
clf_b = BernoulliNB()
#(a) Gaussian distribution\

from sklearn.cross_validation import KFold
def evaluate(X,Y):
    kf = KFold(len(Y), n_folds=3)
    accu=0.0
    sub_accu_list=[]
    for train_indices, test_indices in kf:
        X_train = X[train_indices,:]
        Y_train = Y[train_indices]
        X_test = X[test_indices,:]
        Y_test = Y[test_indices]
        subaccu = clf_g.fit(X_train,Y_train).score(X_test,Y_test)
        sub_accu_list.append(subaccu)
    accu+=subaccu
    print(sub_accu_list)
    return accu/3.0
evaluate(ingredient_matrix.values,train['cuisine'])

[0.37901644290239855, 0.38293860310755767, 0.37758334590435966]
0.37984613063810529

```

Figure 12:

Secondly, we tried the Bernoulli distribution prior assumption. The accuracy by 3 fold cross validation is around 0.68.

```

#(b) Bernoulli distribution
def evaluateB(X,Y):
    kf =KFold(len(Y), n_folds=3)
    accu=0.0
    sub_accu_list=[]
    for train_indices, test_indices in kf:
        X_train = X[train_indices,:]
        Y_train = Y[train_indices]
        X_test = X[test_indices,:]
        Y_test = Y[test_indices]
        subaccu = clf_B.fit(X_train,Y_train).score(X_test,Y_test)
        sub_accu_list.append(subaccu)
        accu+=subaccu
    print(sub_accu_list)
    return accu/3.0
evaluateB(ingredient_matrix.values,train['cuisine'])

[0.68419067732689698, 0.67951425554382261, 0.68690601900739179]
0.68353698395937046

```

Figure 13:

5. For Gaussian prior and Bernoulli prior, which performs better in terms of cross-validation accuracy? Why? Please give specific accuracy.

Compared Gaussian prior with Bernoulli prior, we found that Bernoulli prior performs better with accuracy of around 0.68. The reason of this outcome is that we represent each dish by a binary ingredients feature vector. There are only two possible value for each ingredient, either 0 or 1. Apparently assuming the feature follows a Bernoulli distribution is more reasonable than assuming they follows Gaussian distribution. So compared with Gaussian distribution, Bernoulli distribution fits better.

6. Using logistic regression model to perform 3 fold cross validation on the training set and report you average classification accuracy.

After we run the logistic regression, we get the average classification accuracy as around 0.78. The accuracy of logistic regression is higher than Naive Bayes' result. This is reasonable. Because in Naive Bayes, we assume that each feature is independent to each other. But for each dish, the correlation of each ingredient should be very high. This high correlation is omitted by the Naive Bayes.


```

# 1. (f) Using logistic regression perform 3 fold cross validation on training
# average classification accuracy
from sklearn.linear_model import LogisticRegression
ls = LogisticRegression()
def evaluateLogistic(X,Y):
    kf =KFold(len(Y), n_folds=3)
    accu=0.0
    sub_accu_list=[]
    for train_indices, test_indices in kf:
        X_train = X[train_indices,:]
        Y_train = Y[train_indices]
        X_test = X[test_indices,:]
        Y_test = Y[test_indices]
        subaccu = ls.fit(X_train,Y_train).score(X_test,Y_test)
        sub_accu_list.append(subaccu)
        accu+=subaccu
    print(sub_accu_list)
    return accu/3.0
evaluateLogistic(ingredient_matrix.values,train['cuisine'])

[0.77583345904359635, 0.77213757731181176, 0.77869965303967414]

0.7755689646502746

```

Figure 14:

7. Train your best-performed classification with all the training data and generate test labels on test set. Submit your results to Kaggle and report the accuracy.

After we train Logistic regression with all the training data and generate test labels on test set. We got accuracy of around 0.78 in Kaggle.

Your most recent submission				
Name	Submitted	Wait time	Execution time	Score
dish_test_res.csv	a few seconds ago	0 seconds	0 seconds	0.78338
Complete				
Jump to your position on the leaderboard				

Figure 15: