

CS498 Applied Machine Learning HW1

Name: Yu Che Wang / Netid: yuchecw2

Three accuracies for 1A, 1B, 1D

1A

```
accuracy_vec = [0.6493506 0.7337662 0.6753247 0.6688312 0.6883117 0.6883117 0.7272727  
0.6948052 0.7142857 0.6428571]
```

⇒accuracy = 0.6883117

1B

```
accuracy_vec = [0.6428571 0.6883117 0.6688312 0.6753247 0.6818182 0.6493506 0.6103896  
0.7272727 0.6493506 0.6038961]
```

⇒accuracy = 0.6597403

1D

```
accuracy_vec = [0.6078431 0.7189542 0.6862745 0.6732026 0.6928105 0.6601307 0.6928105  
0.6797386 0.6862745 0.6078431]
```

⇒accuracy = 0.6705882

Code

Evaluations

```
# Part 1B
# Read CSV into R
data = read.csv(file="Desktop/pima-indians-diabetes.csv")
accuracy_vec = seq(1, 10, by=1)
train_ratio = 0.8
attributes = c(3, 4, 6, 8) #blood pressure, skin thickness, BMI, age
param = matrix(), nrow=8, ncol=4)

for (iter in seq(1, 10, by=1)) {
  # Test-train splits
  train_test_list = test_train_split(data, train_ratio, iter)
  train_data = train_test_list[[1]]
  test_data = train_test_list[[2]]
  # Calculate prior probabilities
  outcome = train_data[,9]
  n_positive = length(outcome[outcome==1])
  n_negative = length(outcome[outcome==0])
  P_positive = n_positive/length(outcome)
  P_negative = n_negative/length(outcome)
  # Calculate likelihood
  # Calculate normal distribution parameter for each feature
  for (attribute in attributes) {
    feature = train_data[,attribute]
    # Ignore zero entries aka missing values
    # I don't change 0 entries to NA, ignore those entries in this step
    remain_outcome = outcome[feature!=0]
    feature = feature[feature!=0]
    param[attribute, 1] = mean(feature[remain_outcome==1]) # positive_feature_mean
    param[attribute, 2] = var(feature[remain_outcome==1]) # positive_feature_var
    param[attribute, 3] = mean(feature[remain_outcome==0]) # negative_feature_mean
    param[attribute, 4] = var(feature[remain_outcome==0]) # negative_feature_var
  }
  confusion_matrix = estimate(test_data, param, attribute, P_positive, P_negative)
  n_correct = confusion_matrix[1,1] + confusion_matrix[2,2]
  accuracy = n_correct / nrow(test_data)
  accuracy_vec[iter] = accuracy
}
print(accuracy_vec)
print(mean(accuracy_vec))
```

Probability calculations

```
# Evaluate log probability
evaluate_log_p = function(feature, feature_mean, feature_var) {
  log_p = - ((feature - feature_mean)^2)/(2*feature_var)
  log_p = log_p + log(1/sqrt(2*pi*feature_var))
  return(log_p)
}
```

Calculating the confusion matrix and test-train split

```
# Evaluate on test data and return the confusion matrix
estimate = function(test_data, param, attribute, P_positive, P_negative) {
  true_positive = 0
  false_positive = 0
  true_negative = 0
  false_negative = 0
  logP_positive = log(P_positive)
  logP_negative = log(P_negative)
  for(index in seq(1, nrow(test_data), by=1)) {
    positive_log_p = 0
    negative_log_p = 0
    # Calculate log likelihood
    for (attribute in attributes) {
      if (test_data[index, attribute] == 0) {
        # Ignore this attribute
        next
      } else {
        positive_log_p = positive_log_p + evaluate_log_p(test_data[index,attribute], param[attribute, 1], param[attribute, 2])
        negative_log_p = negative_log_p + evaluate_log_p(test_data[index,attribute], param[attribute, 3], param[attribute, 4])
      }
    }
    if (logP_positive + positive_log_p > logP_negative + negative_log_p) {
      if (test_data[index, 9] == 1) {
        true_positive = true_positive + 1
      } else {
        false_positive = false_positive + 1
      }
    } else {
      if (test_data[index, 9] == 0) {
        true_negative = true_negative + 1
      } else {
        false_negative = false_negative + 1
      }
    }
  }
  confusion_matrix = matrix(
    c(true_positive, false_positive, false_negative, true_negative),
    nrow = 2,
    ncol = 2,
    byrow = TRUE
  )
  return(confusion_matrix)
}

# Split the data
test_train_split = function(data, train_ratio, seed) {
  set.seed(seed)
  sample = sample(seq_len(nrow(data)), size=floor(train_ratio*nrow(data)))
  train_data = data[sample,]
  test_data = data[-sample,]
  return(list(train_data, test_data))
}
```

Table of accuracies for all 12 cases

Accuracies are obtained from validation

	Model	Accuracy
1	Gaussian + untouched	0.558107
2	Gaussian + stretched	0.811631
3	Bernoulli + untouched	0.833274
4	Bernoulli + stretched	0.810357
5	10 trees + 4 depth + untouched	0.753107
6	10 trees + 4 depth + stretched	0.751345
7	10 trees + 16 depth + untouched	0.937464
8	10 trees + 16 depth + stretched	0.942143
9	30 trees + 4 depth + untouched	0.797893
10	30 trees + 4 depth + stretched	0.784762
11	30 trees + 16 depth + untouched	0.955452
12	30 trees + 16 depth + stretched	0.958488

yuchecw2_12.csv a minute ago by Jeffrey Wang 30 trees + 16 depth + stretched	0.96071	<input type="checkbox"/>
yuchecw2_11.csv 2 minutes ago by Jeffrey Wang 30 trees + 16 depth + untouched	0.95600	<input type="checkbox"/>
yuchecw2_10.csv a day ago by Jeffrey Wang 30 trees + 4 depth + stretched	0.76042	<input type="checkbox"/>
yuchecw2_9.csv a day ago by Jeffrey Wang 30 trees + 4 depth + untouched	0.80414	<input type="checkbox"/>
yuchecw2_8.csv a day ago by Jeffrey Wang 10 trees + 16 depth + stretched	0.94414	<input type="checkbox"/>
yuchecw2_7.csv a day ago by Jeffrey Wang 10 trees + 16 depth + untouched	0.93928	<input type="checkbox"/>
yuchecw2_5.csv a day ago by Jeffrey Wang 10 trees + 4 depth + untouched	0.76971	<input type="checkbox"/>
yuchecw2_6.csv 3 days ago by Jeffrey Wang 10 trees + 4 depth + stretched	0.74514	<input type="checkbox"/>
yuchecw2_4.csv 3 days ago by Jeffrey Wang Bernoulli + stretched	0.75942	<input type="checkbox"/>
yuchecw2_3.csv 3 days ago by Jeffrey Wang Bernoulli + untouched	0.81700	<input type="checkbox"/>
yuchecw2_2.csv 3 days ago by Jeffrey Wang Gaussian + stretched	0.80800	<input type="checkbox"/>
yuchecw2_1.csv 3 days ago by Jeffrey Wang Gaussian + untouched	0.54971	<input type="checkbox"/>

Random forest with 30 trees and 16 depth has the best performance since it detects more features.

40 mean images

1: Gaussian + untouched accuracy: 0.558107



2: Gaussian + stretched accuracy: 0.811631



3: Bernoulli + untouched accuracy: 0.833274



4: Bernoulli + stretched accuracy: 0.810357



Code

Main function

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.naive_bayes import BernoulliNB
4 from sklearn.naive_bayes import GaussianNB
5 from sklearn.ensemble import RandomForestClassifier
6 import matplotlib.pyplot as plt
7 import random
8 import math
9 import argparse
10 from helper import *
11
12 parser = argparse.ArgumentParser()
13 parser.add_argument('--model_name', default='Gaussian', type=str)
14 parser.add_argument('--preprocess', default=False, type=bool)
15 parser.add_argument('--num_epochs', default=10, type=int)
16 parser.add_argument('--test', default=True, type=bool)
17 parser.add_argument('--path', default='yuhcw2.1.csv', type=str)
18 parser.add_argument('--n_estimators', default=30, type=int)
19 parser.add_argument('--max_depth', default=16, type=int)
20 args = parser.parse_args()
21
22 train_data = pd.read_csv("MNIST/train.csv")
23 test_data = pd.read_csv("MNIST/test.csv")
24 train_data = pd.DataFrame.as_matrix(train_data)
25 test_data = pd.DataFrame.as_matrix(test_data)
26 if args.preprocess:
27     train_data = preprocess(train_data, False)
28 if not args.test:
29     ratio = 0.2
30     num_data = train_data.shape[0]
31     num_val = (int)(num_data*ratio)
32     accuracy_list = []
33
34     for epoch in range(args.num_epochs):
35         random.seed(epoch)
36         val_samples = random.sample(range(num_data), num_val)
37         mask_array = np.ones((num_data, 1), dtype=bool)
38         for i in val_samples:
39             mask_array[i] = False
40
41         label = np.empty((num_data-num_val, 1))
42         feature = np.empty((num_data-num_val, train_data.shape[1]-1))
43         idx = 0
44         for i in range(num_data):
45             if (mask_array[i]):
46                 label[idx] = train_data[i,0]
47                 feature[idx, :] = train_data[i,1:]
48                 idx = idx + 1
49
50         val_label = train_data[val_samples, 0]
51         val_feature = train_data[val_samples, 1:]
52
53         # Create a model and train it
54         if args.model_name == 'Gaussian':
55             model = GaussianNB()
56         elif args.model_name == 'Bernoulli':
57             # Thresholding
58             feature[feature<127] = 255
59             feature[feature<127] = 0
60             val_feature[val_feature<127] = 255
61             val_feature[val_feature<127] = 0
62             model = BernoulliNB()
63         elif args.model_name == 'Random_Forest':
64             model = RandomForestClassifier(n_estimators=args.n_estimators, max_depth=args.max_depth)
65
66         val_feature[val_feature<127] = 255
67         val_feature[val_feature<127] = 0
68         model = BernoulliNB()
69         label_prediction = fit_model.predict(val_feature)
70         num_correct = (val_label == label_prediction).sum()
71         accuracy = num_correct/num_val
72         accuracy_list.append(accuracy)
73         print('Training and cross validation: uses %s model, preprocess: %d, accuracy: %f'
74               % (args.model_name, args.preprocess, sum(accuracy_list)/len(accuracy_list)))
75     else:
76         label = train_data[:,0]
77         feature = train_data[:,1:]
78         # Create a model and train it
79         if args.model_name == 'Gaussian':
80             model = GaussianNB()
81         elif args.model_name == 'Bernoulli':
82             feature[feature<127] = 255
83             feature[feature<127] = 0
84             model = BernoulliNB()
85         elif args.model_name == 'Random_Forest':
86             model = RandomForestClassifier(n_estimators=args.n_estimators, max_depth=args.max_depth)
87
88         label = np.ravel(label)
89         fit_model = model.fit(feature, label)
90
91         if args.preprocess:
92             test_data = preprocess(test_data, True)
93             test_prediction = fit_model.predict(test_data)
94             df = pd.DataFrame(test_prediction)
95             df_test_data = pd.DataFrame(test_data)
96             df_to_csv(args.path)
97
98         # Calculate class means
99         if args.preprocess:
100             original_size = 20 # box_size
101         else:
102             original_size = 28
103
104         num_classes = 10
105         mean = [None]*num_classes
106         for label in range(num_classes):
107             mean[label] = df_test_data[df[0]==label].mean()
108
109         mean_matrix = []
110         for label in range(num_classes):
111             a = np.empty((original_size, original_size))
112             for j in range(original_size):
113                 for i in range(original_size):
114                     if (mean[label][original_size*j+i] < 127):
115                         a[j, i] = 255
116                     else:
117                         a[j, i] = 0
118             mean_matrix.append(a)
119         plt.figure()
120         for label in range(10):
121             plt.subplot(1, 10, label+1)
122             plt.imshow(mean_matrix[label], cmap='Greys')
123             plt.axis('off')
124         plt.show()
```

Helper functions

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.naive_bayes import BernoulliNB
4 from sklearn.naive_bayes import GaussianNB
5 import random
6 import math
7
8 def preprocess(data, test):
9     box_size = 20
10     original_size = 28
11     if not test:
12         stretched_data = np.empty((data.shape[0], box_size*box_size+1))
13     else:
14         # No label column
15         stretched_data = np.empty((data.shape[0], box_size*box_size))
16     for k in range(data.shape[0]):
17         # Label
18         if not test:
19             stretched_data[k, 0] = data[k, 0]
20         # Find the margin
21         bottom = find_bottom(data, k, original_size)
22         top = find_top(data, k, original_size)
23         num_rows = bottom-top+1
24         right = find_right(data, k, original_size, test)
25         left = find_left(data, k, original_size, test)
26         num_cols = right-left+1
27         # Find the bounding box
28         bounding_box = np.empty((num_rows, num_cols))
29         for j in range(top, bottom+1):
30             for i in range(left, right+1):
31                 if not test:
32                     bounding_box[j-top, i-left] = data[k, original_size*j+i+1]
33                 else:
34                     bounding_box[j-top, i-left] = data[k, original_size*j+i]
35         # plt.subplot(1,2,1)
36         # plt.matshow(bounding_box)
37         # plt.title("#%d bounding box"%(k))
38         # Resize the bounding box
39         stretched_bounding_box = resize_matrix(bounding_box, box_size, box_size)
40         for i in range(box_size):
41             for j in range(box_size):
42                 if not test:
43                     stretched_data[k, box_size*i+j+1] = stretched_bounding_box[i, j]
44                 else:
45                     stretched_data[k, box_size*i+j] = stretched_bounding_box[i, j]
46         return stretched_data
47
48 def resize_matrix(mat, n_row_out, n_col_out):
49     n_row_in = mat.shape[0]
50     n_col_in = mat.shape[1]
51     mat_out = np.empty((n_row_out, n_col_out))
52     row_ratio = n_row_in/n_row_out
53     col_ratio = n_col_in/n_col_out
54     for i in range(n_row_out):
55         for j in range(n_col_out):
56             mat_out[i, j] = mat(math.floor(i*row_ratio), math.floor(j*col_ratio))
57     return mat_out
58
59 def find_top(data, num, original_size):
60     for i in range(1, len(data[num, :])):
61         if data[num, i].item() is not 0:
62             return math.floor(i/original_size)
63
64 def find_bottom(data, num, original_size):
65     for i in range(len(data[num, :]) - 1, 0, -1):
66         if data[num, i].item() is not 0:
67             return math.floor(i/original_size)
68
69 def find_left(data, num, original_size, test):
70     for left in range(original_size):
71         for j in range(original_size):
72             if not test:
73                 if data[num, original_size*j+left+1].item() is not 0:
74                     return left
75             else:
76                 if data[num, original_size*j+left].item() is not 0:
77                     return left
78
79 def find_right(data, num, original_size, test):
80     for right in range(original_size-1, -1, -1):
81         for j in range(original_size):
82             if not test:
83                 if data[num, original_size*j+right+1].item() is not 0:
84                     return right
85             else:
86                 if data[num, original_size*j+right].item() is not 0:
87                     return right
```