

Personal Information

Name: Chengzhi Fan
Student ID: 67883553
CS ID: p2k2b

Template for Question 1

1 Individual Classifiers

1.1 Result

Report the train and test result for each classifier in the given table. You should use the following hyperparameters,

1. Random Forest: no max cap on depth, and a forest size of 15 trees.
2. KNN: $k = 3$.
3. Continuous Naive Bayes: has no hyperparameters.

Model	Your Train Error (%)	Your Test Error (%)
Random Forest	0.3	13.8
KNN	7.8	13.4
Naive Bayes	19.6	17.6

Explain in one paragraph why you think a particular classifier works better on this dataset.

Answer:

KNN classifier works the best on this dataset.

The Naive Bayes runs faster on this dataset, but it's both train and test errors are considerably higher than the other two. For Random Forest, it has relatively good test error, but it runs significantly slow. Only KNN has the lowest test error and performs fast as well.

1.2 Code

Include the code you have written for each particular classifier.

1. Random Forest

Answer:

```

86
87 def Gini_impurity(p):
88     return np.sum(p*(1-p))
89
90 class DecisionStumpGiniIndex(DecisionStumpErrorRate):
91
92     def fit(self, X, y, split_features=None, thresholds=None):
93         N, D = X.shape
94         # Address the trivial case where we do not split.
95         count = np.bincount(y)
96         p = count / np.sum(count) # Convert counts to probabilities.
97         minGiniIndex = Gini_impurity(p)
98
99         self.splitVariable = None
100        self.splitValue = None
101        self.splitSat = np.argmax(count)
102        self.splitNot = None
103
104        # Check if labels are not all equal.
105        if np.unique(y).size <= 1:
106            return
107
108        if split_features is None:
109            split_features = range(D)
110        for d in split_features:
111            column = np.unique(X[:, d]) if (thresholds is None) else np.unique(thresholds[:, d])
112            for value in column:
113                # Count number of class labels where the feature is greater than threshold.
114                y_vals = y[X[:, d] > value]
115                # Avoid "RuntimeWarning: invalid value encountered".
116                if len(y_vals) != 0 and len(y_vals) != N:
117                    count1 = np.bincount(y_vals, minlength=len(count))
118                    count0 = count - count1
119
120                    # Compute Gini Index.
121                    p1 = count1 / np.sum(count1)
122                    p0 = count0 / np.sum(count0)
123                    G1 = Gini_impurity(p1)
124                    G0 = Gini_impurity(p0)
125                    prob1 = np.sum(X[:, d] > value) / N
126                    prob0 = 1 - prob1
127
128                    giniIndex = prob0 * G0 + prob1 * G1
129
130                    # Compare to minimum error so far.
131                    if giniIndex < minGiniIndex:
132                        # This is the lowest Gini Index, store this value.
133                        minGiniIndex = giniIndex
134                        self.splitVariable = d
135                        self.splitValue = value
136                        self.splitSat = np.argmax(count1)
137                        self.splitNot = np.argmax(count0)
138

```

```

286 ...
287 def create_splits(self, X):
288     K = 15
289     N, D = X.shape
290     thresholds = np.zeros((K, D))
291
292     for j in range(D):
293         model = Kmeans(k=K)
294         column = X[:, j:j+1]
295         # I modified Kmeans.fit() to return its means.
296         means = model.fit(column)
297         thresholds[:, j:j+1] = means
298
299     self.thresholds = thresholds
300

```

2. KNN

Answer:

```

38
39 def cosine_distance(self, X1, X2):
40     "Insert your code here"
41
42     # We will first compute cosine similarity, then transform it to cosine distance.
43     # cosine_similarity(a, b) = (a . b) / (|a| * |b|)
44
45     # numerator:    a . b
46     numeratorMatrix = X1 @ (X2.T)
47
48     # denominator: |a| * |b|
49     norms1 = np.linalg.norm(X1, axis=1).reshape((X1.shape[0], 1))
50     norms2 = np.linalg.norm(X2, axis=1).reshape((X2.shape[0], 1))
51     denominatorMatrix = norms1 @ (norms2.T)
52
53     # if there is any zero row in X1 or X2, set the distance between any zero row in X1 to all the rows X2
54     # to zero or vice versa. So we can play a trick here by setting denominator to inf in the case above.
55     denominatorMatrix = np.where(denominatorMatrix == 0, np.inf, denominatorMatrix)
56
57     # transform cosine similarity to cosine distance
58     cosineSimilarity = (numeratorMatrix / denominatorMatrix)
59     cosineDistance = 1 - cosineSimilarity
60     return cosineDistance
61

```

3. Naive Bayes

Answer:

```

3
4 class NaiveBayes:
5     def __init__(self):
6         pass
7
8     def fit(self, X, y):
9         # Important: First clean dataset by removing its duplicate data. Otherwise, the duplicate data will affect
10        # the expected shape of normal distribution (mean stdDev), thus increases training and testing errors.
11        Xy = np.hstack((X, y.reshape((len(y), 1))))
12        Xy_unique = np.unique(Xy, axis=0)
13        N, D = Xy_unique.shape
14        X = Xy_unique[:, 0:D-1]
15        y = Xy_unique[:, D-1].astype(y.dtype)
16
17        # Compute the probability of each class i.e p(y==c).
18        N, D = X.shape
19        counts = np.bincount(y)
20        p_y = counts / N
21
22        # Compute the number of class labels.
23        C = len(counts)
24
25        # Let x_dc represent data entry from X with feature d and class c.
26        # Compute the mean and stdDev for each x_dc.
27        meanMatrix = np.zeros((D, C), float)
28        stdDevMatrix = np.zeros((D, C), float)
29        for d in range(D):
30            for c in range(C):
31                x_dc = X[:, d][(y == c)]
32                meanMatrix[d, c] = np.mean(x_dc)
33                stdDevMatrix[d, c] = np.std(x_dc)
34
35        self.c = C
36        self.p_y = p_y
37        self.meanMatrix = meanMatrix
38        self.stdDevMatrix = stdDevMatrix
39
40    def predict(self, X):
41        C = self.c
42        p_y = self.p_y
43        mean = self.meanMatrix
44        stdDev = self.stdDevMatrix
45
46        N, D = X.shape
47        y_pred = np.zeros(N, int)
48        for i in range(N):
49            # make predictions p(y) in log space.
50            log_p_y = np.log(p_y)
51            # make predictions p(x|y) in log space.
52            x = np.array([X[i],])*C).T
53            log_p_xy = np.sum((-0.5 * (((x - mean) / stdDev) ** 2) - np.log(stdDev * np.sqrt(2 * np.pi))), axis=0)
54            # p(x|y) * p(y) in log space
55            result = log_p_xy + log_p_y
56
57            y_pred[i] = np.argmax(result)
58
59    return y_pred
60

```

2 Stacking

2.1 Result

Report the test error and training error of the stacking classifier.

Answer:

Stacking

Training error: 0.078

Testing error: 0.119

Runtime: 0:01:22.775136

2.2 Code

Include all the code you have written for stacking classifier

Answer:

```
7
8 class Stacking():
9
10     def __init__(self):
11         # classifiers
12         self.randomForest = RandomForest(max_depth=np.inf, num_trees=15)
13         self.knn = KNN(k=3)
14         self.naiveBayes = NaiveBayes()
15
16         # meta-classifier
17         self.decisionTree = DecisionTree(max_depth=np.inf) # stump_class=DecisionStumpErrorRate by default.
18
19     def fit(self, X, y):
20         self.randomForest.fit(X, y)
21         self.knn.fit(X, y)
22         self.naiveBayes.fit(X, y)
23
24     def predict(self, X):
25         pred0 = self.randomForest.predict(X)
26         pred1 = self.knn.predict(X)
27         pred2 = self.naiveBayes.predict(X)
28
29         # Use predictions of the classifiers to form X y for meta-classifier.
30         X = np.vstack((pred0, pred1, pred2)).T
31         y = stats.mode(X, axis=1)[0].flatten()
32
33         # Train meta-classifier and return predictions.
34         self.decisionTree.fit(X, y)
35         return self.decisionTree.predict(X)
36
```