

Homework 4

2021年11月6日 星期六 下午7:47

Problem 1

$$1. \quad \epsilon_1 = \frac{5}{16} = 0.3125$$

$$\alpha_1 = \frac{1}{2} \log \left(\frac{1 - \epsilon_1}{\epsilon_1} \right) = 0.39$$

$$2. \quad \frac{1}{16} e^{-0.39 \cdot (+1)(+1)} = 0.04231605466$$

$$\frac{1}{16} e^{-0.39 \cdot (-1)(+1)} = 0.09231129962$$

$$w_1' = \frac{1}{11 \times 0.04231605466 + 5 \times 0.09231129962} \cdot 0.04231605466$$

$$= 0.04564675705 \quad (\text{the instances predicted correct})$$

$$w_2' = \frac{1}{11 \times 0.04231605466 + 5 \times 0.09231129962} \cdot 0.09231129962$$

$$= 0.0995771345 \quad (\text{the instances predicted incorrect})$$

3. Yes. In this case, $\alpha_i = +\infty$, and the weights of all examples are 0.

Problem 2

1. **Boosted Decision Trees:** Boosted decision trees use an efficient implementation of the MART gradient boosting algorithm. Gradient boosting is a machine learning technique for regression problems. It builds each regression tree in a stepwise fashion, using a predefined loss function to measure the error in each step and correct for it in the next. Gradient boosting works by building simpler (weak) prediction models sequentially where each model tries to predict the error left over by the previous model. Because of this, the algorithm tends to overfit rather quick.

Training methodology

Load the prepared data -> `load_svmlight_file('a9a.txt')`

Train the XGBoost Model -> split the training data and use the K-fold CV

Tune the hyperparameters -> use grid search

Decide the final classifiers

Make Predictions with XGBoost Model -> `model.score(X_test, y_test)`

List of hyperparameter: eta, gamma, max_depth, min_child_weight, max_delta_step, subsample, sampling_method, colsample_bytree, colsample_bylevel, colsample_bynode, lambda, alpha, tree_method, sketch_eps, scale_pos_weight, updater, refresh_leaf, process_type, grow_policy, max_leaves, max_bin, predictor, num_parallel_tree, monotone_constraints, interaction_constraints

Hyperparameter I tuned:

tree method [default = auto]: For training boosted tree models

max_depth [default = 6]: Maximum depth of a tree.

min_child_weight [default = 1]: Minimum sum of instance weight (hessian) needed in a child.

eta [default = .3]: Step size shrinkage used in update to prevents overfitting.

colsample_bytree [default = 1]: is the subsample ratio of columns when constructing each tree.

final hyperparameter settings:

tree method [default = auto]

max_depth = 7

min_child_weight [default = 1]:

eta [default = .3]:

colsample_bytree = 5:

Training error rates: 0.1542

Cross-validation error rates: 0.1535626536

test error rates: 0.1518

2. Random Forest: As in bagging, we build several decision trees on bootstrapped training samples. But when building these decision trees, each time a split in a tree is considered, a random sample of $p < d$ features (predictors) is chosen as split candidates from the full set of all d features. The split is allowed to use only one of those p features. A fresh sample of features is taken at each split, and typically we choose $p < \sqrt{d}$

Training methodology

Load the prepared data -> `load_svmlight_file('a9a.txt')`

Train the RandomForest Model -> split the training data and use the K-fold CV
`clf.fit(X_train, y_train)`

Tune the hyperparameters -> use grid search

Decide the final classifiers

Make Predictions with RandomForest Model -> `clf.score(X_test, y_test)`

List of hyperparameter: {'bootstrap': True, 'ccp_alpha': 0.0, 'class_weight': None, 'criterion': 'gini', 'max_depth': None, 'max_features': 'auto', 'max_leaf_nodes': None, 'max_samples': None, 'min_impurity_decrease': 0.0, 'min_impurity_split': None, 'min_samples_leaf': 1, 'min_samples_split': 2, 'min_weight_fraction_leaf': 0.0, 'n_estimators': 100, 'n_jobs': None, 'oob_score': False, 'random_state': None, 'verbose': 0, 'warm_start': False}

Hyperparameter I tuned:

`n_estimators` [default = 100]: number of trees in the forest

`max_features` [default = auto] = max number of features considered for splitting a node

`max_depth` [default = None] = max number of levels in each decision tree

`min_samples_split` [default = 2] = min number of data points placed in a node before the node is split

`min_samples_leaf` [default = 1] = min number of data points allowed in a leaf node

`bootstrap` [default = True] = method for sampling data points (with or without replacement)

final hyperparameter settings:

```
{'n_estimators': 200,  
'min_samples_split': 5,  
'min_samples_leaf': 2,  
'max_features': 'auto',  
'bootstrap': False}
```

Training error rates: 0.1533

Cross-validation error rates: 0.1565520066

test error rates: 0.155

- Support Vector Machines with Gaussian Kernel: In SVM, kernels are used for solving nonlinear problems in higher dimensional where linear separation is not possible. Generally, SVM is a simple dot product operation. Therefore, we should choose a kernel such that, $K(x,y)=K(x).K(y)$. Gaussian is one such kernel giving good linear separation in higher dimension for many nonlinear problems. How it works? $\Phi: R \rightarrow R^\infty$ $\Phi(x) =$

$$\exp(-x^2)(1, \sqrt{\frac{2^1}{1!}}x, \sqrt{\frac{2^2}{2!}}x^2, \sqrt{\frac{2^3}{3!}}x^3, \dots)$$

Training methodology

Load the prepared data -> `load_svmlight_file('a9a.txt')`

Train the SVM Gaussian Model -> split the training data and use the K-fold CV
`clf.fit(X_train, y_train)`

Tune the hyperparameters -> use grid search

Decide the final classifiers

Make Predictions with SVM Gaussian Model -> `clf.score(X_test, y_test)`

List of hyperparameter: {'C': 1.0, 'break_ties': False, 'cache_size': 200, 'class_weight': None, 'coef0': 0.0, 'decision_function_shape': 'ovr', 'degree': 3, 'gamma': 'scale', 'kernel': 'rbf', 'max_iter': -1, 'probability': False, 'random_state': None, 'shrinking': True, 'tol': 0.001, 'verbose': False}

Hyperparameter I tuned:

Kernel[rbf]: The function of kernel is to take data as input and transform it into the required form.

final hyperparameter settings:

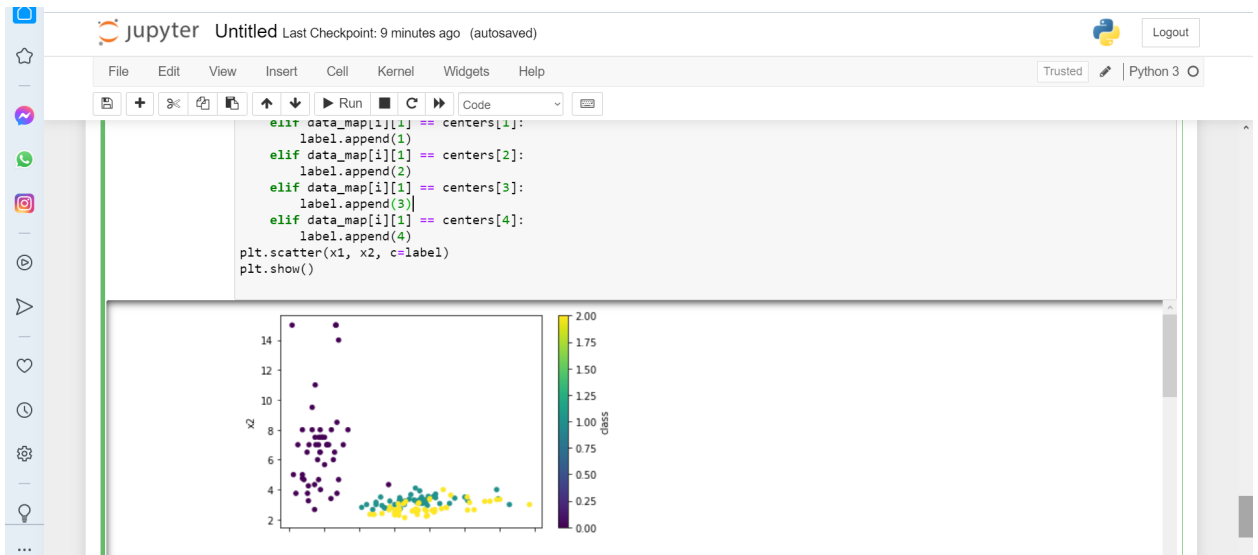
kernel = linear

Training error rates: 0.1569831716

Cross-validation error rates: 0.152047502

test error rates: 0.1502

Problem 3



jupyter Problem3.py 10 minutes ago

```
1 import numpy as np
2 import scipy as sp
3 import pandas as pd
4 import matplotlib
5 import matplotlib.pyplot as plt
6
7 data = np.loadtxt('iris.data', delimiter = ",", dtype='str')
8 #print(data)
9
10 temp = pd.DataFrame(data, columns=list('abcde'))
11 temp['e'] = [0 if item == 'Iris-setosa' else 1 if item == 'Iris-versicolor' else 2 for item in temp['e']]
12 temp = temp.astype(float)
13 x1 = temp['a']/temp['b']
14 x2 = temp['c']/temp['d']
15
16 new_data = pd.DataFrame()
17 new_data['x1'] = x1
18 new_data['x2'] = x2
19 new_data['class'] = temp['e']
20
21 new_data.plot.scatter('x1', 'x2', c='class', colormap='viridis')
22 plt.show()
23
24
25 def k_init(X, k):
26     """ k-means++ initialization algorithm
27
```

jupyter Problem3.py 10 minutes ago

```
37
38 init_centers: array (k, d)
39     The initialize centers for kmeans++
40 """
41 # initialize k centroids
42 centroids = np.zeros((k, X.shape[1]))
43 #print(centroids[0])
44 centroids[0] = X[np.random.randint(X.shape[0])]
45 #print(X[0])
46 #centroids[0].append(X[0])
47 for i in range(1, k):
48     # compute distances from centroids
49     distances = np.zeros(X.shape[0])
50     for j in range(X.shape[0]):
51         distances[j] = np.linalg.norm(X[j] - centroids[i-1])
52     # compute probabilities
53     probabilities = distances / np.sum(distances)
54     # choose centroid
55     centroids[i] = X[np.random.choice(X.shape[0], 1, p=probabilities)]
56     print(X.shape[0])
57 return centroids
58
59
60 def assign_data2clusters(X, C):
61     """ Assignments of data to the clusters
62     Parameters
63     """
```

jupyter Problem3.py 11 minutes ago Logout

File Edit View Language Python

```
73     """ The binary matrix A which shows the assignments of data points (X) to
74     the input centers (C).
75     """
76     data_map = []
77     for i in range(len(X)):
78         cluster = None
79         min_dist = np.inf
80         for c in C:
81             # Calculate distance between data point and cluster center
82             dist = np.linalg.norm(X[i]-c)
83             if dist < min_dist:
84                 min_dist = dist
85                 cluster = c
86         data_map.append(cluster)
87     return data_map
88
89
90 def compute_objective(X, C):
91     """ Compute the clustering objective for X and C
92     Parameters
93     -----
94     X: array, shape(n,d)
95         Input array of n samples and d features
96
97     C: array, shape(k,d)
98         The final cluster centers
99
```

jupyter Problem3.py 11 minutes ago Logout

File Edit View Language Python

```
103     """ The objective for the given assignments
104     """
105     data_map = assign_data2clusters(X,C)
106     total_dist = 0
107     for c in centers:
108         for item in data_map:
109             dist = 0
110             if item[1][0] == c[0] and item[1][1] == c[1]:
111                 for j in range(len(c)):
112                     dist += (c[j] - item[0][j])**2
113                 total_dist += dist
114     return total_dist
115
116
117 def k_means_pp(X, k, max_iter):
118     """ k-means++ clustering algorithm
119
120     step 1: call k_init() to initialize the centers
121     step 2: iteratively refine the assignments
122
123     Parameters
124     -----
125     X: array, shape(n,d)
126         Input array of n samples and d features
127
128     k: int
129         The number of clusters

```

jupyter Problem3.py 11 minutes ago

File Edit View Language Python

```
139 """ The objective value at each iteration
140 """
141 # initialize centroids
142 centroids = k_init(X, k)
143 # initialize cluster assignment
144 cluster_assignment = np.zeros(X.shape[0])
145 # initialize cluster centers
146 cluster_centers = np.zeros((k, X.shape[1]))
147 # initialize cluster counts
148 cluster_counts = np.zeros(k)
149 # initialize cluster centers
150 cluster_centers = centroids
151 # initialize cluster counts
152 cluster_counts = np.zeros(k)
153 # initialize iteration counter
154 iter_count = 0
155 # initialize convergence flag
156 converged = False
157 # run kmeans++ algorithm
158 while not converged:
159     # update iteration counter
160     iter_count += 1
161     # update cluster assignment
162     for i in range(X.shape[0]):
163         # compute distances from centroids
164         distances = np.zeros(k)
165         for j in range(k):
```

jupyter Problem3.py 12 minutes ago

File Edit View Language Python

```
163     # compute distances from centroids
164     distances = np.zeros(k)
165     for j in range(k):
166         distances[j] = np.linalg.norm(x[i] - centroids[j])
167     # compute probabilities
168     probabilities = distances / np.sum(distances)
169     # choose centroid
170     cluster_assignment[i] = np.random.choice(k, 1, p=probabilities)
171 # update cluster centers
172 for i in range(k):
173     # compute cluster counts
174     cluster_counts[i] = np.sum(cluster_assignment == i)
175     # compute cluster centers
176     cluster_centers[i] = np.sum(x[cluster_assignment == i], axis=0) / cluster_counts[i]
177 # check for convergence
178 if np.array_equal(centroids, cluster_centers):
179     converged = True
180 else:
181     centroids = cluster_centers
182 return centroids
183
184
185 cluster_data = pd.DataFrame()
186 cluster_data['x1'] = x1
187 cluster_data['x2'] = x2
188
189 #run with k=1,2,3,4,5
```

jupyter Problem3.py 12 minutes ago

File Edit View Language Python

```
187 cluster_data['x2'] = x2
188
189 #run with k=1,2,3,4,5
190 acc = []
191 centers = k_means_pp(cluster_data, 1, 50)
192 acc.append(compute_objective(cluster_data, centers))
193
194 centers = k_means_pp(cluster_data, 2, 50)
195 acc.append(compute_objective(cluster_data, centers))
196
197 centers = k_means_pp(cluster_data, 3, 50)
198 acc.append(compute_objective(cluster_data, centers))
199
200 centers = k_means_pp(cluster_data, 4, 50)
201 acc.append(compute_objective(cluster_data, centers))
202
203 centers = k_means_pp(cluster_data, 5, 50)
204 acc.append(compute_objective(cluster_data, centers))
205
206 plt.plot([1,2,3,4,5], acc)
207 plt.show()
208 plt.clf()
209
210 #my best was k=5 so now ill change the number fo iterations with k=5
211 acc = []
212 centers = k_means_pp(new_set, 5, 1)
213 acc.append(compute_objective(new_set, centers))
```

jupyter Problem3.py 12 minutes ago

File Edit View Language Python

```
217
218 centers = k_means_pp(new_set, 5, 50)
219 acc.append(compute_objective(new_set, centers))
220
221 centers = k_means_pp(new_set, 5, 100)
222 acc.append(compute_objective(new_set, centers))
223
224 centers = k_means_pp(new_set, 5, 200)
225 acc.append(compute_objective(new_set, centers))
226
227 plt.plot([1,20,50,100,200], acc)
228 plt.show()
229 plt.clf()
230
231 #plot with data colored by cluster
232 centers = k_means_pp(new_set, 5, 200)
233 data_map = assign_data2clusters(new_set, centers)
234 label = []
235 for i in range(len(data_map)):
236     if data_map[i][1] == centers[0]:
237         label.append(0)
238     elif data_map[i][1] == centers[1]:
239         label.append(1)
240     elif data_map[i][1] == centers[2]:
241         label.append(2)
242     elif data_map[i][1] == centers[3]:
243         label.append(3)
```

jupyter Problem3.py 13 minutes ago

File Edit View Language Python

```
222 acc.append(compute_objective(new_set, centers))
223
224 centers = k_means_pp(new_set, 5, 200)
225 acc.append(compute_objective(new_set, centers))
226
227 plt.plot([1,20,50,100,200], acc)
228 plt.show()
229 plt.clf()
230
231 #plot with data colored by cluster
232 centers = k_means_pp(new_set, 5, 200)
233 data_map = assign_data2clusters(new_set, centers)
234 label = []
235 for i in range(len(data_map)):
236     if data_map[i][1] == centers[0]:
237         label.append(0)
238     elif data_map[i][1] == centers[1]:
239         label.append(1)
240     elif data_map[i][1] == centers[2]:
241         label.append(2)
242     elif data_map[i][1] == centers[3]:
243         label.append(3)
244     elif data_map[i][1] == centers[4]:
245         label.append(4)
246 plt.scatter(x1, x2, c=label)
247 plt.show()
248
```


My final classifier is $k = 5$

All the code is in Jupyter notebook and problem3.