

hw10

thefinalhw/	1
└── ch8p10.py	1
└── ch8p8.py	3

ch8p10.py

```
1 from numpy import *
2 import matplotlib.pyplot as plt
3 from sklearn.ensemble import GradientBoostingRegressor
4 from sklearn.metrics import mean_squared_error
5 from sklearn import tree
6 from sklearn.linear_model import LinearRegression, Ridge
7 from sklearn.ensemble import BaggingRegressor
8 #AtBat Hits HmRun Runs RBI Walks Years CatBat Chits ChmRun CRuns CRBI CWalks League Division PutOuts Assists Errors Salary NewLeague
9
10 hitters dt = dtype([('atbat', float64), ('hits', float64), ('hmrn', float64), ('runs', float64), ('rbi', float64), ('walks', float64), ('years', float64), ('catbat', float64)
11 , ('chits', float64), ('chmrun', float64), ('cruns', float64), ('crbi', float64), ('cwalks', float64), ('league', 'U10'), ('division', 'U10'), ('putouts', float64)
12 , ('assits', float64), ('errors', float64), ('salary', object), ('newleague', 'U10')])
13
14
15 def data_loader(fname):
16     data_data_a = loadtxt(fname, skiprows=1, usecols=(14,15,20), delimiter=',', dtype=str)
17     quant_data_a = loadtxt(fname, skiprows=1, usecols=(1,2,3,4,5,6,7,8,9,10,11,12,13,16,17,18), delimiter=',', dtype=float64)
18     mix_data_a = loadtxt(fname, skiprows=1, usecols=(19), delimiter=',', dtype=object)
19     unfil_data_a = empty(quant_data_a.shape[0], dtype=hitters_dt)
20
21     unfil_data_a['atbat'] = quant_data_a[:,0]
22     unfil_data_a['hits'] = quant_data_a[:,1]
23     unfil_data_a['hmrn'] = quant_data_a[:,2]
24     unfil_data_a['runs'] = quant_data_a[:,3]
25     unfil_data_a['rbi'] = quant_data_a[:,4]
26     unfil_data_a['walks'] = quant_data_a[:,5]
27     unfil_data_a['years'] = quant_data_a[:,6]
28     unfil_data_a['catbat'] = quant_data_a[:,7]
29     unfil_data_a['chits'] = quant_data_a[:,8]
30     unfil_data_a['chmrun'] = quant_data_a[:,9]
31     unfil_data_a['cruns'] = quant_data_a[:,10]
32     unfil_data_a['crbi'] = quant_data_a[:,11]
33     unfil_data_a['cwalks'] = quant_data_a[:,12]
34     #unfil_data_a['league'] = data_data_a[:,0]
35     #unfil_data_a['division'] = data_data_a[:,1]
36     unfil_data_a['putouts'] = quant_data_a[:,13]
37     unfil_data_a['assits'] = quant_data_a[:,14]
38     unfil_data_a['errors'] = quant_data_a[:,15]
39     unfil_data_a['salary'] = mix_data_a
40     #unfil_data_a['newleague'] = data_data_a[:,2]
41     print(unfil_data_a.shape[0])
42     return unfil_data_a
43
44 def filter_salary(unfil_data_a):
45     deleted_l = []
46     name_l = ['atbat', 'hits', 'hmrn', 'rbi', 'walks', 'years', 'catbat', 'chits', 'chmrun', 'cruns', 'crbi', 'cwalks', 'passouts', 'assits', 'errors', 'salary']
47
48     for i in range(unfil_data_a.shape[0]):
49         salary_v = unfil_data_a[i]['salary']
50         if 'NA' in salary_v:
51             deleted_l.append(i)
52     print("deleted ppl", deleted_l)
53
54     data_a = delete(unfil_data_a, deleted_l, axis=0)
55
56     for i in range(data_a.shape[0]):
57         salary_v = data_a[i]['salary']
58         data_a[i]['salary'] = log(float(salary_v))
59     #print("new list of ppl with salary log change", filter_data_a)
60
61     #salary_data = data_a['salary']
62
63     return data_a
64
65 def make_x_y(data_a):
66     name_l = ['atbat', 'hits', 'hmrn', 'rbi', 'walks', 'years', 'catbat', 'chits', 'chmrun', 'cruns', 'crbi', 'cwalks', 'putouts', 'assits', 'errors']
67     # 'hmrn', 'rbi', 'walks', 'years', 'catbat', 'chits', 'chmrun', 'cruns', 'crbi', 'cwalks', 'passouts', 'assits', 'errors']
68     x_a = empty((data_a.shape[0], len(name_l)))
69     for i, name in enumerate(name_l):
70         x_a[:,i] = data_a[name]
71     y_v = data_a['salary']
72
73     return x_a, y_v
74
75 def des_tree(xtrain_a, ytrain_v, ytest_v):
76     clf = tree.DecisionTreeClassifier()
77     clf = clf.fit(xtrain_a, ytrain_v)
78     clf.predict(ytrain_v)
79
80
81 def train_grad_boost(xtrain_a, ytrain_v):
82     #def grad_boost(xtrain_a, ytrain_v, xtest_a, ytest_v):
83
84     lamda shrinkage = [0.001, 0.005, 0.01, 0.05, 0.1, 0.15, 0.5, 1]
85     train_error_l = []
86     for i in lamda shrinkage:
87         model = GradientBoostingRegressor(n_estimators=1000, learning_rate=i)
88         model.fit(xtrain_a, ytrain_v)
89         ypred_v = model.predict(xtrain_a)
90         mse = mean_squared_error(ytrain_v, ypred_v)
91         train_error_l.append(mse)
92     print("train mse boost", train_error_l)
93     plt.plot(lamda shrinkage, train_error_l, marker='o')
94     plt.xlabel('Shrinkage Parameter (λ)')
95     plt.ylabel('Training Set MSE')
96     plt.title('Effect of Shrinkage on Training Set MSE')
97     plt.grid(True)
98     plt.show()
99
100 def test_grad_boost(xtrain_a, ytrain_v, xtest_a, ytest_v):
101     #def grad_boost(xtrain_a, ytrain_v, xtest_a, ytest_v):
102
103     lamda shrinkage = [0.001, 0.005, 0.01, 0.05, 0.1, 0.15, 0.5, 1]
104     test_error_l = []
105     for i in lamda shrinkage:
106         model = GradientBoostingRegressor(n_estimators=1000, learning_rate=i)
107         model.fit(xtrain_a, ytrain_v)
108         ytestpred_v = model.predict(xtest_a)
109         mse = mean_squared_error(ytest_v, ytestpred_v)
110         feature_importances = model.feature_importances_
111         test_error_l.append(mse)
112     print("test mse boost", test_error_l)
113     plt.plot(lamda shrinkage, test_error_l, marker='o')
114     plt.xlabel('Shrinkage Parameter (λ)')
115     plt.ylabel('Training Set MSE')
```

```

116 plt.title('Effect of Shrinkage on Test Set MSE')
117 plt.grid(True)
118 plt.show()
119
120 plt.show()
121
122 def find_imp_pred(xtrain_a, ytrain_v, name_l):
123     model = GradientBoostingRegressor(n_estimators=1000, learning_rate=0.227625)
124     #use this as learning rate cause it is the average of what was given earlier
125     model.fit(xtrain_a, ytrain_v)
126     feat_imp = model.feature_importances_
127     feat_imp_d = {f: imp for f, imp in zip(name_l, feat_imp)}
128     sort_feat = sorted(feat_imp_d.items(), key=lambda x: x[1], reverse=True)
129     print(sort_feat)
130     return sort_feat
131
132
133 def lin_reg(xtrain_a, ytrain_v, xtest_a, ytest_v):
134     lin_model = LinearRegression()
135     lin_model.fit(xtrain_a, ytrain_v)
136     y_pred_lin = lin_model.predict(xtest_a)
137     test_mse_linear = mean_squared_error(ytest_v, y_pred_lin)
138     print("test_mse_linear", test_mse_linear)
139
140 def rid_reg(xtrain_a, ytrain_v, xtest_a, ytest_v):
141     #like in grad boost and using different lamda values i will use dif alpha values for ridge
142     alpha= [0.001, 0.005,0.01, 0.05,0.1,0.5, 1, 5, 10]
143     rid_error = []
144     for i in alpha:
145         rid_model = Ridge(alpha=i)
146         rid_model.fit(xtrain_a, ytrain_v)
147         y_pred_rid = rid_model.predict(xtest_a)
148         rid_error.append(mean_squared_error(ytest_v, y_pred_rid))
149
150     print("test mse ridge", rid_error)
151
152 def bag(xtrain_a, ytrain_v, xtest_a, ytest_v):
153     model = BaggingRegressor(n_estimators=1000, random_state=1)
154     model.fit(xtrain_a, ytrain_v)
155     y_pred_test = model.predict(xtest_a)
156     mse_bag = mean_squared_error(ytest_v, y_pred_test)
157     print("test mse bagging", mse_bag)
158     return mse_bag
159
160 def main():
161     unfdata_a = data_loader("Hitters.csv")
162     data_a = filter_salary(unfdata_a)
163     x_a, y_v = make_x_y(data_a)
164     xtrain_a = x_a[:200]
165     ytrain_v = y_v[:200]
166     xtest_a = x_a[200:]
167     ytest_v = y_v[200:]
168     train_grad_boost(xtrain_a,ytrain_v)
169     test_grad_boost(xtrain_a, ytrain_v, xtest_a,ytest_v)
170     lin_reg(xtrain_a, ytrain_v, xtest_a,ytest_v)
171     rid_reg(xtrain_a, ytrain_v, xtest_a,ytest_v)
172     name_l = ['atbat', 'hits', 'hmrn', 'rbi', 'walks', 'years', 'catbat', 'chits', 'chmrun', 'cruns', 'crbi', 'cwalks', 'putouts', 'assits', 'errors']
173     find_imp_pred(xtrain_a, ytrain_v, name_l)
174     bag(xtrain_a, ytrain_v, xtest_a,ytest_v)
175     raise SystemExit
176     xtrain_a = filtered_log_a['']
177     ytrain_v = filtered_log_a[200]['salary']
178     #xtest_a = filtered_log_a[200:][xtrain_dtype]
179     #ytest_v = filtered_log_a[200:]['salary']
180 if __name__ == '__main__':
181     main()

```

```

1 from numpy import *
2 import matplotlib.pyplot as plt
3 from sklearn.ensemble import GradientBoostingRegressor
4 from sklearn.metrics import mean_squared_error
5 from sklearn import tree
6 from sklearn.linear_model import LinearRegression, Ridge
7 from sklearn.ensemble import BaggingRegressor
8
9 from sklearn.tree import DecisionTreeRegressor, plot_tree
10 from sklearn.model_selection import cross_val_score
11 from sklearn.metrics import mean_squared_error
12 from sklearn.model_selection import cross_val_predict
13 from sklearn.ensemble import BaggingRegressor
14 from sklearn.ensemble import RandomForestRegressor
15
16 def data_loader(fname):
17     x_a = loadtxt(fname, skiprows=1, usecols=(1,2,3,4,5,7,8), delimiter=',', dtype=float64)
18     y_v = loadtxt(fname, skiprows=1, usecols=0, delimiter=',', dtype=float64)
19
20     return x_a, y_v
21
22
23 def tree_things(train_x_a, train_y_v):
24     tree_reg = DecisionTreeRegressor(max_depth=None)
25     tree_reg.fit(train_x_a, train_y_v)
26     ytrain_pred = tree_reg.predict(train_x_a)
27     train_mse = mean_squared_error(train_y_v, ytrain_pred)
28     plt.figure(figsize=(20, 15))
29     plot_tree(tree_reg, filled=True)
30     plt.show()
31
32     return tree_reg, train_mse
33
34 def crossval_for_tree(xtrain, ytrain, xtest, ytest, maxrange):
35     train_mse_l = []
36     test_mse_l = []
37     cross_val_l = []
38     for i in maxrange:
39         tree_reg = DecisionTreeRegressor(max_depth=i)
40         ypred = cross_val_predict(tree_reg, xtrain, ytrain, cv=50)
41         cross_val_mse = mean_squared_error(ytrain, ypred)
42         cross_val_l.append(cross_val_mse)
43
44         tree_reg.fit(xtrain, ytrain)
45         ytrainpred = tree_reg.predict(xtrain)
46         trainmse = mean_squared_error(ytrain, ytrainpred)
47         train_mse_l.append(trainmse)
48
49         ytestpred = tree_reg.predict(xtest)
50         testmse = mean_squared_error(ytest, ytestpred)
51         test_mse_l.append(testmse)
52     print("training mse", train_mse_l)
53     print("test mse", test_mse_l)
54     print("cross val mse", cross_val_l)
55     plt.figure(figsize=(20, 15))
56     plt.plot(maxrange, cross_val_l, marker='o', label='CV MSE')
57     plt.plot(maxrange, train_mse_l, marker='o', label='Train MSE')
58     plt.plot(maxrange, test_mse_l, marker='o', label='Test MSE')
59     plt.xlabel('Max Depth')
60     plt.ylabel('Mean Squared Error')
61     plt.title('Cross-Validation, Training, Test')
62     plt.legend()
63     plt.show()
64
65     return train_mse_l, test_mse_l, cross_val_l
66
67 def doin_bagging(xtrain, ytrain, xtest, ytest):
68     ''' so basically do a tree then bagging?'''
69     tree_reg = DecisionTreeRegressor()
70     bag_reg = BaggingRegressor(base_estimator=tree_reg, n_estimators=1000)
71     bag_reg.fit(xtrain, ytrain)
72     ytestpred = bag_reg.predict(xtest)
73     testmse = mean_squared_error(ytest, ytestpred)
74     print("bagging mse", testmse)
75     #test_mse, feature_importances = bagging_regression(xtrain, ytrain, xtest, ytest)
76
77 def forest_stuff(xtrain, ytrain, xtest, ytest, n_estimators=1000, max_features_values=None):
78     test_mse_l = []
79     feature_importances = []
80
81     for max_features in max_features_values:
82         rf_reg = RandomForestRegressor(n_estimators=n_estimators, max_features=max_features, random_state=42)
83         rf_reg.fit(xtrain, ytrain)
84         y_test_pred = rf_reg.predict(xtest)
85         test_mse = mean_squared_error(ytest, y_test_pred)
86         test_mse_l.append(test_mse)
87         feature_importances.append(rf_reg.feature_importances_)
88
89
90

```

```

91 plt.figure(figsize=(10, 6))
92 plt.plot(max_features_values, test_mse_l, marker='o')
93 plt.xlabel('Max Features')
94 plt.ylabel('Test Mean Squared Error')
95 plt.title('Effect of Max Features on Test MSE')
96 plt.show()
97 for idx, max_features in enumerate(max_features_values):
98     print(f"Max Features: {max_features}")
99     print("Feature Importances:", feature_importances[idx])
100     print()
101 return test_mse_l, feature_importances
102
103
104
105 def bart(train_x, train_y, test_x, test_y, num_trees=100, num_burn_in=100, num_iterations=1000):
106     pass
107     ''' note: never really used it before: heavily influenced from online tutorials such as https://allenai.github.io/pybart/'''
108     train_x = train_x.astype(float32)
109     train_y = train_y.astype(float32)
110     test_x = test_x.astype(float32)
111     test_y = test_y.astype(float32)
112     model = Model()
113     model.num_trees = num_trees
114     model.num_burn_in = num_burn_in
115     model.num_iterations = num_iterations
116
117     model.fit(train_x, train_y)
118
119     y_test_pred = model.predict(test_x)
120
121     test_mse = np.mean((test_y - y_test_pred) ** 2)
122     print("Bart Test MSE:", test_mse)
123
124     return test_mse
125
126
127
128 def main():
129     x_a, y_v = data_loader("Carseats.csv")
130
131     train_x_a = x_a[:300]
132     train_y_v = y_v[:300]
133     test_x_a = x_a[300:]
134     test_y_v = y_v[300:]
135
136     #turn this back on later for decs tree pic
137     tree_things(train_x_a, train_y_v)
138
139
140     maxrange = range(1, 101)
141     crossval_for_tree(train_x_a, train_y_v, test_x_a, test_y_v, maxrange)
142     doin_bagging(train_x_a, train_y_v, test_x_a, test_y_v)
143     max_features_values = [None, 0.2, 0.4, 0.6, 0.8]
144     test_mse, feature_importances = forest_stuff(train_x_a, train_y_v, test_x_a, test_y_v, max_features_values=max_features_values)
145     #test_mse = bart(train_x_a, train_y_v, test_x_a, test_y_v)
146
147
148 if __name__ == '__main__':
149     main()

```