

Relatório do Terceiro Exercício

Délio Gomes Soares - 188947

24/10/2016

1 Atividade Proposta

Use os dados do **dataset SECOM do UCI**. O arquivo **secom.data** contém os dados. O arquivo **secom_labels.data** contém (na 1ª coluna) a classe de cada dado.

Usando um 5-fold externo para calcular a acurácia, e um 3-fold interno para a escolha dos hiperparâmetros, determine qual algoritmo entre *kNN*, *SVM* com *kernel RBF*, *redes neurais*, *Random Forest*, e *Gradient Boosting Machine* tem a maior acurácia.

1. Preprocesse os dados do arquivo: **Substitua os dados faltantes pela média da coluna (imputação pela média)**. Finalmente padronize as colunas para média 0 e desvio padrão 1.
2. Para o *kNN*, faça um *PCA* que mantém 80% da variância. Busque os valores do **k** entre os valores 1, 5, 11, 15, 21, 25.
3. Para o *SVM RBF* teste para **C** = 2**(-5), 2**(0), 2**(5), 2**(10) e **gamma** = 2**(-15), 2**(-10), 2**(-5), 2**(0), 2**(5).
4. Para a *rede neural*, teste com 10, 20, 30 e 40 **neurônios na camada escondida**.
5. Para o *RF*, teste com **mtry** ou **n_feats** = 10, 15, 20, 25 e **ntrees** = 100, 200, 300 e 400.
6. Para o *GBM* (ou *XGB*) teste para **número de árvores** = 30, 70, e 100, com **learning rate** de 0.1 e 0.05, e **profundidade da árvore** =

5. Você pode tanto usar alguma versão do *gbm* para R ou *SKlearn*, ou usar o *XGBoost* (para ambos).
7. Você não precisa fazer os *loops* da validação cruzada explicitamente. Pode usar as funções como *tunegrid* (do *caret*) ou *tuneParams* (do *mlr*) ou *GridSearchCV* do *SKlearn*.
8. Reporte a acurácia de cada algoritmo calculada pelo 5-fold CV externo.

2 Solução

2.1 Resposta

A acurácia média de cada algoritmo foi:

1. kNN: 92,98%
2. SVM: 93,36%
3. Rede Neural: 78,96%
4. Random Forest: 93,36%
5. GBM: 84,44%

2.2 Saída do código

A saída produzida pelo código em *python* com a acurácia média de cada algoritmo.

```
Acuracia media de cada um dos Classificadores
knn:      0.929811912242
SVM:      0.933633568293
NN:       0.792793904406
RF:       0.932996625618
GBM:      0.843163249759
```

2.3 Código fonte em *python*

```
1  # -*- coding: utf-8 -*-
2  """
3  Exercício 3
4  Autor: Delio Gomes Soares
5  RA: 188947
6  Data: 24/10/2016
7  """
8
9  import numpy as np
10 from sklearn.neighbors import KNeighborsClassifier
11 from sklearn.svm import SVC
12 from sklearn.neural_network import MLPClassifier
13 from sklearn.ensemble import RandomForestClassifier
14 from sklearn.ensemble import GradientBoostingClassifier
15 from sklearn.decomposition import PCA
16 from sklearn.model_selection import StratifiedKFold
17 from sklearn.model_selection import GridSearchCV
18 from sklearn.preprocessing import Imputer
19 from sklearn import preprocessing
20
21 #-----#
22 data = np.genfromtxt('secom.txt', delimiter=' ')
23 data_labels = np.genfromtxt('secom_labels.txt', delimiter=' ')
24 labels = data_labels[:,0]
25 #-----#
26 def featureImputation(X):
27     imp = Imputer(missing_values='NaN', strategy='mean', axis=0)
28     X_new = imp.fit_transform(X)
29     return X_new
30 #-----#
31 def featureStandardization(X):
32     X_new = preprocessing.scale(X)
33     return X_new
34 #-----#
35 def PCATransform(X):
36     pca = PCA(n_components=0.80)
37     XTransform = pca.fit_transform(X)
38     return XTransform
39 #-----#
40 def knn(X, y):
41     neigh = KNeighborsClassifier()
42     parameters = {'n_neighbors':[1, 5, 11, 15, 21, 25]}
43     best_n = 0
44     clf = GridSearchCV(neigh, parameters, cv=3)
```

```

45     clf.fit(X, y)
46     best_n = clf.best_estimator_.n_neighbors
47     return best_n
48 #-----#
49 def svm(X, y):
50     svc = SVC(kernel='rbf')
51     parameters = {'C': [2**-5, 2**0, 2**5, 2**10],
52                  'gamma': [2**-15, 2**-10, 2**-5, 2**0, 2**5]}
53     best_c = best_gamma = 0
54     clf = GridSearchCV(svc, parameters, cv=3)
55     clf.fit(X, y)
56     best_c = clf.best_estimator_.C
57     best_gamma = clf.best_estimator_.gamma
58     return best_c, best_gamma
59 #-----#
60 def neuralNetwork(X, y):
61     mlp = MLPClassifier(solver='lbfgs')
62     parameters = {'hidden_layer_sizes': [10, 20, 30, 40]}
63     best_neuron = 0
64     clf = GridSearchCV(mlp, parameters, cv=3)
65     clf.fit(X, y)
66     best_neuron = clf.best_estimator_.hidden_layer_sizes
67     return best_neuron
68 #-----#
69 def randomForest(X, y):
70     rf = RandomForestClassifier()
71     parameters = {'max_features': [10, 15, 20, 25],
72                  'n_estimators': [100, 200, 300, 400]}
73     best_max = 0
74     best_nEst = 0
75     clf = GridSearchCV(rf, parameters, cv=3)
76     clf.fit(X, y)
77     best_max = clf.best_estimator_.max_features
78     best_nEst = clf.best_estimator_.n_estimators
79     return best_max, best_nEst
80 #-----#
81 def gbm(X, y):
82     gbmc = GradientBoostingClassifier(max_depth=5)
83     parameters = {'learning_rate': [0.1, 0.05], 'n_estimators':
84                  [30, 70, 100]}
85     best_lR = best_nEst = 0
86     clf = GridSearchCV(gbmc, parameters, cv=3)
87     clf.fit(X, y)
88     best_lR = clf.best_estimator_.learning_rate
89     best_nEst = clf.best_estimator_.n_estimators

```

```

    return best_lR , best_nEst
89 #-----#
def main(X, y):
91     skf5PCA = StratifiedKFold(n_splits=5)
    XPCA = PCATransform(X)
93
    resultKNN = 0
95     for train_indexInt , test_indexInt in skf5PCA.split(XPCA, y):
        n_Neig = knn(PCATransform(XPCA[train_indexInt]), y[
train_indexInt])
97         neigh = KNeighborsClassifier(n_neighbors= n_Neig)
        neigh.fit(XPCA[train_indexInt], y[train_indexInt])
99         resultKNN += neigh.score(XPCA[test_indexInt], y[
test_indexInt])

101     skf5 = StratifiedKFold(n_splits=5)
    resultSVM = resultNeural = resultRF = resultGBM = 0
103     for train_indexInt , test_indexInt in skf5.split(X, y):

105         C_svm, Gamma_svm = svm(X[train_indexInt], y[
train_indexInt])
        neuron = neuralNetwork(X[train_indexInt], y[
train_indexInt])
107         maxFea, nEst = randomForest(X[train_indexInt], y[
train_indexInt])
        lR, nEstimators = gbM(X[train_indexInt], y[
train_indexInt])

109
        svc = SVC(kernel='rbf',C= C_svm, gamma= Gamma_svm)
111         svc.fit(X[train_indexInt], y[train_indexInt])
        resultSVM += svc.score(X[test_indexInt], y[test_indexInt
])

113
        mlp = MLPClassifier(solver='lbfgs', hidden_layer_sizes=
neuron)
115         mlp.fit(X[train_indexInt], y[train_indexInt])
        resultNeural += mlp.score(X[test_indexInt], y[
test_indexInt])

117
        rf = RandomForestClassifier(max_features= maxFea,
119                                     n_estimators= nEst)
        rf.fit(X[train_indexInt], y[train_indexInt])
121         resultRF += rf.score(X[test_indexInt], y[test_indexInt])

123
        gbmc = GradientBoostingClassifier(max_depth=5,

```

```

125                                     learning_rate = lR,
                                     n_estimators =
nEstimators)
    gbmc.fit(X[train_indexInt], y[train_indexInt])
127    resultGBM += gbmc.score(X[test_indexInt], y[
test_indexInt])

129    print("Acuracia media de cada um dos Classificadores")
    mediaKNN = resultKNN/5
131    print("\tkNN:\t", mediaKNN)
    mediaSVM = resultSVM/5
133    print("\tSVM:\t", mediaSVM)
    mediaNeural = resultNeural/5
135    print("\tNN:\t", mediaNeural)
    mediaRF = resultRF/5
137    print("\tRF:\t", mediaRF)
    mediaGBM = resultGBM/5
139    print("\tGBM:\t", mediaGBM)

141 dataImp = featureImputation(data)
    dataStan = featureStandardization(dataImp)
143 main(dataStan, labels)

```