

Enhancing Accuracy with AdaBoost Ensemble: A Comparison of Decision Tree and Fused Model Approaches"

Report:

TASK 1 -

The task is to compare the accuracies of two classifiers: a fused model and an AdaBoost ensemble with Decision Tree as the base learner. The classifiers should be trained using the first 1000 instances and tested using the remaining 3601 instances. The results of the comparison should be reported along with observations and conclusions, supported by evidence.

To complete this task, first, the training and testing data files should be created by splitting the data into two parts, with 1000 instances used for training and 3601 instances used for testing. Then, the two classifiers should be trained and tested on the respective data sets. Finally, the accuracy of each classifier should be calculated and compared, and observations and conclusions should be drawn based on the results.

Below is the accuracy of the fused model.

```
#clf.fit(spam_training_data,spam_training_target)
[118] eclf.fit(spam_training_data,spam_training_target)

VotingClassifier(estimators=[('DT',
                             DecisionTreeClassifier(class_weight='balanced',
                                                    criterion='entropy',
                                                    max_depth=1000,
                                                    max_features='sqrt',
                                                    min_samples_split=10,
                                                    random_state=42)),
                             ('LR',
                             LogisticRegression(class_weight='none',
                                                fit_intercept=1.0,
                                                intercept_scaling=1.0,
                                                max_iter=3000,
                                                random_state=42)),
                             ('GNB', GaussianNB())])

[119] spam_test_target_predict=eclf.predict(spam_test_data)

[120] accuracy_score(spam_test_target,spam_test_target_predict)
0.9130797000833102

❶ confusion_matrix(spam_test_target,spam_test_target_predict)
❷ array([[1963,  222],
       [ 91, 1325]])

❸ classification_report(spam_test_target,spam_test_target_predict)
          precision    recall  f1-score   support
ham      0.96      0.90      0.93    2185\n    spam
0.86      0.94      0.89    1416\n    accuracy
0.91      3601\nweighted avg      0.92      0.91      3601\n
```

Below is the accuracy of the adaboost ensemble with decision tree as the base learner (Note: To use decision tree as the base learner, the base estimator should be set as none, If None, then the base estimator is DecisionTreeClassifier)

```

[109] from sklearn.naive_bayes import GaussianNB
from sklearn.ensemble import AdaBoostClassifier, VotingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
from sklearn.metrics import accuracy_score
#clf = RandomForestClassifier(n_estimators=1000)
clf_dt = DecisionTreeClassifier(criterion = "entropy")
#clf_dt = DecisionTreeClassifier(min_samples_split=10, max_depth = 1000, criterion = "entropy", splitter = "best", random_state=42, class_weight = "balanced", min_weight_fraction_leaf=0.001)
#clf = KNeighborsClassifier(n_neighbors=7)
clf = AdaBoostClassifier(n_estimators = 200, base_estimator=None)
#LRI = LogisticRegression()
#clf_lr = LogisticRegression(solver='lbfgs', penalty = "l2", dual = False, tol = 1e-4, C = 1.0, fit_intercept = 1.0, intercept_scaling = 1.0)
#clf_gnb = GaussianNB()
#clf = AdaBoostClassifier(n_estimators = 200,base_estimator=LRI)
#clf = VotingClassifier(estimators = [('DT', clf_dt),('LR',clf_lr),('GNB',clf_gnb)], voting = 'hard')
clf.fit(spam_training_data,spam_training_target)
#clf.fit(spam_training_data,spam_training_target)

AdaBoostClassifier(n_estimators=200)

[110] spam_test_target_predict=clf.predict(spam_test_data)

[111] accuracy_score(spam_test_target,spam_test_target_predict)
0.9208553179672313

[112] confusion_matrix(spam_test_target,spam_test_target_predict)
array([[2007, 178],
       [107, 1309]])

```

TASK 2 -

The task is to compare the accuracy of two classifiers: a fused model and a Random Forest classifier with 1000 base learners.

The classifiers will be trained using the first 1000 instances and tested using the remaining 3601 instances.

Below is the accuracy of the fused model.

```

+ Code + Text
[118] #clf.fit(spam_training_data,spam_training_target)
eclf.fit(spam_training_data,spam_training_target)

VotingClassifier( estimators=[('DT',
                               DecisionTreeClassifier(class_weight='balanced',
                                                      criterion='entropy',
                                                      max_depth=1000,
                                                      max_features='sqrt',
                                                      min_samples_split=10,
                                                      random_state=42)),
                           ('LR',
                               LogisticRegression(class_weight='none',
                                                      fit_intercept=1.0,
                                                      intercept_scaling=1.0,
                                                      max_iter=3000,
                                                      random_state=42)),
                           ('GNB',
                               GaussianNB())))

[119] spam_test_target_predict=eclf.predict(spam_test_data)

[120] accuracy_score(spam_test_target,spam_test_target_predict)
0.9130797000833102

[121] confusion_matrix(spam_test_target,spam_test_target_predict)
array([[1963, 222],
       [ 91, 1325]])

[122] classification_report(spam_test_target,spam_test_target_predict)
          precision    recall  f1-score   support
ham      0.96      0.90      0.93    2185
spam     0.86      0.94      0.89    1416
          accuracy         0.91      0.91    3601
          weighted avg     0.91      0.92      0.91    3601

```

Below is the accuracy of the random forest classifier with 1000 base learners:

The screenshot shows a Jupyter Notebook interface with the following code and output:

```
clf = RandomForestClassifier(n_estimators=1000)
#clf = DecisionTreeClassifier(criterion = "entropy")
#clf_dt = DecisionTreeClassifier(min_samples_split=10, max_depth = 1000, criterion = "entropy", splitter = "best", random_state=42,
#splitter = "best", max_depth = 10, min_samples_split=100, min_samples_leaf=1, random_state=42, class_weight = "balanced", min_weight_fraction_leaf=0.001)
#clf = KNeighborsClassifier(n_neighbors=7)
#clf = AdaBoostClassifier(n_estimators = 200, base_estimator=None)
#LRI = LogisticRegression()
#clf_lr = LogisticRegression(solver='lbfgs', penalty = "l2", dual = False, tol = 1e-4, C = 1.0, fit_intercept = 1.0, intercept_scaling = 1.0, class_weight = "balanced", random_state=42)
#clf_gnb = GaussianNB()
#clf = AdaBoostClassifier(n_estimators = 200,base_estimator=LRI)
#clfif = VotingClassifier(estimators = [('DT', clf_dt),('LR',clf_lr),('GNB',clf_gnb)], voting ='hard')
clf.fit(spam_training_data,spam_training_target)
#clfif.fit(spam_training_data,spam_training_target)
```

[126] spam_test_target_predict=clf.predict(spam_test_data)

[127] accuracy_score(spam_test_target,spam_test_target_predict)
0.9327964454318245

[128] confusion_matrix(spam_test_target,spam_test_target_predict)
array([[2061, 124],
 [118, 1298]])

[129] classification_report(spam_test_target,spam_test_target_predict)

	precision	recall	f1-score	support	ham	0.95	0.94	0.94	2185	spam	0.93
0.91	0.92	0.91	1416	accuracy	0.93	3601	3601	3601	0.93	0.93	0.93
0.93	3601	weighted avg	0.93	0.93	3601						

TASK 3-

The task is to study the impact of training sample size on the accuracy of two classifiers: a fused classifier and an AdaBoost ensemble with Decision Tree as the base learner. The study will involve comparing their accuracy with four different training-test splits: 50%-50%, 60%-40%, 70%-30%, and 80%-20%.

Below is the accuracy of the fused classifier with the 50%-50% training test split

The screenshot shows a Jupyter Notebook interface with the following code and output:

```
eclf = VotingClassifier(estimators = [( 'DT', clf_dt),('LR',clf_lr),('GNB',clf_gnb)], voting ='hard')
clf.fit(spam_training_data,spam_training_target)
#clfif.fit(spam_training_data,spam_training_target)
```

VotingClassifier(estimators=[('DT',
DecisionTreeClassifier(class_weight='balanced',
criterion='entropy',
max_depth=1000,
max_features='sqrt',
min_samples_split=10,
random_state=42)),
('LR',
LogisticRegression(class_weight='none',
fit_intercept=1.0,
intercept_scaling=1.0,
max_iter=3000,
random_state=42)),
('GNB', GaussianNB())))

[193] spam_test_target_predict=eclf.predict(spam_test_data)

[194] accuracy_score(spam_test_target,spam_test_target_predict)
0.9243478260869565

[195] confusion_matrix(spam_test_target,spam_test_target_predict)
array([[11276, 125],
 [49, 850]])

[196] classification_report(spam_test_target,spam_test_target_predict)

	precision	recall	f1-score	support	ham	0.96	0.91	0.94	1401	spam	0.93
0.87	0.95	0.91	0.99	accuracy	0.92	2300	2300	2300	0.92	0.92	0.93
0.92	2300	weighted avg	0.93	0.92	0.92						

Below is the accuracy of the Ada boost Ensemble with decision tree as the base learner with the 50%-50% training test split

```

Untitled0.ipynb
File Edit View Insert Runtime Tools Help Saving... Comment Share RAM Disk M
+ Code + Text
[28]: from sklearn.metrics import accuracy_score
#clf = RandomForestClassifier(n_estimators=1000)
clf = DecisionTreeClassifier(criterion = "entropy")
clf_dt = DecisionTreeClassifier(min_samples_split=10, max_depth = 1000, criterion = "entropy", splitter = "best", random_state=42, criterion="entropy")
#splited = "best", max_depth = 10, min_samples_split=100, min_samples_leaf=1, random_state=42, class_weight = "balanced", min_weight_fraction_leaf=0.001
#clf = KNeighborsClassifier(n_neighbors=7)
clf = AdaBoostClassifier(n_estimators = 200, base_estimator=None)
#LRI = LogisticRegression()
#clf_lr = LogisticRegression(solver='lbfgs', penalty = "l2", dual = False, tol = 1e-4, C = 1.0, fit_intercept = 1.0, intercept_scaling = 1.0, class_weight = "balanced")
#clf_gnb = GaussianNB()
#clf = AdaBoostClassifier(n_estimators = 200,base_estimator=LRI)
#clf = VotingClassifier(estimators = [('DT', clf_dt),('LR',clf_lr),('GNB',clf_gnb)], voting = 'hard')
clf.fit(spam_training_data,spam_training_target)
#eclf.fit(spam_training_data,spam_training_target)

AdaBoostClassifier(n_estimators=200)

[144]: spam_test_target_predict=clf.predict(spam_test_data)

[145]: accuracy_score(spam_test_target,spam_test_target_predict)
0.9304347826086956

[146]: confusion_matrix(spam_test_target,spam_test_target_predict)
array([[11303,    98],
       [   62,  837]])

[147]: classification_report(spam_test_target,spam_test_target_predict)
      precision    recall  f1-score   support
ham       0.90      0.93      0.91     899
accuracy                           0.93
macro avg       0.90      0.93      0.91    2300
weighted avg       0.90      0.93      0.91    2300

1s completed at 4:37 PM

```

Below is the accuracy of the fused classifier with the 60%-40% training test split

```

File Edit View Insert Runtime Tools Help
+ Code + Text
[164]: eclf = VotingClassifier(estimators = [('DT', clf_dt),('LR',clf_lr),('GNB',clf_gnb)], voting = 'hard')
#clf.fit(spam_training_data,spam_training_target)
eclf.fit(spam_training_data,spam_training_target)

VotingClassifier(estimators=[('DT',
DecisionTreeClassifier(class_weight='balanced',
criterion='entropy',
max_depth=1000,
max_features='sqrt',
min_samples_split=10,
random_state=42)),
('LR',
LogisticRegression(class_weight='none',
fit_intercept=1.0,
intercept_scaling=1.0,
max_iter=3000,
random_state=42)),
('GNB', GaussianNB())))

[165]: spam_test_target_predict=eclf.predict(spam_test_data)

[166]: accuracy_score(spam_test_target,spam_test_target_predict)
0.9266702878870179

[167]: confusion_matrix(spam_test_target,spam_test_target_predict)
array([[1037,    97],
       [   38,  669]])

[168]: classification_report(spam_test_target,spam_test_target_predict)
      precision    recall  f1-score   support
ham       0.96      0.91      0.94     707
accuracy                           0.93
macro avg       0.96      0.91      0.94    1134
weighted avg       0.96      0.91      0.94    1134

1s completed at 4:40 PM

```

Below is the accuracy of the Ada boost Ensemble with decision tree as the base learner with the 60%-40% training test split

```
#clf = RandomForestClassifier(n_estimators=1000)
#clf_dt = DecisionTreeClassifier(criterion = "entropy")
#clf_dt = DecisionTreeClassifier(min_samples_split=10, max_depth = 1000, criterion = "entropy", splitter = "best", random_state=42, 
#splitter = "best", max_depth = 10, min_samples_split=100, min_samples_leaf=1, random_state=42, class_weight = "balanced", min_weight_fraction_leaf=0.001)
#clf = KNeighborsClassifier(n_neighbors=7)
clf = AdaBoostClassifier(n_estimators = 200, base_estimator=None)
#LRI = LogisticRegression()
#clf_lr = LogisticRegression(solver='lbfgs', penalty = 'l2', dual = False, tol = 1e-4, C = 1.0, fit_intercept = 1.0, intercept_scaling = 1.0, class_weight = "balanced")
#clf_gnb = GaussianNB()
#clf = AdaBoostClassifier(n_estimators = 200,base_estimator=LRI)
#clf = VotingClassifier(estimators = [('DT', clf_dt),('LR',clf_lr),('GNB',clf_gnb)], voting = 'hard')
clf.fit(spam_training_data,spam_training_target)
#eclf.fit(spam_training_data,spam_training_target)

AdaBoostClassifier(n_estimators=200)

spam_test_target_predict=clf.predict(spam_test_data)
accuracy_score(spam_test_target,spam_test_target_predict)
0.9342748506246605

confusion_matrix(spam_test_target,spam_test_target_predict)
array([[1061,    73],
       [ 48,  659]])

classification_report(spam_test_target,spam_test_target_predict)
          precision    recall  f1-score   support
ham      0.96      0.94      0.95     1134
spam     0.93      0.93      0.93     1841
macro avg  0.94      0.93      0.93     3134
weighted avg  0.94      0.93      0.93     3134
```

Below is the accuracy of the fused classifier with the 70%-30% training test split

```
#clf.fit(spam_training_data,spam_training_target)
#eclf.fit(spam_training_data,spam_training_target)

VotingClassifier(estimators=[('DT',
                             DecisionTreeClassifier(class_weight='balanced',
                                                    criterion='entropy',
                                                    max_depth=1000,
                                                    max_features='sqrt',
                                                    min_samples_split=10,
                                                    random_state=42)),
                            ('LR',
                             LogisticRegression(class_weight='none',
                                                fit_intercept=1.0,
                                                intercept_scaling=1.0,
                                                max_iter=3000,
                                                random_state=42)),
                            ('GNB', GaussianNB())])

spam_test_target_predict=eclf.predict(spam_test_data)
accuracy_score(spam_test_target,spam_test_target_predict)
0.925673090649536

confusion_matrix(spam_test_target,spam_test_target_predict)
array([[773,    72],
       [ 32,  524]])

classification_report(spam_test_target,spam_test_target_predict)
          precision    recall  f1-score   support
ham      0.96      0.91      0.94     845
spam     0.92      0.93      0.93     1401
macro avg  0.94      0.92      0.93     2246
weighted avg  0.94      0.92      0.93     2246
```

Below is the accuracy of the Ada boost Ensemble with decision tree as the base learner with the 70%-30% training test split

```
+ Code + Text
[153] from sklearn.metrics import accuracy_score
#clf = RandomForestClassifier(n_estimators=1000)
clf = DecisionTreeClassifier(criterion = "entropy")
clf_dt = DecisionTreeClassifier(min_samples_split=10, max_depth = 1000, criterion = "entropy", splitter = "best", random_state=42, c
#splitter = "best", max_depth = 10, min_samples_split=100, min_samples_leaf=1, random_state=42, class_weight = "balanced", min_weigl
#clf = KNeighborsClassifier(n_neighbors=7)
clf = AdaBoostClassifier(n_estimators = 200, base_estimator=None)
#LRI = LogisticRegression()
#clf_lr = LogisticRegression(solver='lbfgs', penalty = "l2", dual = False, tol = 1e-4, C = 1.0, fit_intercept = 1.0, intercept_sca
#clf_gnb = GaussianNB()
#clf = AdaBoostClassifier(n_estimators = 200,base_estimator=LRI)
#clf = VotingClassifier(estimators = [('DT', clf_dt),('LR',clf_lr),('GNB',clf_gnb)], voting ='hard')
clf.fit(spam_training_data,spam_training_target)
#clf.fit(spam_training_data,spam_training_target)

AdaBoostClassifier(n_estimators=200)

[154] spam_test_target_predict=clf.predict(spam_test_data)

[155] accuracy_score(spam_test_target,spam_test_target_predict)
0.9307637401855817

[156] confusion_matrix(spam_test_target,spam_test_target_predict)
array([[783,  62],
       [ 35, 521]])

[157] classification_report(spam_test_target,spam_test_target_predict)
          precision    recall  f1-score   support
ham      0.96      0.93      0.94      845
spam     0.89      0.94      0.91      556
           accuracy         macro avg
           0.93      1401
           weighted avg
               0.93      0.93      1401
```

Below is the accuracy of the fused classifier with the 80%-20% training test split

```
[186] #clf = VotingClassifier(estimators = [('DT', clf_dt),('LR',clf_lr),('GNB',clf_gnb)], voting = 'hard')
#clf.fit(spam_training_data,spam_training_target)
clf.fit(spam_training_data,spam_training_target)

votingClassifier(estimators=[('DT',
    DecisionTreeClassifier(class_weight='balanced',
        criterion='entropy',
        max_depth=1000,
        max_features='sqrt',
        min_samples_split=10,
        random_state=42)),
    ('LR',
    LogisticRegression(class_weight='none',
        fit_intercept=1.0,
        intercept_scaling=1.0,
        max_iter=3000,
        random_state=42)),
    ('GNB', GaussianNB())))

[187] spam_test_target_predict=eclf.predict(spam_test_data)

[188] accuracy_score(spam_test_target,spam_test_target_predict)
0.9294245385450597

[189] confusion_matrix(spam_test_target,spam_test_target_predict)
array([[512,  43],
       [ 22, 344]])

[190] classification_report(spam_test_target,spam_test_target_predict)
          precision    recall  f1-score   support
ham      0.96      0.92      0.94      555
spam     0.89      0.94      0.91      366
           accuracy         macro avg
           0.93      921
           weighted avg
               0.93      0.93      921
```

Below is the accuracy of the Ada boost Ensemble with decision tree as the base learner with the 80%-20% training test split

```
[180] from sklearn.metrics import accuracy_score
      #clf = RandomForestClassifier(n_estimators=1000)
      clf = DecisionTreeClassifier(criterion = "entropy")
      clf_dt = DecisionTreeClassifier(min_samples_split=10, max_depth = 1000, criterion = "entropy", splitter = "best", random_state=42, class_weight="balanced", min_weight_fraction_in_leaf=0.001)
      #splitter = "best", max_depth = 10, min_samples_split=100, min_samples_leaf=1, random_state=42, class_weight = "balanced", min_weight_fraction_in_leaf=0.001
      #clf = KNeighborsClassifier(n_neighbors=7)
      clf = AdaBoostClassifier(n_estimators = 200, base_estimator=None)
      #LRI = LogisticRegression()
      #clf_lr = LogisticRegression(solver='lbfgs', penalty = "l2", dual = False, tol = 1e-4, C = 1.0, fit_intercept = 1.0, intercept_scaling=1, class_weight="balanced", random_state=42)
      #clf_gnb = GaussianNB()
      #clf = AdaBoostClassifier(n_estimators = 200,base_estimator=LRI)
      #eclf = VotingClassifier(estimators = [('DT', clf_dt),('LR',clf_lr),('GNB',clf_gnb)], voting ='hard')
      clf.fit(spam_training_data,spam_training_target)
      #eclf.fit(spam_training_data,spam_training_target)

      AdaBoostClassifier(n_estimators=200)

[181] spam_test_target_predict=clf.predict(spam_test_data)

[182] accuracy_score(spam_test_target,spam_test_target_predict)
> 0.9283387622149837

[183] confusion_matrix(spam_test_target,spam_test_target_predict)
> array([[512,  43],
       [ 23, 348]])

[184] classification_report(spam_test_target,spam_test_target_predict)
>
      precision    recall  f1-score   support
ham        0.96      0.92      0.94      555
spam       0.89      0.94      0.91      366
          avg       0.93      0.93      0.93      921
      accuracy         0.93
      weighted avg       0.93      0.93      0.93      921
```