

CS 573: Assignment 4

Tiantu Xu

1. Preprocessing

In Question 1, data pre-processing is run by the command line below:

```
$ python preprocess-assg4.py
```

The output is trainingSet.csv and testSet.csv.

2. Implement Logistic Regression and Linear SVM

To train and test decision tree, specify `sys.argv[3] = 1`:

```
$ python trees.py trainingSet.csv testSet.csv 1
```

The output from my code is

```
Training Accuracy DT: 0.77
Test Accuracy DT: 0.72
```

To train and test bagging, specify `sys.argv[3] = 2`:

```
$ python trees.py trainingSet.csv testSet.csv 2
```

The output from my code is

```
Training Accuracy BT: 0.78
Test Accuracy BT: 0.75
```

To train and test random forest, specify `sys.argv[3] = 3`:

```
$ python trees.py trainingSet.csv testSet.csv 3
```

The output from my code is

```
Training Accuracy RF: 0.76
Test Accuracy RF: 0.73
```

3. The Influence of Tree Depth on Classifier Performance

(a) K-fold cross-validation is run on the command line below.

```
$ python cv_depth.py
```

The model performance on 3 models (DT, BT, and RF) is shown below.

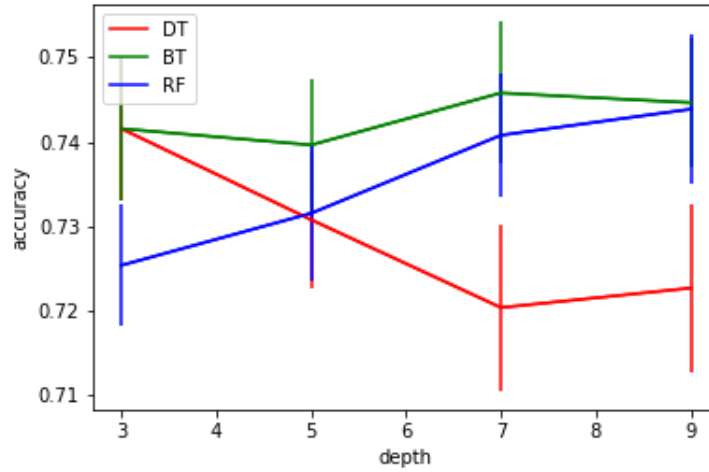


Figure 1: The model performance of DT, BT, and RF

(b) The hypothesis testing is formulated as

H_0 : DT and RF model performances do not differ significantly.

H_1 : DT and RF model performances differ significantly.

Assume I have a significance level of $\alpha = 0.05$. ttest is run on the performance numbers obtained in the above cross-validation. The output from the ttest is shown below

```
d = 3 Ttest_indResult(statistic=1.37313899108, pvalue=0.186573940678)
d = 5 Ttest_indResult(statistic=-0.0644900370588, pvalue=0.94929085751)
d = 7 Ttest_indResult(statistic=-1.566293430098, pvalue=0.1346910261630)
d = 9 Ttest_indResult(statistic=-1.502921737574, pvalue=0.1502010577598)
```

It turns out that for every tree depth, $p\text{-value} > \alpha$, so that we **fail to reject** the null hypothesis H_0 that DT and RF performances do not differ significantly.

4. Compare Performance of Different Models

(a) K-fold cross-validation is run on the command line below.

```
$ python cv_frac.py
```

The model performance on 3 models (DT, BT, and RF) is shown below.

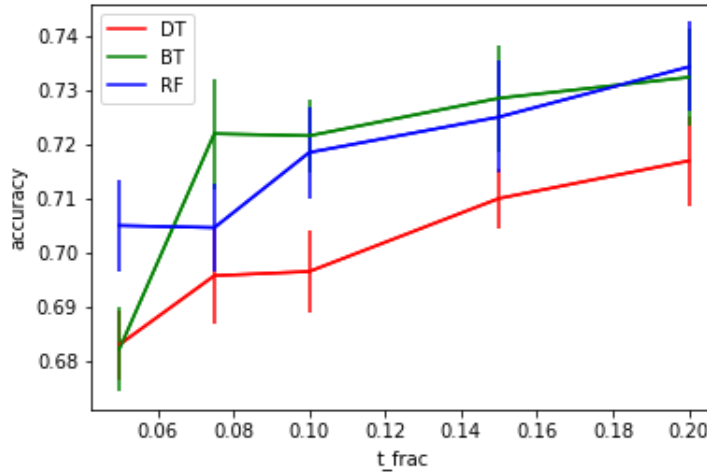


Figure 2: The model performance of DT, BT, and RF on the fraction of training data

(b) The hypothesis testing is formulated as

H_0 : BT and RF model performances do not differ significantly.

H_1 : BT and RF model performances differ significantly.

Assume I have a significance level of $\alpha = 0.05$. ttest is run on the performance numbers obtained in the above cross-validation. The output from the ttest is shown below

```
t_frac = 0.05 Ttest_indResult(statistic=-1.867300552201, pvalue=0.0782322885128)
t_frac = 0.075 Ttest_indResult(statistic=1.273528235141, pvalue=0.2190313955400)
t_frac = 0.1 Ttest_indResult(statistic=0.27350538350, pvalue=0.787578344338)
t_frac = 0.15 Ttest_indResult(statistic=0.2307692307692, pvalue=0.820096694805)
t_frac = 0.2 Ttest_indResult(statistic=-0.14930732766, pvalue=0.882971277311)
```

It turns out that $p\text{-value} > \alpha$, so that we **fail to reject** the null hypothesis H_0 that BT and RF performances do not differ significantly.

5. The Influence of Number of Trees on Classifier Performance

(a) K-fold cross-validation is run on the command line below.

```
$ python cv_numtrees.py
```

The model performance on 3 models (DT, BT, and RF) is shown below.

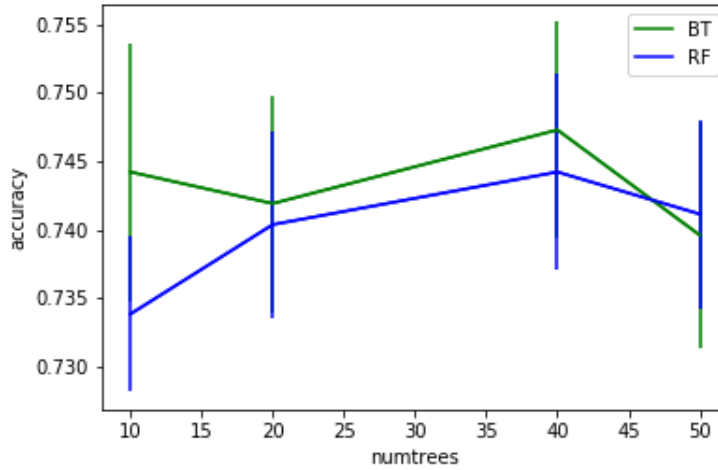


Figure 3: The model performance of DT, BT, and RF

(b) The hypothesis testing is formulated as

H_0 : BT and RF model performances do not differ significantly.

H_1 : BT and RF model performances differ significantly.

Assume I have a significance level of $\alpha = 0.05$. ttest is run on the performance numbers obtained in the above cross-validation. The output from the ttest is shown below

```
t = 10 Ttest_indResult(statistic=0.897730833027, pvalue=0.3811802797843)
t = 20 Ttest_indResult(statistic=0.139553753811, pvalue=0.890562201897)
t = 40 Ttest_indResult(statistic=0.274469502343, pvalue=0.786849103547)
t = 50 Ttest_indResult(statistic=-0.135164620935, pvalue=0.89398173860)
```

It turns out that $p\text{-value} > \alpha$, so that we **fail to reject** the null hypothesis H_0 that BT and RF performances do not differ significantly.