

CS-GY 6923 Machine Learning

Professor: Dr. Raman Kannan

Homework 3: Ensemble techniques

Hansheng Li

Contents

1. Start up	3
2. Cross Validation	3
3. Leave One out Cross Validation	7
4. Boosting	12
5. Original classification result	16
6. Comparison	17

1. Start up

We will continue to use the code from HW 02, so no more extra import needed. In this assignment, I have chosen three Ensemble Techniques which are Cross Validation, Leave One out Cross Validation, and Bagging. And we going to compare Variance, bias, and accuracy to see how much improvement have been made by Ensemble Techniques.

2. Cross Validation

Multinational Logistic Regression:

```
1. > #####
2. > ## cross - validation methods##
3. > #####
4. >
5. > # Multinational Logistic Regression #
6. >
7. > # training the model by assigning Class column
8. > # as target variable and rest other column
9. > # as independent variable
10.> log.cv <- train(Class ~., data = data,
11.+                 method = "multinom",
12.+                 trControl = trainControl(
13.+                 method = "CV",
14.+                 number = 10,
15.+                 verboseIter = T
16.+                 ))
17.+ Fold01: decay=0e+00
18.# weights: 126 (102 variable)
19.initial value 23839.345236
20.iter 10 value 18372.055625
21.iter 20 value 15025.622376
22.iter 30 value 10919.859504
23.....
24.iter 80 value 2793.901318
25.iter 90 value 2779.793130
26.iter 100 value 2772.876570
27.final value 2772.876570
28.stopped after 100 iterations
29.> print(log.cv)
30.Penalized Multinomial Regression
31.
32.13611 samples
```

```

33. 16 predictor
34. 7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROZ', 'SEKER', 'SIRA'
35.
36.No pre-processing
37.Resampling: Cross-Validated (10 fold)
38.Summary of sample sizes: 12251, 12251, 12249, 12249, 12250, 12248, ...
39.Resampling results across tuning parameters:
40.
41. decay Accuracy Kappa
42. 0e+00 0.9260157 0.9105411
43. 1e-04 0.9255748 0.9100060
44. 1e-01 0.9212403 0.9047508
45.
46.Accuracy was used to select the optimal model using the largest value.
47.The final value used for the model was decay = 0.
48.> log.cv_pred <- factor(predict(log.cv , data.valid), ordered = T)
49.> log.cv_cM <- table(log.cv_pred, data.valid$Class)
50.> log.cv_summary <-
  data.frame( Bias = mean(predict(log.cv, data.train) != data.train$Class
  ),
51.+                               Variance = mean(log.cv_pred != data.valid$Class),
52.+                               Accuracy = sum(diag(log.cv_cM)) / sum(log.cv_cM))
53.> log.cv_summary
54.      Bias Variance Accuracy
55. 1 0.07179595 0.07639569 0.9236043

```

SVM:

```

1. > # Support Vector Machine Classification
2. > svm.cv <- train(Class ~., data = data,
3. +               method = "svmLinear",
4. +               trControl = trainControl(
5. +                 method = "CV",
6. +                 number = 10,
7. +                 verboseIter = T
8. +               ))
9. + Fold01: C=1
10. - Fold01: C=1
11. ....
12. - Fold10: C=1
13. Aggregating results
14. Fitting final model on full training set
15. > print(svm.cv)

```

```

16. Support Vector Machines with Linear Kernel
17.
18. 13611 samples
19. 16 predictor
20. 7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROZ', 'SEKE
    R', 'SIRA'
21.
22. No pre-processing
23. Resampling: Cross-Validated (10 fold)
24. Summary of sample sizes: 12250, 12251, 12250, 12249, 12250, 12250, ...
25. Resampling results:
26.
27. Accuracy Kappa
28. 0.9262365 0.9107779
29.
30. Tuning parameter 'C' was held constant at a value of 1
31. > svm.cv_pred <- factor(predict(svm.cv , data.valid), ordered = T)
32. > svm.cv_cM <- table(svm.cv_pred, data.valid$Class)
33. > svm.cv_summary <-
    data.frame( Bias = mean(predict(svm.cv, data.train) != data.t
    rain$Class),
34. + Variance = mean(svm.cv_pred!= data.vali
    d$Class),
35. + Accuracy = sum(diag(svm.cv_cM)) / sum(s
    vm.cv_cM))
36. > svm.cv_summary
37. Bias Variance Accuracy
38. 1 0.07148105 0.07492654 0.9250735

```

Decision tree:

```

1. > # Decision Tree
2. > tree.cv <- train(Class ~., data = data,
3. + method = "rpart",
4. + trControl = trainControl(
5. + method = "CV",
6. + number = 10,
7. + verboseIter = T
8. + ))
9. + Fold01: cp=0.1647
10. - Fold01: cp=0.1647
11. - Fold10: cp=0.1647
12. Aggregating results
13. Selecting tuning parameters
14. Fitting cp = 0.165 on full training set
15. > print(tree.cv)

```

```

16. CART
17.
18. 13611 samples
19. 16 predictor
20. 7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROZ', 'SEKE
    R', 'SIRA'
21.
22. No pre-processing
23. Resampling: Cross-Validated (10 fold)
24. Summary of sample sizes: 12248, 12250, 12250, 12249, 12250, 12247, ...
25. Resampling results across tuning parameters:
26.
27.  cp          Accuracy   Kappa
28.  0.1647293  0.5931715  0.4864191
29.  0.1705912  0.4694109  0.3208308
30.  0.1992052  0.3471001  0.1409502
31.
32. Accuracy was used to select the optimal model using the largest value.
33. The final value used for the model was cp = 0.1647293.
34. > data.train$Class <- factor(data.train$Class, ordered = FALSE)
35. > data.valid$Class <- factor(data.valid$Class, ordered = FALSE)
36. > tree.cv_pred <- predict(tree.cv , data.valid, type = "raw")
37. > tree.cv_cM <- table(tree.cv_pred, data.valid$Class)
38. > tree.cv_summary <-
    data.frame( Bias = mean(predict(tree.cv, data.train, type = "
    raw") != data.train$Class),
39. +                               Variance = mean(tree.cv_pred != data.v
    alid$Class),
40. +                               Accuracy = sum(diag(tree.cv_cM)) / sum
    (tree.cv_cM))
41. > tree.cv_summary
42.      Bias  Variance  Accuracy
43. 1 0.4682481 0.4608227 0.5391773

```

KNN:

```

1. > # K-Nearest Neighbor
2. > knn.cv <- train(Class ~., data = data,
3. +               method = "knn",
4. +               trControl = trainControl(
5. +               method = "CV",
6. +               number = 10,
7. +               verboseIter = T
8. +               ))
9. + Fold01: k=5
10. - Fold01: k=5

```

```

11. + Fold01: k=7
12. - Fold01: k=7
13. ....
14. + Fold10: k=9
15. - Fold10: k=9
16. Aggregating results
17. Selecting tuning parameters
18. Fitting k = 5 on full training set
19. > print(knn.cv)
20. k-Nearest Neighbors
21.
22. 13611 samples
23.    16 predictor
24.    7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROZ', 'SEKE
      R', 'SIRA'
25.
26. No pre-processing
27. Resampling: Cross-Validated (10 fold)
28. Summary of sample sizes: 12248, 12249, 12250, 12250, 12250, 12251, ...
29. Resampling results across tuning parameters:
30.
31.  k  Accuracy  Kappa
32.  5  0.7321261  0.6746749
33.  7  0.7253680  0.6662656
34.  9  0.7228692  0.6630517
35.
36. Accuracy was used to select the optimal model using the largest value.
37. The final value used for the model was k = 5.
38. > knn.cv_pred <- predict(knn.cv , data.valid)
39. > knn.cv_cM <- table(knn.cv_pred, data.valid$Class)
40. > knn.cv_summary <-
      data.frame( Bias = mean(predict(knn.cv, data.train) != data.t
rain$Class),
41. +                                     Variance = mean(knn.cv_pred!= data.vali
d$Class),
42. +                                     Accuracy = sum(diag(knn.cv_cM)) / sum(k
nn.cv_cM))
43. > knn.cv_summary
44.      Bias Variance Accuracy
45. 1 0.1783353 0.1767875 0.8232125

```

3. Leave One out Cross Validation

Multinational Logistic Regression:

```

1. > log.loocv <- train(Class ~., data = data,

```

```

2. +             method = "multinom",
3. +             trControl = trainControl(
4. +                 method = "LGOCV",
5. +                 number = 10,
6. +                 p = 0.1,
7. +                 verboseIter = T
8. +             ))
9. + Resample01: decay=0e+00
10. # weights: 126 (102 variable)
11. initial value 2654.221443
12. iter 10 value 2031.935250
13. ....
14. final value 2780.792544
15. stopped after 100 iterations
16. > print(log.loocv)
17. Penalized Multinomial Regression
18.
19. 13611 samples
20. 16 predictor
21. 7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROS', 'SEKE
    R', 'SIRA'
22.
23. No pre-processing
24. Resampling: Repeated Train/Test Splits Estimated (10 reps, 10%)
25. Summary of sample sizes: 1364, 1364, 1364, 1364, 1364, 1364, ...
26. Resampling results across tuning parameters:
27.
28.  decay  Accuracy  Kappa
29.  0e+00  0.9182739  0.9011821
30.  1e-04  0.9184617  0.9014083
31.  1e-01  0.9153997  0.8976753
32.
33. Accuracy was used to select the optimal model using the largest value.
34. The final value used for the model was decay = 1e-04.
35. > log.loocv_pred <- predict(log.loocv , data.valid)
36. > log.loocv_cM <- table(log.loocv_pred, data.valid$Class)
37. > log.loocv_summary <-
    data.frame( Bias = mean(predict(log.loocv, data.train) != dat
    a.train$Class),
38. +             Variance = mean(log.loocv_pred != da
    ta.valid$Class),
39. +             Accuracy = sum(diag(log.loocv_cM)) /
    sum(log.loocv_cM))
40. > log.loocv_summary
41.      Bias  Variance  Accuracy
42. 1 0.07179595 0.07664055 0.9233595

```


SVM:

```
1. > # Support Vector Machine Classification
2. > svm.loocv <- train(Class ~., data = data,
3. +                     method = "svmLinear",
4. +                     trControl = trainControl(
5. +                         method = "LGOVCV",
6. +                         number = 10,
7. +                         p = 0.1,
8. +                         verboseIter = T
9. +                     ))
10. + Resample01: C=1
11. ....
12. - Resample10: C=1
13. Aggregating results
14. Fitting final model on full training set
15. > print(svm.loocv)
16. Support Vector Machines with Linear Kernel
17.
18. 13611 samples
19.    16 predictor
20.    7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROZ', 'SEKE
    R', 'SIRA'
21.
22. No pre-processing
23. Resampling: Repeated Train/Test Splits Estimated (10 reps, 10%)
24. Summary of sample sizes: 1364, 1364, 1364, 1364, 1364, 1364, ...
25. Resampling results:
26.
27. Accuracy   Kappa
28. 0.9200131  0.9032572
29.
30. Tuning parameter 'C' was held constant at a value of 1
31. > svm.loocv_pred <- predict(svm.loocv , data.valid)
32. > svm.loocv_cM <- table(svm.loocv_pred, data.valid$Class)
33. > svm.loocv_summary <-
    data.frame( Bias = mean(predict(svm.loocv, data.train) != dat
    a.train$Class),
34. +             Variance = mean(svm.loocv_pred != da
    ta.valid$Class),
35. +             Accuracy = sum(diag(svm.loocv_cM)) /
    sum(svm.loocv_cM))
36. > svm.loocv_summary
37.      Bias      Variance      Accuracy
38. 1 0.07148105 0.07492654 0.9250735
```

Decision tree:

```
1. > # Decision Tree
2. > tree.loocv <- train(Class ~., data = data,
3. +                       method = "rpart",
4. +                       trControl = trainControl(
5. +                         method = "LGOCV",
6. +                         number = 10,
7. +                         p = 0.1,
8. +                         verboseIter = T
9. +                       ))
10. + Resample01: cp=0.1647
11. ....
12. - Resample10: cp=0.1647
13. Aggregating results
14. Selecting tuning parameters
15. Fitting cp = 0.165 on full training set
16. > print(tree.loocv)
17. CART
18.
19. 13611 samples
20.    16 predictor
21.    7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROZ', 'SEKE
    R', 'SIRA'
22.
23. No pre-processing
24. Resampling: Repeated Train/Test Splits Estimated (10 reps, 10%)
25. Summary of sample sizes: 1364, 1364, 1364, 1364, 1364, 1364, ...
26. Resampling results across tuning parameters:
27.
28.   cp          Accuracy   Kappa
29. 0.1647293  0.6461746  0.5591477
30. 0.1705912  0.5035111  0.3689403
31. 0.1992052  0.3773577  0.1900813
32.
33. Accuracy was used to select the optimal model using the largest value.
34. The final value used for the model was cp = 0.1647293.
35. > tree.loocv_pred <- predict(tree.loocv, data.valid, type = "raw")
36. > tree.loocv_cM <- table(tree.loocv_pred, data.valid$Class)
37. > tree.loocv_summary <-
    data.frame( Bias = mean(predict(tree.loocv, data.train, type
    = "raw") != data.train$Class),
38. +                               Variance = mean(tree.loocv_pred !=
    factor(data.valid$Class, ordered = FALSE )),
39. +                               Accuracy = sum(diag(tree.loocv_cM))
    / sum(tree.loocv_cM))
40. > tree.loocv_summary
```

```

41.          Bias  Variance  Accuracy
42. 1 0.4682481 0.4608227 0.5391773

```

KNN:

```

1. > # K-Nearest Neighbor
2. > knn.loocv <- train(Class ~., data = data,
3. +                   method = "knn",
4. +                   trControl = trainControl(
5. +                   method = "LGOCV",
6. +                   number = 10,
7. +                   p = 0.1,
8. +                   verboseIter = T
9. +                   ))
10. + Resample01: k=5
11. - Resample01: k=5
12. ....
13. + Resample10: k=9
14. - Resample10: k=9
15. Aggregating results
16. Selecting tuning parameters
17. Fitting k = 9 on full training set
18. > print(knn.loocv)
19. k-Nearest Neighbors
20.
21. 13611 samples
22.    16 predictor
23.    7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROZ', 'SEKE
      R', 'SIRA'
24.
25. No pre-processing
26. Resampling: Repeated Train/Test Splits Estimated (10 reps, 10%)
27. Summary of sample sizes: 1364, 1364, 1364, 1364, 1364, 1364, ...
28. Resampling results across tuning parameters:
29.
30.  k  Accuracy  Kappa
31.  5  0.6160284  0.5334295
32.  7  0.6153099  0.5324363
33.  9  0.6165183  0.5337183
34.
35. Accuracy was used to select the optimal model using the largest value.
36. The final value used for the model was k = 9.
37. > knn.loocv_pred <- predict(knn.loocv , data.valid)
38. > knn.loocv_cM <- table(knn.loocv_pred, data.valid$Class)

```

```

39. > knn.loocv_summary <-
      data.frame( Bias = mean(predict(knn.loocv, data.train) != dat
a.train$Class),
40. +              Variance = mean(knn.loocv_pred != da
ta.valid$Class),
41. +              Accuracy = sum(diag(knn.loocv_cM)) /
      sum(knn.loocv_cM))
42. > knn.loocv_summary
43.      Bias Variance Accuracy
44. 1 0.2196914 0.2140059 0.7859941

```

4. Boosting

Multinational Logistic Regression:

```

1. > # Multinational Logistic Regression #
2. > log.boot <- train(Class ~., data = data,
3. +                  method = "multinom",
4. +                  trControl = trainControl(
5. +                  method = "boot",
6. +                  number = 10,
7. +                  verboseIter = T
8. +                  ))
9. + Resample01: decay=0e+00
10. # weights: 126 (102 variable)
11. initial value 26485.783039
12. iter 10 value 19351.454575
13. ....
14. final value 2780.792544
15. stopped after 100 iterations
16. > print(log.boot)
17. Penalized Multinomial Regression
18.
19. 13611 samples
20. 16 predictor
21. 7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROS', 'SEKE
R', 'SIRA'
22.
23. No pre-processing
24. Resampling: Bootstrapped (10 reps)
25. Summary of sample sizes: 13611, 13611, 13611, 13611, 13611, ...
26. Resampling results across tuning parameters:
27.
28. decay Accuracy Kappa
29. 0e+00 0.9268254 0.9114270

```

```

30. 1e-04 0.9268865 0.9115006
31. 1e-01 0.9224944 0.9061742
32.
33. Accuracy was used to select the optimal model using the largest value.
34. The final value used for the model was decay = 1e-04.
35. > log.boot_pred <- predict(log.boot , data.valid)
36. > log.boot_cM <- table(log.boot_pred, data.valid$Class)
37. > log.boot_summary <-
      data.frame( Bias = mean(predict(log.boot, data.train) != data
      .train$Class),
38. +                               Variance = mean(log.boot_pred != data
      .valid$Class),
39. +                               Accuracy = sum(diag(log.boot_cM)) / s
      um(log.boot_cM))
40. > log.boot_summary
41.      Bias  Variance  Accuracy
42. 1 0.07179595 0.07664055 0.9233595

```

SVM:

```

1. > # Support Vector Machine Classification
2. > svm.boot <- train(Class ~., data = data,
3. +                   method = "svmLinear",
4. +                   trControl = trainControl(
5. +                       method = "boot",
6. +                       number = 10,
7. +                       verboseIter = T
8. +                   ))
9. + Resample01: C=1
10. - Resample01: C=1
11. ....
12. - Resample10: C=1
13. Aggregating results
14. Fitting final model on full training set
15. > print(svm.boot)
16. Support Vector Machines with Linear Kernel
17.
18. 13611 samples
19. 16 predictor
20. 7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROZ', 'SEKE
      R', 'SIRA'
21.
22. No pre-processing
23. Resampling: Bootstrapped (10 reps)
24. Summary of sample sizes: 13611, 13611, 13611, 13611, 13611, 13611, ...

```

```

25. Resampling results:
26.
27.   Accuracy   Kappa
28.   0.9254708  0.9098593
29.
30. Tuning parameter 'C' was held constant at a value of 1
31. > svm.boot_pred <- predict(svm.boot , data.valid)
32. > svm.boot_cM <- table(svm.boot_pred, data.valid$Class)
33. > svm.boot_summary <-
      data.frame( Bias = mean(predict(svm.boot, data.train) != data
      .train$Class),
34. +               Variance = mean(svm.boot_pred != data
      .valid$Class),
35. +               Accuracy = sum(diag(svm.boot_cM)) / s
      um(svm.boot_cM))
36. > svm.boot_summary
37.      Bias   Variance Accuracy
38. 1 0.07148105 0.07492654 0.9250735

```

Decision tree:

```

1. > # Decision Tree
2. > tree.boot <- train(Class ~., data = data,
3. +                   method = "rpart",
4. +                   trControl = trainControl(
5. +                   method = "boot",
6. +                   number = 10,
7. +                   verboseIter = T
8. +                   ))
9. + Resample01: cp=0.1647
10. - Resample01: cp=0.1647
11. .....
12. - Resample10: cp=0.1647
13. Aggregating results
14. Selecting tuning parameters
15. Fitting cp = 0.165 on full training set
16. > print(tree.boot)
17. CART
18.
19. 13611 samples
20.   16 predictor
21.    7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROZ', 'SEKE
      R', 'SIRA'
22.
23. No pre-processing

```

```

24. Resampling: Bootstrapped (10 reps)
25. Summary of sample sizes: 13611, 13611, 13611, 13611, 13611, 13611, ...
26. Resampling results across tuning parameters:
27.
28.   cp          Accuracy   Kappa
29. 0.1647293 0.5558049 0.4369662
30. 0.1705912 0.4681566 0.3200987
31. 0.1992052 0.3311703 0.1178575
32.
33. Accuracy was used to select the optimal model using the largest value.
34. The final value used for the model was cp = 0.1647293.
35. > tree.boot_pred <- predict(tree.boot , data.valid, type = "raw")
36. > tree.boot_cM <- table(tree.boot_pred, data.valid$Class)
37. > tree.boot_summary <-
      data.frame( Bias = mean(predict(tree.boot, data.train, type =
      "raw") != data.train$Class),
38. +               Variance = mean(tree.boot_pred != fa
      ctor(data.valid$Class, ordered = FALSE )),
39. +               Accuracy = sum(diag(tree.boot_cM)) /
      sum(tree.boot_cM))
40. > tree.boot_summary
41.      Bias Variance Accuracy
42. 1 0.4682481 0.4608227 0.5391773

```

KNN:

```

1. > # K-Nearest Neighbor
2. > knn.boot <- train(Class ~., data = data,
3. +                   method = "knn",
4. +                   trControl = trainControl(
5. +                     method = "boot",
6. +                     number = 10,
7. +                     verboseIter = T
8. +                   ))
9. + Resample01: k=5
10. - Resample01: k=5
11. + Resample01: k=7
12. ....
13. + Resample10: k=9
14. - Resample10: k=9
15. Aggregating results
16. Selecting tuning parameters
17. Fitting k = 5 on full training set
18. > print(knn.boot)
19. k-Nearest Neighbors

```

```

20.
21. 13611 samples
22. 16 predictor
23. 7 classes: 'BARBUNYA', 'BOMBAY', 'CALI', 'DERMASON', 'HOROS', 'SEKE
    R', 'SIRA'
24.
25. No pre-processing
26. Resampling: Bootstrapped (10 reps)
27. Summary of sample sizes: 13611, 13611, 13611, 13611, 13611, 13611, ...
28. Resampling results across tuning parameters:
29.
30. k Accuracy Kappa
31. 5 0.6998387 0.6358444
32. 7 0.6993312 0.6350086
33. 9 0.6990990 0.6345673
34.
35. Accuracy was used to select the optimal model using the largest value.
36. The final value used for the model was k = 5.
37. > knn.boot_pred <- predict(knn.boot , data.valid)
38. > knn.boot_cM <- table(knn.boot_pred, data.valid$Class)
39. > knn.boot_summary <-
    data.frame( Bias = mean(predict(knn.boot, data.train) != data
    .train$Class),
40. + Variance = mean(knn.boot_pred != data
    .valid$Class),
41. + Accuracy = sum(diag(knn.boot_cM)) / s
    um(knn.boot_cM))
42. > knn.boot_summary
43. Bias Variance Accuracy
44. 1 0.1792799 0.1755632 0.8244368

```

5. Original classification result

Here is the result that have not use ensemble techniques.

```

1. > data.valid$Class <- factor(data.valid$Class, ordered = T)
2. > log.sum <-
    data.frame( Bias = mean(predict(log.fit, data.train) != data.
    train$Class),
3. + Variance = mean(log.pred != data.valid$Class),
4. + Accuracy = sum(diag(log.cM)) / sum(log.cM))

```



```

5. > svm.sum <-
      data.frame( Bias = mean(predict(svm.fit, data.train) != data.
train$Class),
6. +           Variance = mean(svm.pred != data.valid$Class),
7. +           Accuracy = sum(diag(svm.cM)) / sum(svm.cM))
8. > tree.sum <-
      data.frame( Bias = mean(predict(log.fit, data.train) != data.
train$Class),
9. +           Variance = mean(tree.pred != data.valid$Class),
10. +           Accuracy = sum(diag(tree.cM)) / sum(tree.cM))
11. > knn.sum <- data.frame( Variance = mean(knn.pred != data.valid$Class),
12. +           Accuracy = sum(diag(knn.cM)) / sum(knn.cM))
13. > log.sum
14.      Bias  Variance  Accuracy
15. 1 0.07232077 0.07713026 0.9228697
16. > svm.sum
17.      Bias  Variance  Accuracy
18. 1 0.06623281 0.07590597 0.924094
19. > tree.sum
20.      Bias  Variance  Accuracy
21. 1 0.07232077 0.130999 0.869001
22. > knn.sum
23.      Variance  Accuracy
24. 1 0.3574927 0.6425073

```

6. Comparison

Here we put all the ensemble techniques result together to compare with the classification without ensemble techniques.

In Logistic Regression we can see all three ensemble techniques have improved performance compared to the original one slightly, Bias have been decreased, Variance have been decreased, and accuracy have been increased. It also shows that Logistic Regression is easy to get overfitting, by using ensemble techniques can be very helpful for prevent Logistic Regression get over fitting. All three ensemble techniques performed very close, cross validation did slightly better.

In Support vector machines we get similar result as Logistic Regression, the interesting thing is the bias have been increased, but the variance has been decreased, and accuracy have been increased which is pretty good.

Ensemble techniques have made a better tread off between Bias and Variance. Another interesting thing is all three ensemble techniques gives the same result. This may be because of fold have been set to 10 which is not enough to bring a big difference between them. Or maybe because of this SVM model completion is close to saturation.

For the decision tree, the performance have significantly worse, this may be due to ensemble techniques error or need more optimized calculations and need adjustment to the parameters.

In KNN we can see the ensemble techniques have been very effectively in improving performance in all the aspects. The boosting gives the best result, this is due to in subsequent bags, boosting has chosen those data instances that had been modeled poorly in the overall system before, so that the performance can keep being improved very effectively.

```
1. > # Comparison
2. > log.cv_summary
3.      Bias  Variance  Accuracy
4. 1 0.07179595 0.07639569 0.9236043
5. > log.loocv_summary
6.      Bias  Variance  Accuracy
7. 1 0.07179595 0.07664055 0.9233595
8. > log.boot_summary
9.      Bias  Variance  Accuracy
10. 1 0.07179595 0.07664055 0.9233595
11. > log.sum
12.      Bias  Variance  Accuracy
13. 1 0.07232077 0.07713026 0.9228697

14. > svm.cv_summary
15.      Bias  Variance  Accuracy
16. 1 0.07148105 0.07492654 0.9250735
```

```

17. > svm.loocv_summary
18.      Bias  Variance  Accuracy
19. 1 0.07148105 0.07492654 0.9250735
20. > svm.boot_summary
21.      Bias  Variance  Accuracy
22. 1 0.07148105 0.07492654 0.9250735
23. > svm.sum
24.      Bias  Variance  Accuracy
25. 1 0.06623281 0.07590597 0.924094

26. > tree.cv_summary
27.      Bias  Variance  Accuracy
28. 1 0.4682481 0.4608227 0.5391773
29. > tree.loocv_summary
30.      Bias  Variance  Accuracy
31. 1 0.4682481 0.4608227 0.5391773
32. > tree.boot_summary
33.      Bias  Variance  Accuracy
34. 1 0.4682481 0.4608227 0.5391773
35. > tree.sum
36.      Bias  Variance  Accuracy
37. 1 0.07232077 0.130999 0.869001

38. > knn.cv_summary
39.      Bias  Variance  Accuracy
40. 1 0.1783353 0.1767875 0.8232125
41. > knn.loocv_summary
42.      Bias  Variance  Accuracy
43. 1 0.2196914 0.2140059 0.7859941
44. > knn.boot_summary
45.      Bias  Variance  Accuracy
46. 1 0.1792799 0.1755632 0.8244368
47. > knn.sum
48.      Variance  Accuracy
49. 1 0.3574927 0.6425073

```