# *QSAR Ready Biodegradability Prediction*

## *EECS E6690 Project Presentation*

**Lei Lyu (ll3433)**

**Yang Yu (yy3102)**

**Wenxiang Zhou (wz2542)**
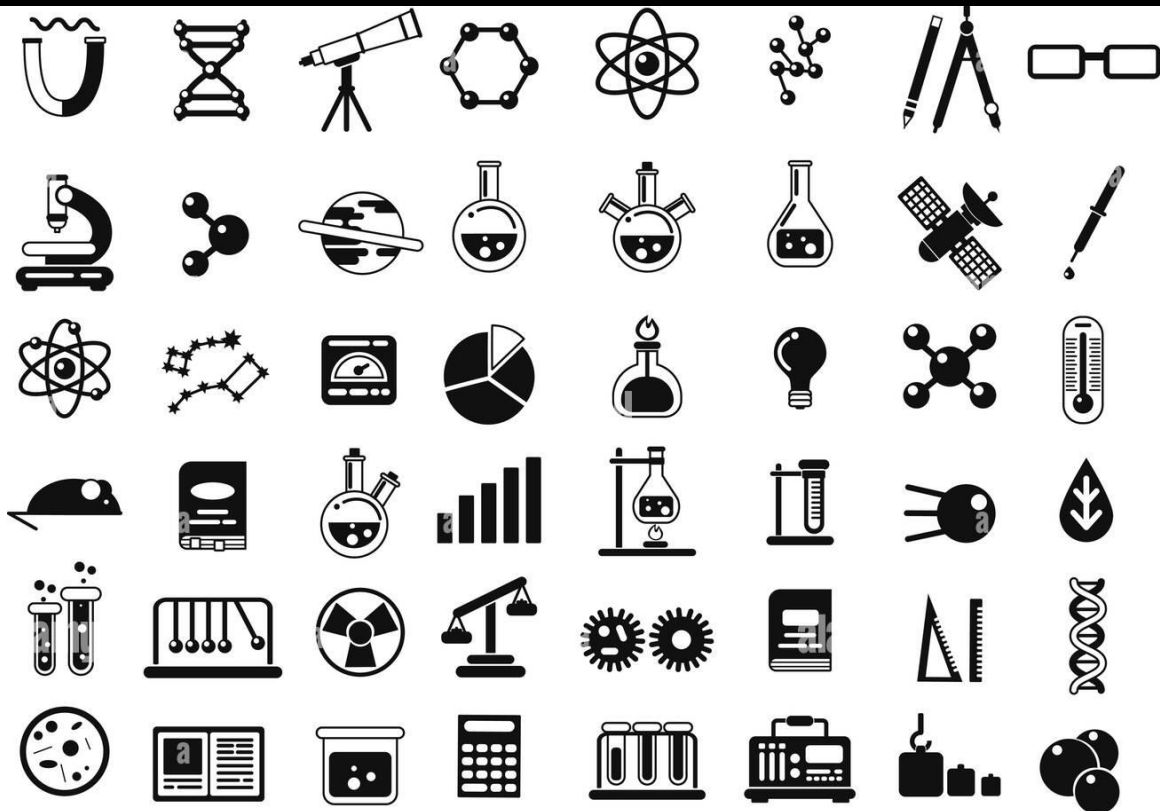
**Yi Chen (yc4029)**

# *Part 1: Dataset description*

*QSAR biodegradation Data Set*

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# Dataset Description

- Various kinds of chemicals are left in the environment by industry

- Many of the chemicals are not biodegradable

- Information about biodegradability of chemicals are not yet abundant



alamy

Image ID: 2BH2FMC
www.alamy.com

COLUMBIA | ENGINEERING
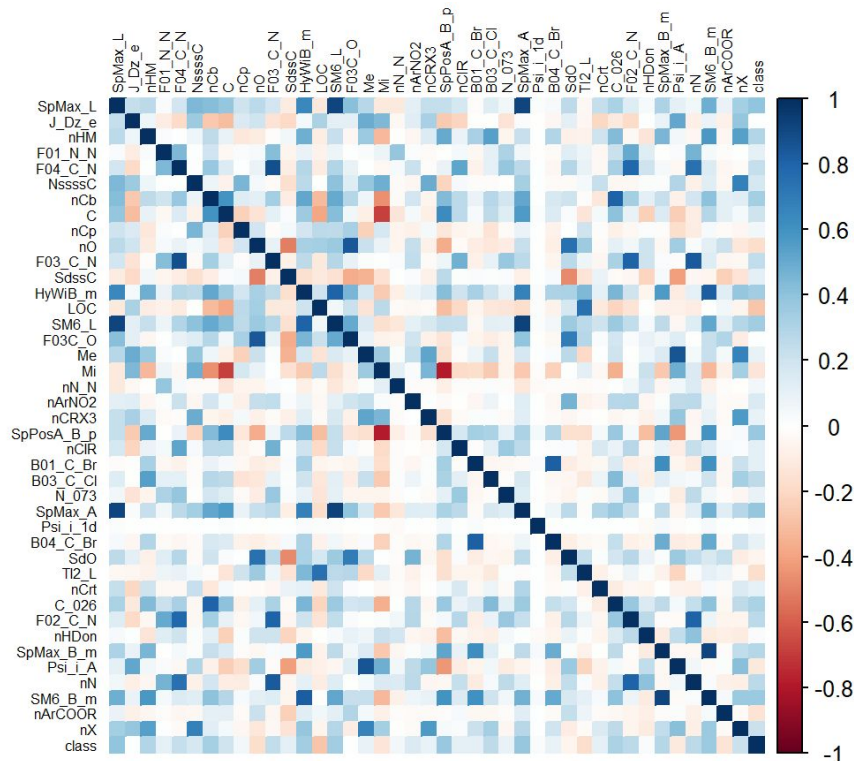The Fu Foundation School of Engineering and Applied Science

# Dataset Description

- QSAR (Quantitative Structure-Activity Relationships) is used to predict the biodegradability of chemicals

- QSAR biodegradation data set was built to develop QSAR models for studying the relationship between chemical structure and biodegradability of molecules

- Experimental values collected from webpage of the National Institute of Technology and Evaluation of Japan (NITE)

COLUMBIA | ENGINEERING
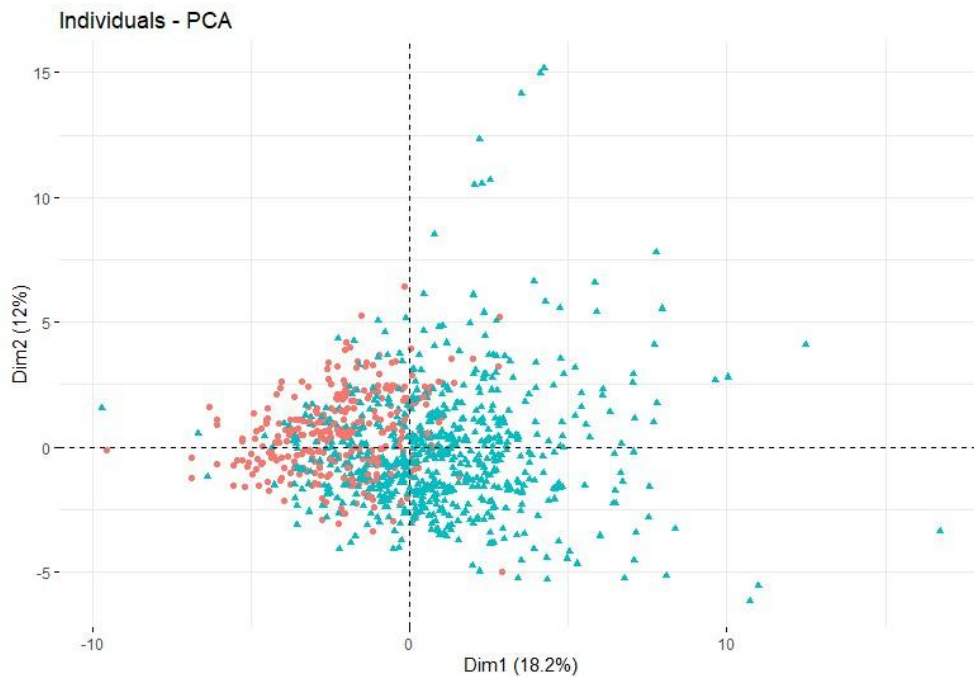The Fu Foundation School of Engineering and Applied Science

- Number of Instances: 1055 where 356 molecules are ready biodegradable (RB) and 699 are not ready biodegradable (NRB)

- Number of Attributes: 41 selected using many classification modeling methods combined with genetic algorithms

- Correlations between descriptors:
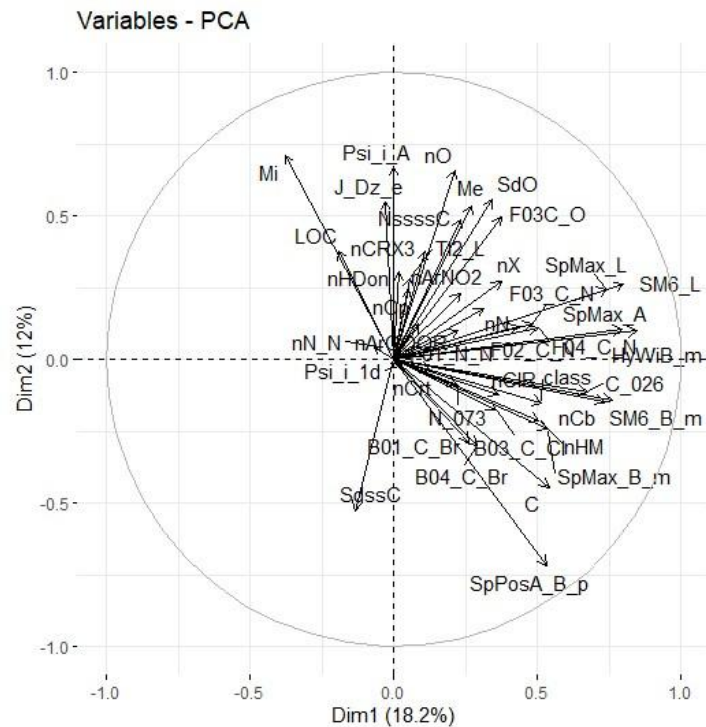
**Heatmap of variable correlations**

Columbia | Engineering
The Fu Foundation School of Engineering and Applied Science

# All Descriptors PCA



Score Plot

Loading Plot

Columbia | ENGINEERING
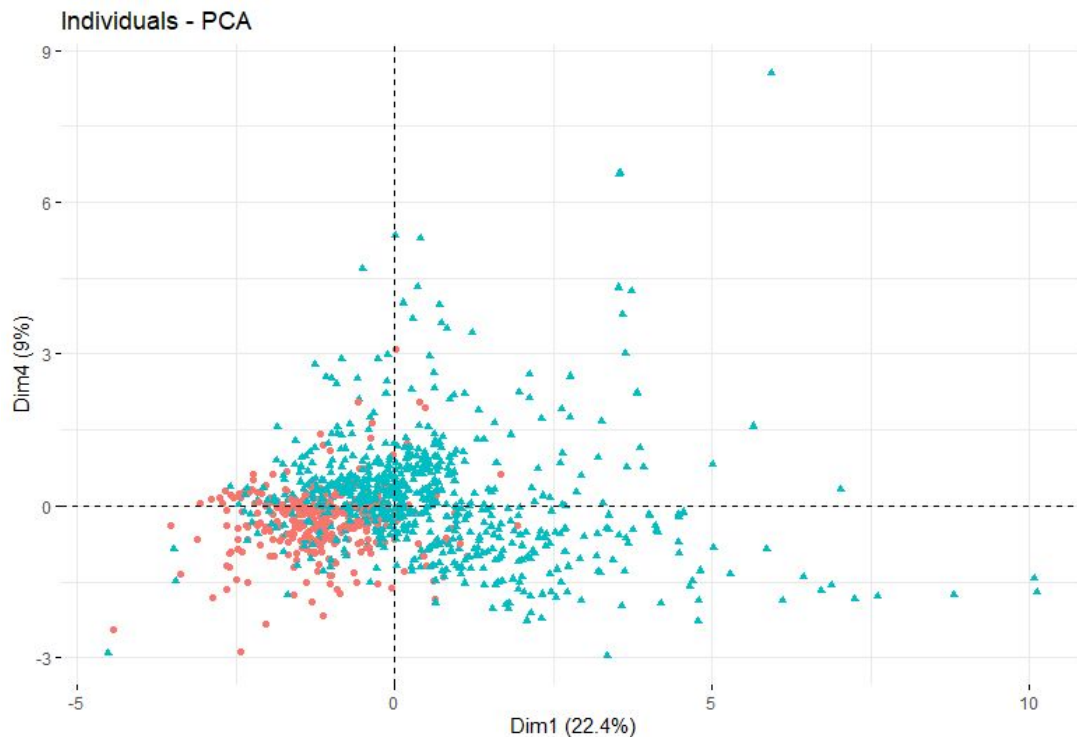The Fu Foundation School of Engineering and Applied Science

## Score Plot

## Loading Plot

# Paper's kNN Descriptors PCA

Score Plot

Loading Plot

Columbia ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# *Part 2: Paper detail and reproduce*

*KNN, PLSDA, SVM*

## Data set dividing:

41 Attributes, 1055 instances
training set: 837     test set: 218
5-fold cross-validation

2 Class:
RB: ready biodegradable
NRB: not ready biodegradable

Columbia | Engineering
The Fu Foundation School of Engineering and Applied Science

# Model Validation

- **Specificity:**
the ability to correctly predict RB molecules

$$Sp = \frac{TN}{TN + FP}$$

TN: # true negatives
FP: # false positives

- **Sensitivity:**
the ability to correctly predict NRB molecules

$$Sn = \frac{TP}{TP + FN}$$

TP: # true positives
FN: # false negatives

- **ER:**
the classification error rate

$$ER = 1 - \frac{Sp + Sn}{2}$$

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# 1. k Nearest Neighbors (kNN)

Our reproduce result:

**CONFUSION MATRIX**



12 descriptors selected

Euclidean distance

| Techniques | ER | Sp | Sn |
|---|---|---|---|
| KNN (paper) | 0.15 | 0.90 | 0.81 |
| KNN (reproduced) | 0.12 | 0.90 | 0.83 |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# 2. Partial Least Squares Discriminant Analysis (PLSDA)

Our reproduce result:

**CONFUSION MATRIX**



**23 descriptors selected**

| Techniques | ER | Sp | Sn |
|---|---|---|---|
| PLSDA (paper) | 0.15 | 0.87 | 0.83 |
| PLSDA (reproduced) | 0.15 | 0.92 | 0.72 |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# 3. support vector machines (SVM)

Our reproduce result:

**CONFUSION MATRIX**

| | Actual | |
|---|---|---|
| | RB | NRB |
| Predicted RB | 52 | 12 |
| Predicted NRB | 20 | 134 |

**DETAILS**

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.722 | 0.918 | 0.812 | 0.722 | 0.765 |
| | **Accuracy** | | **Kappa** | |
| | **0.853** | | **0.659** | |

14 descriptors selected

| Techniques | ER | Sp | Sn |
|---|---|---|---|
| SVM (paper) | 0.14 | 0.91 | 0.82 |
| SVM (reproduced) | 0.15 | 0.92 | 0.72 |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# SVM-improved

Our reproduce result:

all descriptors

5-fold cross-validation ==>7-fold

**CONFUSION MATRIX**

|  | Actual RB | Actual NRB |
|---|---|---|
| Predicted RB | 55 | 10 |
| Predicted NRB | 17 | 136 |

**DETAILS**

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.764 | 0.932 | 0.846 | 0.764 | 0.803 |
|  | Accuracy 0.876 |  | Kappa 0.713 |  |

| Techniques | ER | Sp | Sn |
|---|---|---|---|
| SVM (paper) | 0.14 | 0.91 | 0.82 |
| SVM (reproduced) | 0.15 | 0.92 | 0.72 |
| SVM (improved) | 0.12 | 0.93 | 0.76 |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# *Part 3: Other Techniques Implementation*

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

Methods

- LDA
- Naive Bayes
- Decision Tree
- Bagging
- Random Forest
- Adaboost
- Neural Network

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# LDA(Linear Discriminant Analysis)

**R Code**

```r
1  lda.fit <- lda(class ~ ., data =
   qsar_train)
2  lda.pred <- predict(lda.fit,
   qsar_test)
3  table(lda.class, qsar_test$class)
4  mean(lda.class == qsar_test$class)
```

**CONFUSION MATRIX**

| | Actual | |
|---|---|---|
| | RB | NRB |
| **Predicted** RB | 52 | 12 |
| **Predicted** NRB | 20 | 134 |

**DETAILS**

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.722 | 0.918 | 0.812 | 0.722 | 0.765 |

| | Accuracy | | Kappa | |
|---|---|---|---|---|
| | 0.853 | | 0.659 | |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# Naive Bayes

## R Code

```
1  nb.fit <- naiveBayes(class ~.,data
   = qsar_train)
2  nb.class <- predict(nb.fit,
   qsar_test)
3  table(nb.class, qsar_test$class)
4  mean(nb.class == qsar_test$class)
```
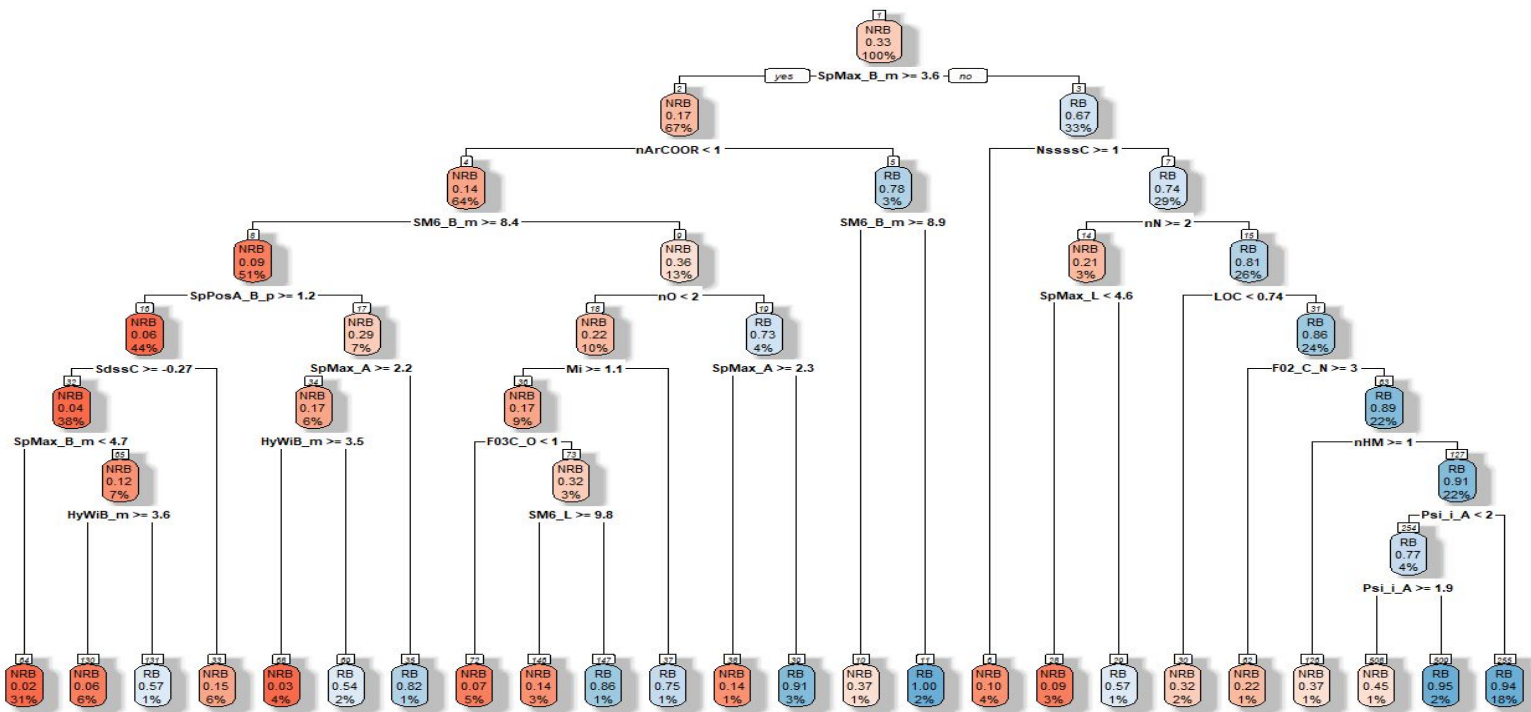
**CONFUSION MATRIX**



**DETAILS**

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.944 | 0.623 | 0.553 | 0.944 | 0.697 |

| | Accuracy | | Kappa | |
|---|---|---|---|---|
| | 0.729 | | 0.481 | |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

```
1  q.rpart <- rpart(class~., data=q_train, method = 'class', cp=0)
```

# Decision Tree(Without pruning)

## CONFUSION MATRIX



**Actual**

|  | RB | NRB |
|---|---|---|
| **Predicted RB** | 114 | 19 |
| **Predicted NRB** | 28 | 57 |

### DETAILS

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.803 | 0.75 | 0.857 | 0.803 | 0.829 |

| Accuracy | Kappa |
|---|---|
| 0.784 | 0.538 |

COLUMBIA ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# Decision Tree(Pruning)

▶ For each $\alpha$:

$$\min_{T \subseteq T_0} \left\{ \sum_{m=1}^{|T|} \sum_{i:\, x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T| \right\},$$
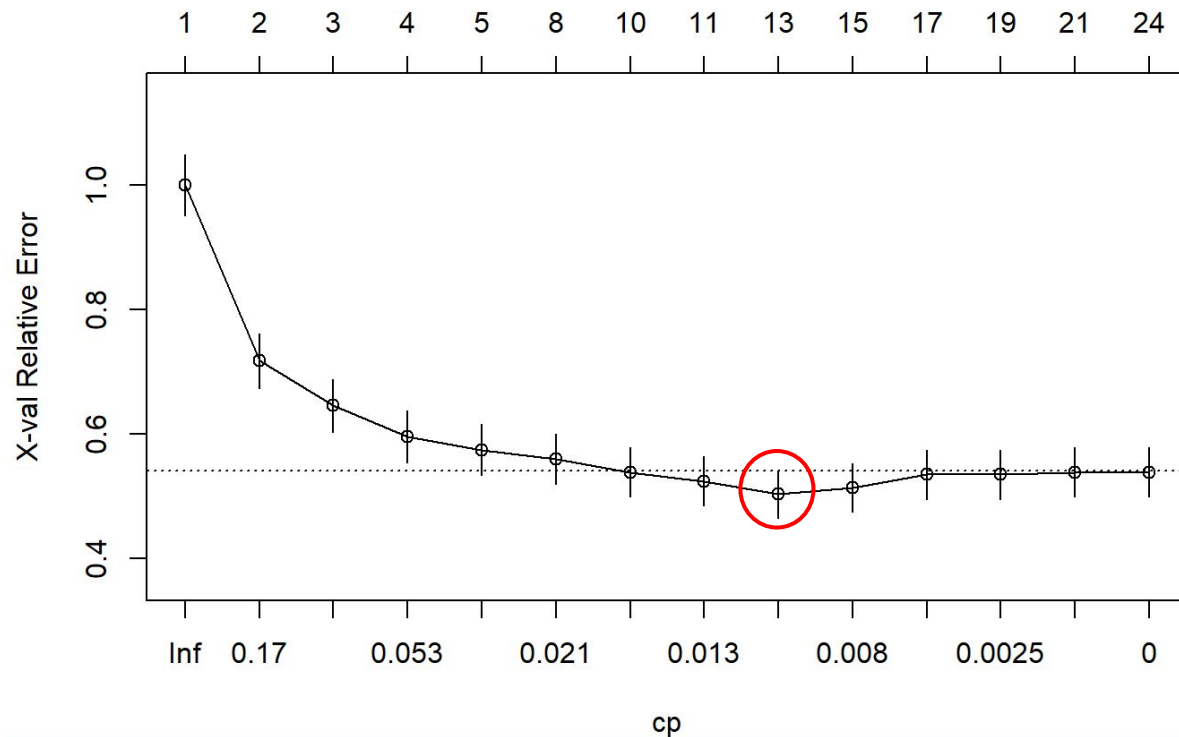
CP(Complexity Parameter) in rpart

$$\sum_{\text{Terminal Nodes}} Misclass_i + \lambda * \left( Splits \right)$$

COLUMBIA | ENGINEERING
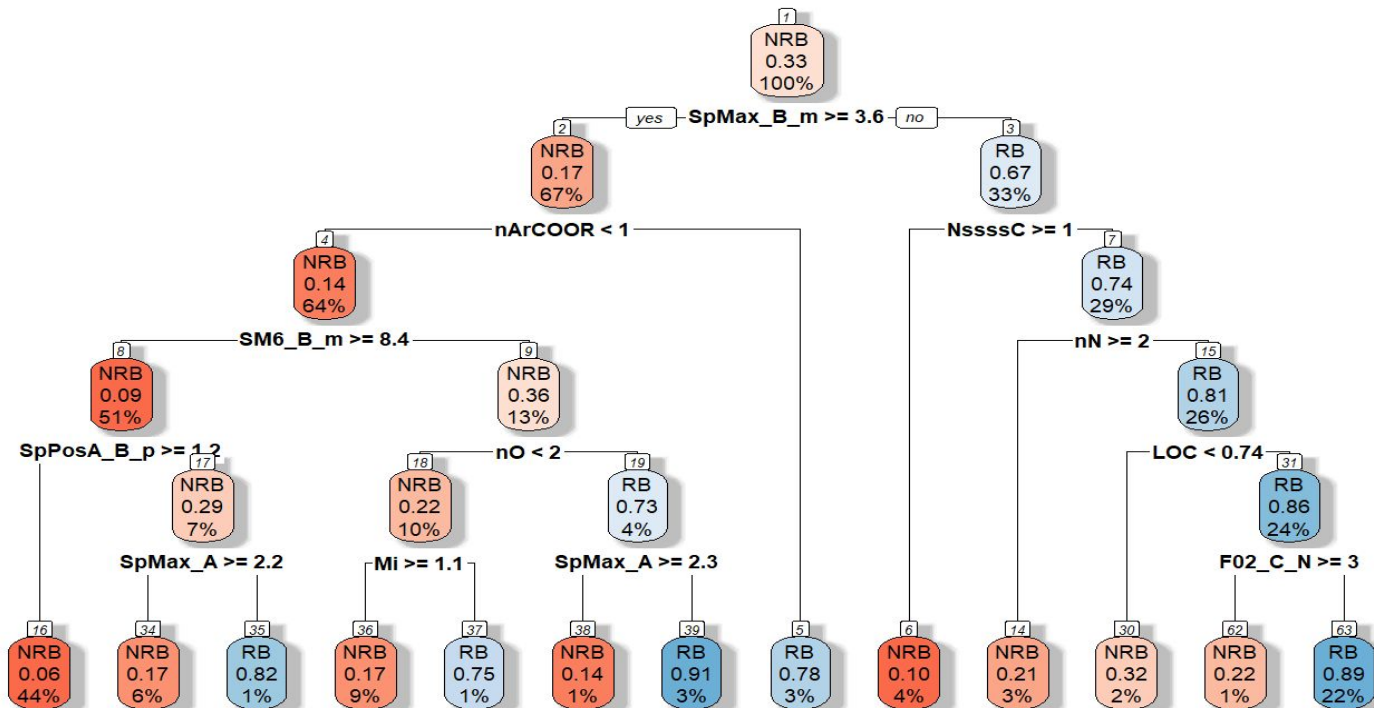The Fu Foundation School of Engineering and Applied Science

# Decision Tree(Pruning)

**R Code**

```
1  printcp(q.rpart)
2  plotcp(q.rpart)
```

```
1   q.prune <- rpart(class~., data=q_train, method = 'class', cp=0.0089286)
```

**CONFUSION MATRIX**

**DETAILS**

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.852 | 0.75 | 0.864 | 0.852 | 0.858 |

| Accuracy | Kappa |
|---|---|
| 0.817 | 0.598 |

Columbia Engineering
The Fu Foundation School of Engineering and Applied Science

# Bagging

## Cross Validation

```
1  num_trees<-c(50, 100, 150, 200, 300, 400)
2  bagging_cv_acc<-seq(0, 0, length = 6)
3  for (epoch in 1:6) {
4    cv_index<-sample(837, 837)
5    for (i in 1:6) {
6      for (k in 1:5) {
7        cv_test_index<-cv_index[(1+round(167.4*(k-
   1))):round(167.4*k)]
8        qsar.bag.cv<-randomForest(as.factor(class)
   ~.,data = qsar_train[-cv_test_index, ],
9                                    ntree =
   num_trees[i], mtry = 41, importance = TRUE)
10       bag.cv.pred <- predict(qsar.bag.cv,
   newdata = qsar_train[cv_test_index, ])
11       cv.acc<-mean(bag.cv.pred ==
   qsar_train[cv_test_index, ]$class)
12       bagging_cv_acc[i]<-bagging_cv_acc[i] +
   cv.acc
13     }
14   }
15 }
16 num_tree = num_trees[which.max(bagging_cv_acc)]
```



**Cross Validation for Bagging**

Columbia Engineering
The Fu Foundation School of Engineering and Applied Science

# Bagging

## R Code

```r
1  qsar.bagging <-
   randomForest(as.factor(class) ~.,data =
   qsar_train, ntree = num_tree, mtry = 41,
   importance = TRUE)
2  bag.pred <- predict(qsar.bagging, newdata =
   qsar_test)
3  mean(bag.pred == qsar_test$class)
```
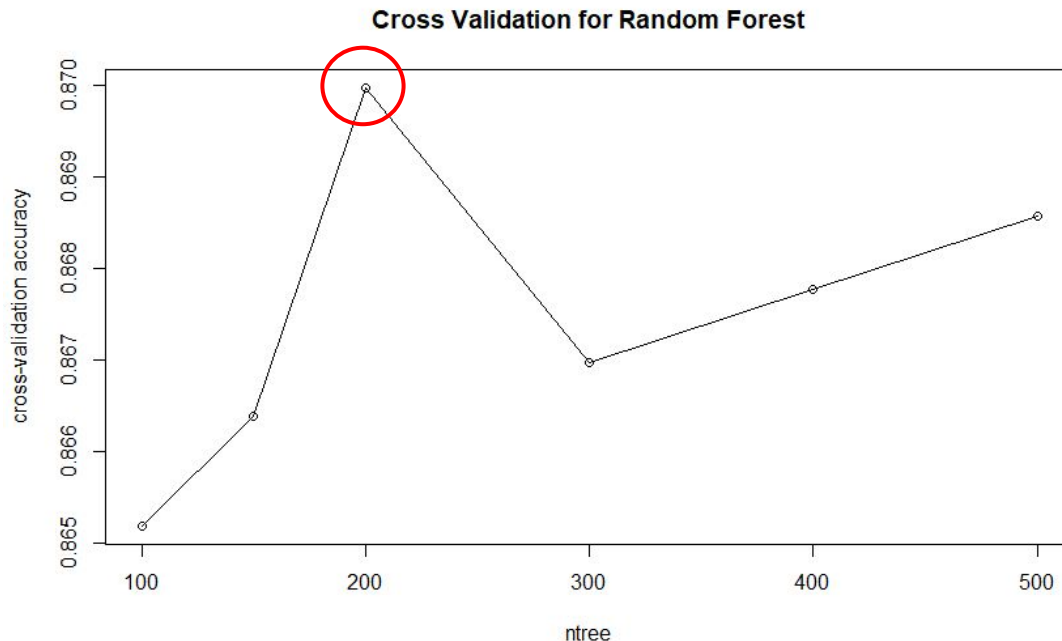
## CONFUSION MATRIX



| | Actual | |
|---|---|---|
| | RB | NRB |
| Predicted RB | 54 | 9 |
| Predicted NRB | 18 | 137 |

## DETAILS

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.75 | 0.938 | 0.857 | 0.75 | 0.8 |
| | **Accuracy** | | **Kappa** | |
| | **0.876** | | **0.711** | |

## Cross Validation

```
1  num_trees_rf<-c(100, 150, 200, 300, 400,
   500)
2  rf_cv_acc<-seq(0, 0, length = 6)
3  for (epoch in 1:6) {
4    cv_index<-sample(837, 837)
5    for (i in 1:6) {
6      for (k in 1:5) {
7        cv_test_index<-
   cv_index[(1+round(167.4*(k-
   1))):round(167.4*k)]
8        qsar.rf.cv<-
   randomForest(as.factor(class) ~.,data =
   qsar_train[-cv_test_index, ],
9                          ntree =
   num_trees_rf[i], importance = TRUE)
10       rf.cv.pred <- predict(qsar.rf.cv,
   newdata = qsar_train[cv_test_index, ])
11       cv.acc<-mean(rf.cv.pred ==
   qsar_train[cv_test_index, ]$class)
12       rf_cv_acc[i]<-rf_cv_acc[i] + cv.acc
13     }
14   }
15 }
16 num_tree_rf =
   num_trees_rf[which.max(rf_cv_acc)]
```



**Cross Validation for Random Forest**

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# Random Forest

## R Code

```r
1  qsar.rf <- randomForest(as.factor(class)
   ~.,data = qsar_train, ntree = num_tree_rf,
   importance = TRUE)
2  rf.pred <- predict(qsar.rf, newdata =
   qsar_test)
3  mean(rf.pred == qsar_test$class)
```
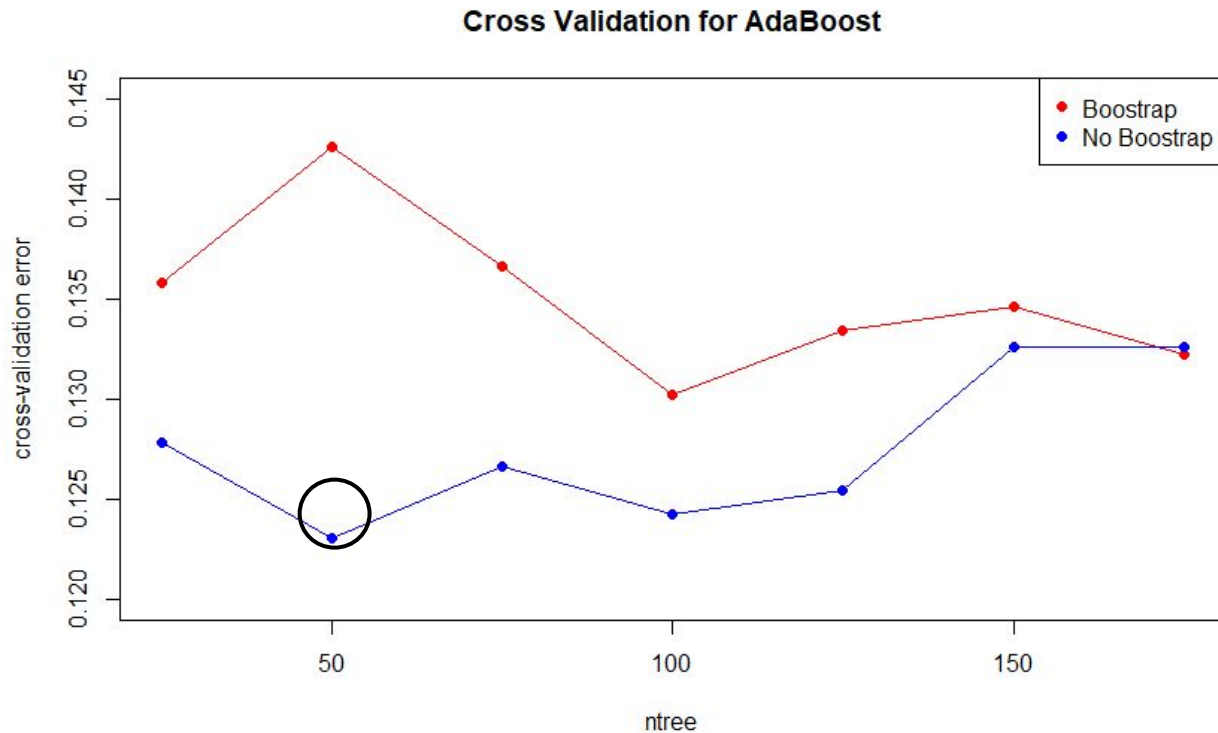
**CONFUSION MATRIX**



**DETAILS**

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.736 | 0.952 | 0.883 | 0.736 | 0.803 |
| | **Accuracy** | | **Kappa** | |
| | **0.881** | | **0.719** | |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

Cross Validation for AdaBoost

## CONFUSION MATRIX



**DETAILS**

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.75 | 0.911 | 0.806 | 0.75 | 0.777 |

| | Accuracy | | Kappa | |
|---|---|---|---|---|
| | 0.858 | | 0.673 | |

Columbia Engineering
The Fu Foundation School of Engineering and Applied Science

# AdaBoost with Breiman coefficient

1. Set $w_i = 1/n, i = 1, 2, \ldots, n$, where $n$ is the number of training points.

2. For $m = 1, \ldots, M$, repeat:

   (a) Fit a (weak) classifier $G_m(x)$ to training data using wights $w_i$.

   (b) Compute the weighted error

   $$e_m = \frac{\sum_{i=1}^{n} w_i \mathbf{1}_{\{y_i \neq G_m(x_i)\}}}{\sum_{i=1}^{n} w_i}$$

   (c) Compute $\alpha_m = \overline{\log((1 - e_m)/e_m)}$.

   (d) Update

   $$w_i \leftarrow w_i \exp(\alpha_m \mathbf{1}_{\{y_i \neq G_m(x_i)\}}), \quad , i = 1, 2, \ldots, n.$$

3. Final classifier

$$G(x) = \mathrm{sign}\left[\sum_{m=1}^{M} \alpha_m G_m(x)\right]$$

$$\alpha_m = \frac{1}{2} \ln \frac{1 - e_m}{e_m}$$

Cross Validation for AdaBoost

## CONFUSION MATRIX

|  | Actual RB | Actual NRB |
|---|---|---|
| Predicted RB | 55 | 11 |
| Predicted NRB | 17 | 135 |

## DETAILS

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.764 | 0.925 | 0.833 | 0.764 | 0.797 |

| Accuracy | Kappa |
|---|---|
| 0.872 | 0.703 |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# Neural Network (NN): 1 hidden layer



Input Layer

↓

Hidden Layer 1: 400 units

↓

ReLu

↓

Dropout Layer

↓

Output Layer: 2 units

↓

Softmax

*QSAR Ready Biodegradability Prediction*

Columbia | Engineering
The Fu Foundation School of Engineering and Applied Science

**CONFUSION MATRIX**

| | Actual | |
|---|---|---|
| | RB | NRB |
| Predicted RB | 64 | 17 |
| Predicted NRB | 8 | 129 |

**DETAILS**

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.889 | 0.884 | 0.79 | 0.889 | 0.837 |

| | Accuracy | | Kappa | |
|---|---|---|---|---|
| | 0.885 | | 0.749 | |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# Deep Neural Network (DNN): 2 hidden layers

# Deep Neural Network (DNN): 2 hidden layers

# *Part 4: Conclusion*

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

D.S. = descriptors selection

| Techniques | Accuracy | Specificity | Sensitivity |
|---|---|---|---|
| SVM+D.S. | 85.3% | 91.8% | 72.2% |
| KNN+D.S. | 87.6% | 89.7% | 83.3% |
| PLSDA+D.S. | 85.3% | 91.8% | 72.2% |
| LDA | 85.3% | 91.8% | 72.2% |
| Naive Bayes | 72.9% | 62.3% | 94.4% |
| Tree | 78.4% | 75% | 80.3% |
| Pruning Tree | 81.7% | 75% | 85.2% |
| SVM | 87.6% | 93.2% | 76.4% |

| Techniques | Accuracy | Specificity | Sensitivity |
|---|---|---|---|
| Bagging | 87.6% | 93.8% | 75.0% |
| Random Forest | 88.1% | 95.2% | 73.6% |
| Adaboost (with Breiman coef.) | 87.2% | 92.5% | 76.4% |
| Adaboost (with Freud coef.) | 85.8% | 91.1% | 75.0% |
| Neural Network (NN) | 88.5% | 88.4% | 88.9% |
| DNN | 87.6% | 92.5% | 77.8% |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

**CONFUSION MATRIX**

**DETAILS**

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.778 | 0.925 | 0.836 | 0.778 | 0.806 |

| | Accuracy | | Kappa | |
|---|---|---|---|---|
| | 0.876 | | 0.715 | |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science

# Consensus Model

| Consensus Model | Accuracy | Specificity | Sensitivity |
|---|---|---|---|
| C1: (SVM+D.S.) + (KNN+D.S.) + (PLSDA+D.S.) | 87.61% | 93.15% | 76.39% |
| C2: NN + SVM + (KNN+D.S.) | 88.99% | 92.47% | 81.94% |
| C3: NN + DNN + (KNN+D.S.) | 89.91% | 92.47% | 84.72% |
| C4: NN + Adaboost (with Breiman coef.) + (KNN+D.S.) | 90.37% | 93.15% | 84.72% |
| C5: NN + Random Forest + (KNN+D.S.) | 89.91% | 93.15% | 83.33% |
| C6: NN + Bagging + (KNN+D.S.) | 88.99% | 91.78% | 83.33% |
| C7: NN + Adaboost (with Freund coef.) + (KNN+D.S.) | 90.37% | 93.84% | 83.33% |
| C8: NN * 2 + Adaboost (with Freud coef.) + Adaboost (with Breiman coef.) + (KNN+D.S.) | 90.83% | 93.84% | 84.72% |

COLUMBIA | ENGINEERING
The Fu Foundation School of Engineering and Applied Science