

Introduction to Data Science

DSA1101

Semester 1, 2018/2019

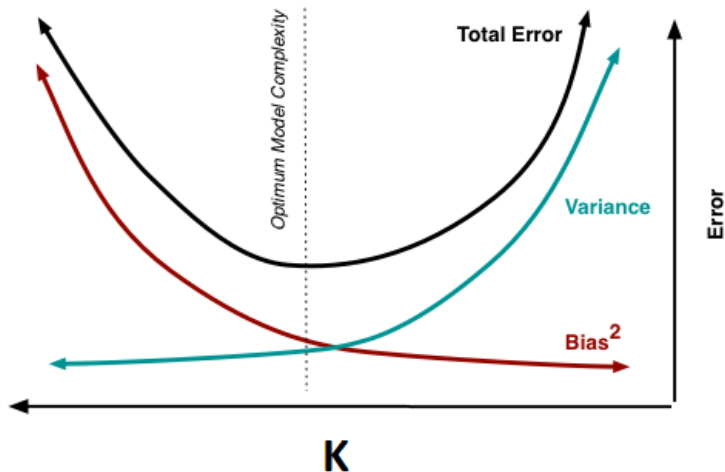
Week 6

***n*-fold cross validation in R**

n -fold cross validation in R

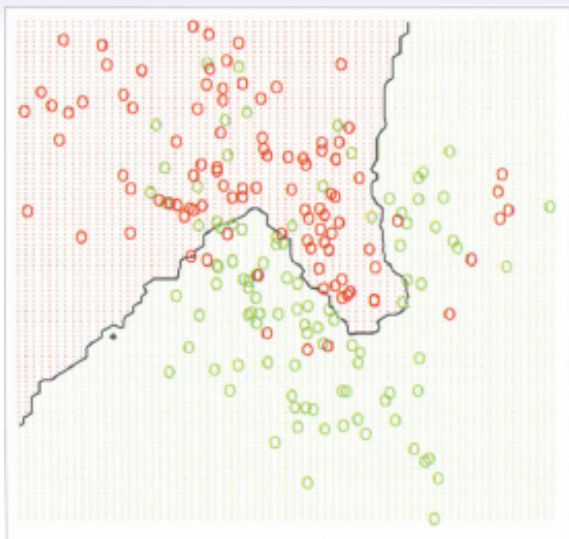
- This week, we will continue to explore n -fold cross validation in R
- We will also examine the R code in more detail

n -fold cross validation in R



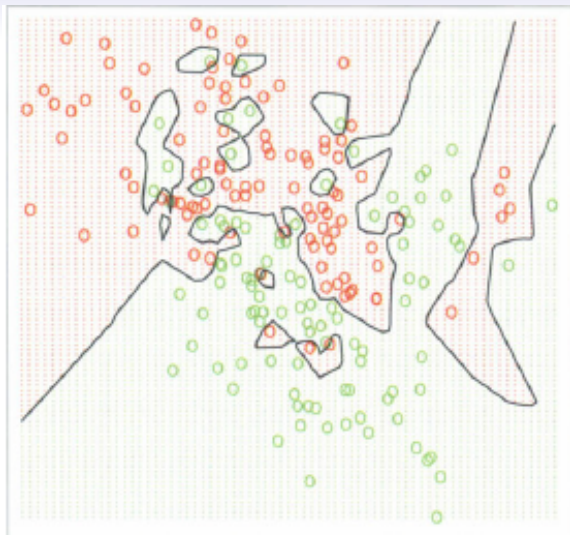
Bias-variance tradeoff. Source: <http://scott.fortmann-roe.com>

n -fold cross validation in R



Prediction by majority vote with 15 nearest neighbors. Source: *The Elements of Statistical Learning*, Hastie et al.

n -fold cross validation in R

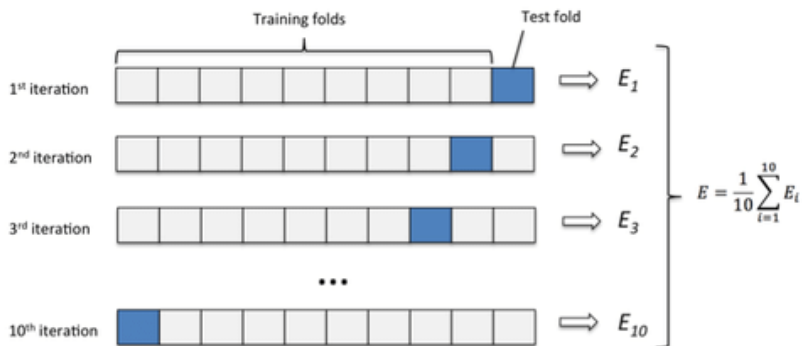


Prediction by majority vote with one nearest neighbor. Source: *The Elements of Statistical Learning*, Hastie et al.

n -fold cross validation in R

- The entire dataset is randomly split into N datasets of approximately equal size.
- $N-1$ of these datasets are treated as the training dataset, while the remaining one is the test dataset. A measure of the model error is obtained.
- This process is repeated across the various combinations of N datasets taken $N - 1$ at a time.
- The observed N models errors are averaged across the N folds.

Diagnostics of Classifiers



n -fold cross validation in R

- We will illustrate how to implement n -fold cross-validation in R to evaluate the performance of the k -nearest neighbor classifier
- In particular, we will attempt to estimate the optimal value of k (or the optimal model complexity) that will give the best classification performance

n -fold cross validation in R

- We will illustrate n -fold cross validation using the famous Iris Flower Dataset which was first introduced in 1936 by the famous statistician Ronald Fisher and consists of 50 observations from each of three species of Iris (*Iris setosa*, *Iris virginica* and *Iris versicolor*).
- Four features (x) were measured from each sample: the length and the width of the sepals and petals.

n -fold cross validation in R



Source: <http://suruchifialoke.com>

n -fold cross validation in R

- The dataset has been posted to IVLE as 'iris.csv'
- Also available at
<https://archive.ics.uci.edu/ml/datasets/Iris>
- 150 data points with the following attribute information:

- (1) sepal length in cm
- (2) sepal width in cm
- (3) petal length in cm
- (4) petal width in cm
- (5) class: *Iris Setosa*, *Iris Versicolour* or *Iris Virginica*

n -fold cross validation in R

- Read in the dataset. Note that variable names are not found in the CSV file

```
1 > iris= read.csv("iris.csv",header=FALSE)
2 > head(iris)
3      V1  V2  V3  V4      V5
4 1  5.1  3.5  1.4  0.2 Iris-setosa
5 2  4.9  3.0  1.4  0.2 Iris-setosa
6 3  4.7  3.2  1.3  0.2 Iris-setosa
7 4  4.6  3.1  1.5  0.2 Iris-setosa
8 5  5.0  3.6  1.4  0.2 Iris-setosa
9 6  5.4  3.9  1.7  0.4 Iris-setosa
```

n -fold cross validation in R

- We can rename the attributes more meaningfully

```
1 > names(iris)
2 [1] "V1" "V2" "V3" "V4" "V5"
3 > names(iris)= c("X1","X2","X3","X4","Y")
4 > names(iris)
5 [1] "X1" "X2" "X3" "X4" "Y"
6 > head(iris)
7      X1  X2  X3  X4      Y
8 1  5.1  3.5  1.4  0.2 Iris-setosa
9 2  4.9  3.0  1.4  0.2 Iris-setosa
10 3  4.7  3.2  1.3  0.2 Iris-setosa
11 4  4.6  3.1  1.5  0.2 Iris-setosa
12 5  5.0  3.6  1.4  0.2 Iris-setosa
13 6  5.4  3.9  1.7  0.4 Iris-setosa
```

n -fold cross validation in R

- Using the same dataset as both the training and testing data can give misleading results
- For example, when $k = 1$, we use each single data point to predict itself:

```
1 > set.seed(1)
2 > library(class)
3 > X=iris[,c("X1", "X2", "X3", "X4")]
4 > Y=iris[,c("Y")]
5 > pred <- knn(train=X, test=X, cl=Y, k=1)
6 > mean(Y != pred)
7 [1] 0
8 > mean(Y == pred)
9 [1] 1
```

- Therefore there is a need for more principled methods to evaluate classifier performance, e.g. n -fold cross validation

n -fold cross validation in R

- Note that the code `mean(Y!=pred)` computes the proportion of data points which label is predicted incorrectly

- A small example:

```
1 > Predicted= c(1,1,0,0,1,1)
2 > Actual= c(1,1,1,0,0,1)
3 > Predicted != Actual
4 [1] FALSE FALSE TRUE FALSE TRUE FALSE
5 > mean(Predicted != Actual)
6 [1] 0.3333333
7 > Predicted == Actual
8 [1] TRUE TRUE FALSE TRUE FALSE TRUE
9 > mean(Predicted == Actual)
10 [1] 0.6666667
```

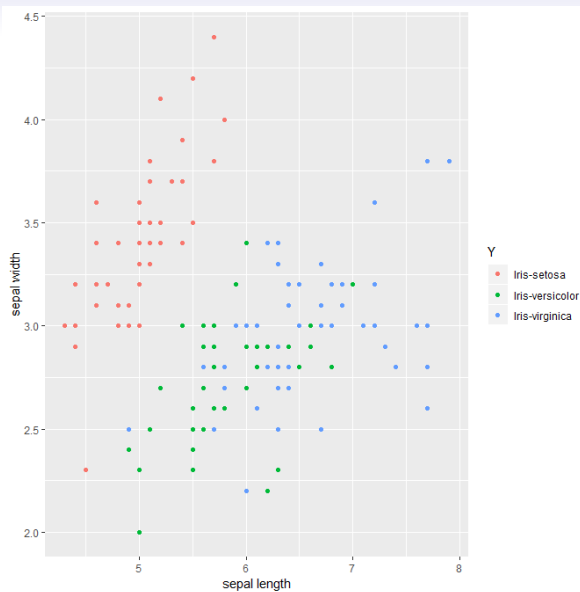
- Similarly, the code `mean(Y==pred)` computes the proportion of data points which label is predicted correctly

n -fold cross validation in R

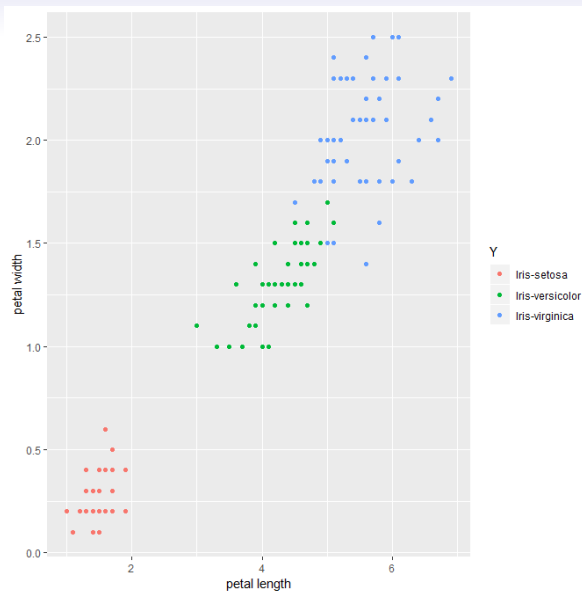
- Before running n -fold cross validation, we can try to visualize the label and feature patterns

```
1 library(ggplot2)
2 library(magrittr)
3 # sepal width vs. sepal length
4 ggplot(aes(x=X[,1], y=X[,2], color=Y)) +
5   geom_point()+
6   labs(x = "sepal length")+   labs(y = "sepal width"
7     )
8 # petal width vs. petal length
9 ggplot(aes(x=X[,3], y=X[,4], color=Y)) +
10  geom_point()+
11  labs(x = "petal length")+   labs(y = "petal width"
12    )
```

n -fold cross validation in R

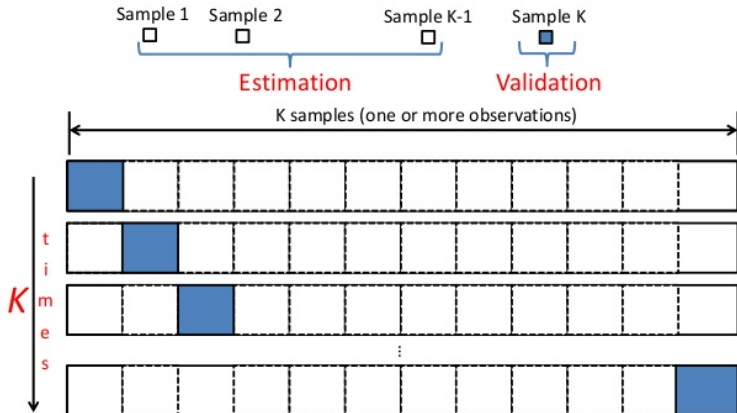


n -fold cross validation in R



Cross-validation: How it works?

- K-fold cross-validation:



n -fold cross validation in R

- We will perform 10-fold cross validation; first randomly split the 150 data points into 10 sets:

```
1 > n_folds=10
2 > folds_i <- sample(rep(1:n_folds, length.out =
   150))
3 > table(folds_i)
4 folds_i
5  1  2  3  4  5  6  7  8  9 10
6 15 15 15 15 15 15 15 15 15 15
```

n -fold cross validation in R

- We illustrate how the 150 data points is randomly divided into 10 sets using the previous code
- A simple example involves 3-fold, 12 data points

```
1 > n_folds=3
2 > Number_of_datapoints=12
3 > rep(1:n_folds,length.out = Number_of_datapoints)
4 [1] 1 2 3 1 2 3 1 2 3 1 2 3
5 > n_folds=3
6 > Number_of_datapoints=12
7 > index=rep(1:n_folds,length.out = Number_of_
8           datapoints)
9 > sample(index)
[1] 3 2 2 1 1 1 2 3 1 3 2 3
```

n -fold cross validation in R

- In this case, `sample(index)` generates a random permutation of the elements of `index`

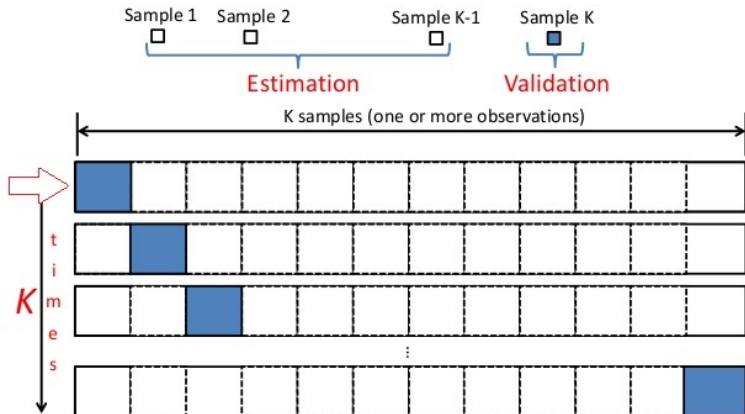
```
1 > n_folds=3
2 > Number_of_datapoints=12
3 > rep(1:n_folds,length.out = Number_of_datapoints)
4 [1] 1 2 3 1 2 3 1 2 3 1 2 3
5 > n_folds=3
6 > Number_of_datapoints=12
7 > index=rep(1:n_folds,length.out = Number_of_
      datapoints)
8 > sample(index)
9 [1] 3 2 2 1 1 1 2 3 1 3 2 3
```

n -fold cross validation in R

- Now, we are ready to run n -fold cross validation in R
- We will first set $k = 1$ for the k -nearest neighbors classifier

Cross-validation: How it works?

- K-fold cross-validation:



n -fold cross validation in R

- We start with $k = 1$, and observe first iteration: using the first dataset as our test data and the remaining 9 datasets as training data

```
1 > test_j <- which(folds_j == 1)
2 > pred <- knn(train=X[ -test_j, ], test=X[test_j,
   ], cl=Y[-test_j ], k=1)
3 > mean(Y[test_j] != pred)
4 [1] 0.06666667
5 > mean(Y[test_j] == pred)
6 [1] 0.9333333
```

n -fold cross validation in R

- Note that the command `test_j<-which(folds_j==1)` generates the indices corresponding to data points in the first data set:

```
1 > folds_j[1:20]
2 [1]  4  9  8  4  7  7  1 10  1  1  6  3
3 [13]  8  1  7  9  5  7  4  3
4 > which(folds_j == 1)
5 [1]  7  9 10 14 22 23 26 37 43
6 [10] 98 107 118 130 132 146
```

n-fold cross validation in R

- The command `X[-test_j,]` removes the data points from first data set from the feature dataframe, keeping the data points from the remaining nine data sets
- Similarly, the command `Y[-test_j]` removes the data points from first data set from the label vector, keeping the data points from the remaining nine data sets
- The command `X[test_j,]` keeps the data points only from first data set from the feature dataframe

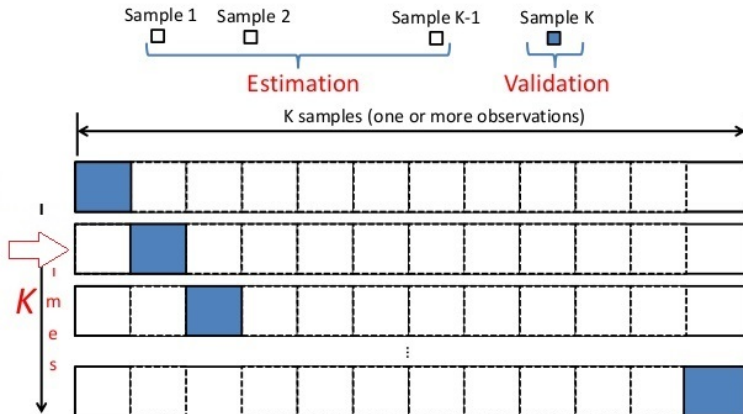
n -fold cross validation in R

- Therefore, the following command takes only the first data set as the test data, and the remaining $n-1$ data sets as the training data

```
1 > pred <- knn(train=X[ -test_j, ], test=X[test_j,  
    ], cl=Y[-test_j ], k=1)
```

Cross-validation: How it works?

- K-fold cross-validation:



n -fold cross validation in R

- Now we perform the second iteration: using the second dataset as our test data and the remaining 9 datasets as training data

```
1 > ## second fold of the iteration
2 > test_j <- which(folds_j == 2)
3 > pred <- knn(train=X[ -test_j, ], test=X[test_j,
      ], cl=Y[-test_j ], k=1)
4 > mean(Y[test_j] != pred)
5 [1] 0.06666667
6 > mean(Y[test_j] == pred)
7 [1] 0.9333333
```

n -fold cross validation in R

- We need to run the same code for $j = 3, 4, 5, \dots, 10$.
- The for loop in R allows us to write compact code for this

```
1 > for (j in 1:10) {  
2 +   print(j)  
3 + }  
4 [1] 1  
5 [1] 2  
6 [1] 3  
7 [1] 4  
8 [1] 5  
9 [1] 6  
10 [1] 7  
11 [1] 8  
12 [1] 9
```


n -fold cross validation in R

- Putting it altogether in a for loop in R:

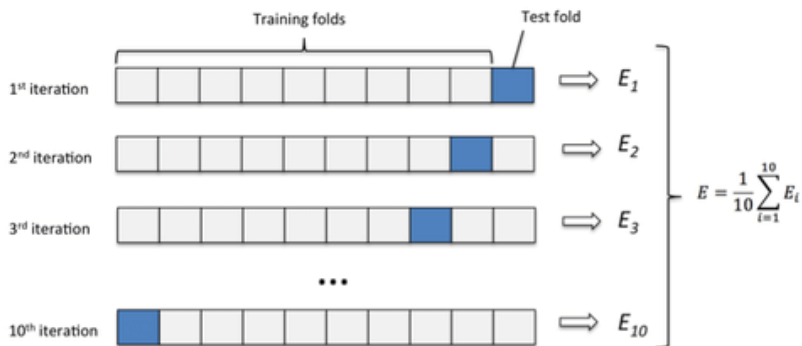
```
1 > err=numeric(10)
2 > acc=numeric(10)
3 >
4 > for (j in 1:10) {
5 +
6 + test_j <- which(folds_j == j)
7 + pred <- knn(train=X[ -test_j, ], test=X[test_j,
8 +           ], cl=Y[-test_j ], k=1)
9 +
10 + err[j]=mean(Y[test_j] != pred)
11 + acc[j]=mean(Y[test_j] == pred)
12 + }
13 > err
14 [1] 0.06666667 0.06666667 0.06666667
15 [4] 0.06666667 0.00000000 0.06666667
16 [7] 0.06666667 0.00000000 0.00000000
17 [10] 0.00000000
```

n -fold cross validation in R

- Here, `err=numeric(10)` declares a numeric vector of length 10
- Therefore, the error rate from each iteration of the 10-fold cross-validation can be saved to the vector `err`
- Similarly, the accuracy from each iteration of the 10-fold cross-validation can be saved to the vector `acc`

```
1 > test=numeric(10)
2 > test
3 [1] 0 0 0 0 0 0 0 0 0 0
```

Diagnostics of Classifiers



n -fold cross validation in R

- For the k -nearest neighbor classifier with $k = 1$, the final estimated error rate is the average of the ten error rates across the 10-fold iterations.
- Similarly, the final estimated accuracy is the average of the ten accuracies across the 10-fold iterations.

```
1 > error=mean(err)
2 > accur=mean(acc)
3 > error
4 [1] 0.04
5 > accur
6 [1] 0.96
```

n -fold cross validation in R

- However, to plot the accuracy (or error rate) against different values of k , we need to repeat the 10-fold cross validation procedure at different values of k
- This suggest another for loop structure, indexed by k
- In this example, we will perform 10-fold cross validation for $k = 1, 2, \dots, 50$.

n-fold cross validation in R

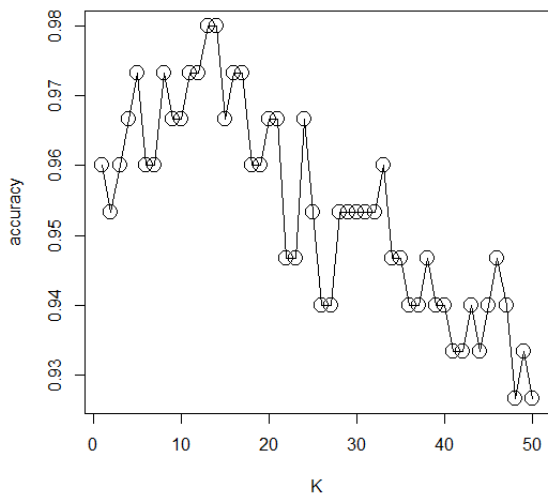
```
1 error_differentK=numeric(50)
2 accur_differentK=numeric(50)
3
4 for (K in 1:50) {
5
6   error=numeric(10)
7   accur=numeric(10)
8   for (j in 1:10) {
9     test_j <- which(folds_j == j)
10    pred <- knn(train=X[ -test_j, ], test=X[test_j,
11      ], cl=Y[-test_j ], k=K)
12    error[j]=mean(Y[test_j] != pred)
13    accur[j]=mean(Y[test_j] == pred)
14  }
15  error_differentK[K]=mean(error)
16  accur_differentK[K]=mean(accur)
17 }
```

n -fold cross validation in R

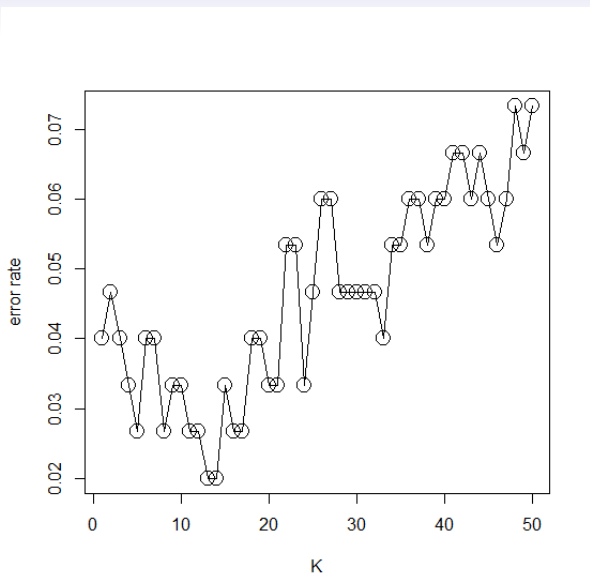
- Then, we can plot the accuracy (or error rate) against different values of k , to determine the optimal k .

```
1 plot(1:50, accur_differentK, type="o", ylab="
    accuracy", xlab="K", cex.axis=1, cex=2)
2 plot(1:50, error_differentK, type="o", ylab="error
    rate", xlab="K", cex.axis=1, cex=2)
```

n -fold cross validation in R



n -fold cross validation in R



n -fold cross validation in R

- Recall the confusion matrix

		Predicted Class	
		Positive	Negative
	Actual Class	True Positives (TP)	False Negatives (FN)
	Positive		
	Negative	False Positives (FP)	True Negatives (TN)

n -fold cross validation in R

- Recall the confusion matrix

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positives (TP)	False Negatives (FN)
	Negative	False Positives (FP)	True Negatives (TN)

- $\text{accuracy} = \frac{TP+TN}{TP+FN+FP+TN}$
- $\text{error rate} = \frac{TP+TN}{TP+FN+FP+TN}$
- Therefore $\text{accuracy} = 1 - \text{error rate}$

n -fold cross validation in R



- To illustrate the n -fold cross validation procedure, we have written our own code in R to implement it
- There are many data analytics packages in R that offers ready-made solutions which make our jobs much easier
- We will use one such package, 'caret', as an example

```
1 library(caret)
```

n -fold cross validation in R

- We will perform 10-fold cross-validation for the k -nearest neighbor classifier:

```
1 fitControl <- trainControl(## 10-fold CV
2                             method = "repeatedcv",
3                             number = 10,
4                             ## repeated ten times
5                             repeats = 10)
6 set.seed(2)
7 knnFit1 <- train(Y, data = iris,
8                  method = "knn",
9                  trControl = fitControl,
10                  preprocess = c("center", "scale")
11                  ,
12                  tuneLength=50)
```

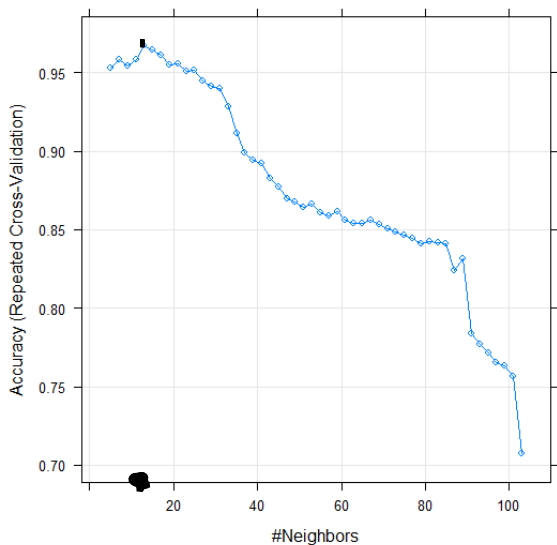


n -fold cross validation in R

- We will perform 10-fold cross-validation for the k -nearest neighbor classifier:

```
1 > knnFit1
2 k-Nearest Neighbors
3 ...
4 Resampling results across tuning parameters:
5
6   k      Accuracy      Kappa
7   5  0.9526667  0.929
8   7  0.9586667  0.938
9   9  0.9540000  0.931
10  ...
11 Accuracy was used to select the optimal model
    using the largest value.
12 The final value used for the model was k = 13.
13
14 plot(knnFit1)
```

n -fold cross validation in R



n -fold cross validation in R: Wine recognition example



- Wines were grown in the same region in Italy but derived from 3 different cultivars.
- The task is to predict wine origin based on 13 attributes having continuous values

n -fold cross validation in R: Wine recognition example

- The 13 features (X) of the dataset are:

- 1 Alcohol
- 2 Malic acid
- 3 Ash
- 4 Alkalinity of ash
- 5 Magnesium
- 6 Total phenols
- 7 Flavanoids
- 8 Nonflavonoids phenols
- 9 Proanthocyanins
- 10 Color intensity
- 11 Hue
- 12 OD280/OD315 of diluted wines
- 13 Proline

n -fold cross validation in R: Wine recognition example

- The data set is available from <https://archive.ics.uci.edu/ml/machine-learning-databases/wine/wine.data>
- Task is to predict label Y of origin: 1,2 or 3

```
1 dataurl <- "https://archive.ics.uci.edu/ml/machine-  
  learning-databases/wine/wine.data"  
2 download.file(url = dataurl, destfile = "wine.data"  
  )  
3 wine_df <- read.csv("wine.data", header = FALSE)
```

n -fold cross validation in R: Wine recognition example

```
1 > head(wine_df)
2   V1    V2    V3    V4    V5    V6    V7    V8
3 1  1 14.23 1.71 2.43 15.6 127 2.80 3.06
4 2  1 13.20 1.78 2.14 11.2 100 2.65 2.76
5 3  1 13.16 2.36 2.67 18.6 101 2.80 3.24
6 4  1 14.37 1.95 2.50 16.8 113 3.85 3.49
7 5  1 13.24 2.59 2.87 21.0 118 2.80 2.69
8 6  1 14.20 1.76 2.45 15.2 112 3.27 3.39
9      V9   V10   V11   V12   V13   V14
10 1 0.28 2.29 5.64 1.04 3.92 1065
11 2 0.26 1.28 4.38 1.05 3.40 1050
12 3 0.30 2.81 5.68 1.03 3.17 1185
13 4 0.24 2.18 7.80 0.86 3.45 1480
14 5 0.39 1.82 4.32 1.04 2.93  735
15 6 0.34 1.97 6.75 1.05 2.85 1450
```

- First column 'V1' contains the wine origin labels

n -fold cross validation in R: Wine recognition example

- We will perform 10-fold cross validation for the k -nearest neighbors classifier on the wine data set:

```
1 wine_df$V1=factor(wine_df$V1)
2 trctrl <- trainControl(method = "repeatedcv",
  number = 10, repeats = 10)
3 set.seed(3)
4 knn_fit <- train(V1 ~., data = wine_df, method = "
  knn",
5 trControl=trctrl,
6 preProcess = c("center", "scale")
7 tuneLength = 20)
8 plot(knn_fit)
```

n -fold cross validation in R: Wine recognition example

- We will perform 10-fold cross validation for the k -nearest neighbors classifier on the wine data set:

```
1 > knn_fit
2 k-Nearest Neighbors
3
4 178 samples
5 13 predictor
6 3 classes: '1', '2', '3'
7 ...
8   k   Accuracy   Kappa
9   5  0.9718920  0.9577857
10  7  0.9689800  0.9534599
11  9  0.9713657  0.9570151
12 11  0.9657086  0.9486180
13 ...
14 Accuracy was used to select the optimal model
    using
15 the largest value.
16 The final value used for the model was k = 35.
```

n -fold cross validation in R

