

Statistisches Data Mining (StDM)

Praktikum Woche 13

Aufgabe 1 Boosting

In dieser Aufgabe werden verschiedenen Klassifikationsmethoden auf simulierte Daten angewandt und die Performance verglichen.

- a) Simulieren Sie 12'000 mal 10 Standardnormalverteilte Beobachtungen X_1, X_2, \dots, X_{10} . Generieren Sie dann die Variable Y_i aufgrund der simulierten Daten folgendermassen:

$$Y_i = \begin{cases} 0 & \text{falls } \sum_{i=1}^{10} X_i^2 \leq 9.34 \\ 1 & \text{sonst} \end{cases}$$

Verwenden Sie 2000 Beispiele zum Trainieren und 10000 zum Testen.

```
n<-12000;p<-10
set.seed(100)
x<-matrix(rnorm(n*p),ncol=p)
y<-as.factor(c(0,1)[as.numeric(apply(x^2,1,sum)>9.34)+1])

indtrain<-sample(1:n,2000,FALSE)
train<-data.frame(y=y[indtrain],x[x[indtrain,]])
indtest<-setdiff(1:n,indtrain)
test<-data.frame(y=y[indtest],x[x[indtest,]])
```

- b) Erstellen Sie einen Klassifikationsbaum auf dem Trainings-Datensatz der nur aus dem Stumpf besteht. Schätzen Sie die Fehlerrate mithilfe des Test-Datensatzes ab. Hinweis:

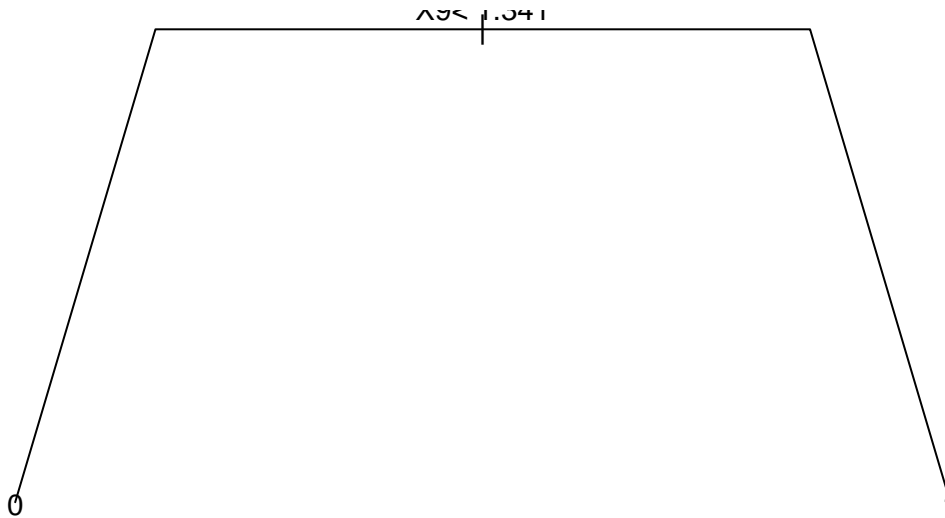
```
treeSingleStump<-rpart(y ~ .,data=train,
  control=list(maxdepth=1, cp=-1,minsplit=0,xval=0),method="class")
```

```
library(rpart)
library(mda)
```

```
## Loading required package: class
```

```
## Loaded mda 0.4-9
```

```
treeSingleStump<-rpart(y~.,data=train,control=list(maxdepth=1,
  cp=-1,minsplit=0,xval=0),method="class")
plot(treeSingleStump, branch=0.7, uniform=T)
text(treeSingleStump, use.n=TRUE, cex=0.9)
```



```
print(treeSingleStump)
```

```
## n= 2000
##
## node), split, n, loss, yval, (yprob)
##      * denotes terminal node
##
## 1) root 2000 990 1 (0.4950000 0.5050000)
##   2) X9< 1.340778 1815 861 0 (0.5256198 0.4743802) *
##   3) X9>=1.340778 185  36 1 (0.1945946 0.8054054) *
```

```
confusion(true=test$y, predict(treeSingleStump,newdata=test,type="class"))
```

```
##           true
## predicted    0    1
##           0 4794 4329
##           1  233  644
```

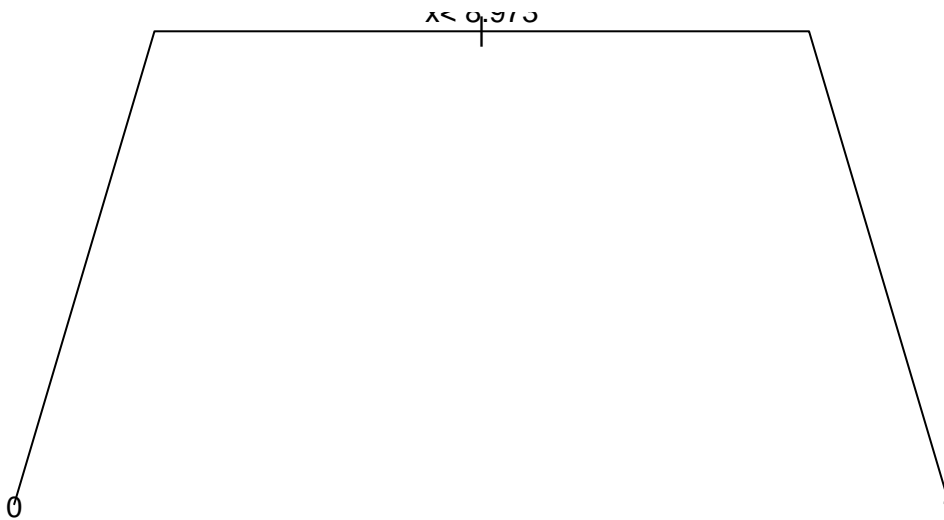
- c) Erstellen Sie einen Klassifikationsbaum basierend auf den internen Stopp-Kriterien, berechnen Sie die Fehlerrate.

```
tree.default<-rpart(y~.,data=train,method="class")
plot(tree.default, branch=0.7, uniform=T)
text(tree.default, use.n=TRUE, cex=0.9)
```


- f) Angenommen Sie kennen den Prozess der Datenerzeugung wie in a) beschrieben bis auf die Konstante 9.34. Überlegen Sie sich eine geschickte Transformation der Variablen $X_1, X_2 \dots X_{10}$ in eine neue Variable X und wenden Sie diese Transformation an. Verwenden Sie einen Klassifikationsbaum wie in c) und bestimmen Sie die Fehlerrate.

```
X = apply(train, MARGIN = 1, FUN=function(x) sum(x[2:10]^2))
train1 = data.frame(y=as.factor(train$y), x=X)
X = apply(test, MARGIN = 1, FUN=function(x) sum(x[2:10]^2))
test1 = data.frame(y=as.factor(test$y), x=X)
```

```
tree.default<-rpart(y~.,data=train1,method="class")
plot(tree.default, branch=0.7, uniform=T)
text(tree.default, use.n=TRUE, cex=0.9)
```



```
confusion(true=test1$y,predict(tree.default,newdata=test1,type="class"))
```

```
##           true
## predicted    0    1
##           0 4900  682
##           1  127 4291
```