

## (8) Prädiktive Modellierung

Ziel dieser Sitzung ist es, das Vorgehen typischer Analysen der prädiktiven Modellierung auf ihrer Implementierungsebene zu verdeutlichen. Dafür wollen wir das Prinzip der  $n$ -fachen *leave-one-out cross-validation* mithilfe eines “Nearest-Neighbour-Klassifikationsansatzes” verdeutlichen: Zu zwei Klassen zweidimensionaler Trainingsfeaturevektoren bestimmen wir jeweils das Stichprobenmittel und die Stichprobenkovarianz und ordnen einen zweidimensionalen Testfeaturevektor dann derjenigen Klasse zu, deren Stichprobenmittel es im Sinne der Mahalanobis-Distanz in Bezug zur jeweiligen Stichprobenkovarianzmatrix näher liegt.

### Datensatz generieren

Wir erstellen zunächst einen simulierten Datensatz mit  $n = 40$  Datenpunkten, wobei je die Hälfte der Datenpunkte auf eine Klasse entfallen und jeder Datenpunkt aus einer bivariaten Normalverteilung mit klassenspezifischen Erwartungswertparameter und klassenspezifischem Kovarianzmatrixparameter kommt.

```
# R-Paket
library(mvtnorm) # multivariate Normalverteilung
set.seed(0)      # reproduzierbare Ergebnisse

# Modellparameter
m = 2 # Featurevektordimension
n = 40 # Anzahl Trainingsdatenpunkte
mu_0 = c(1,1) # w.a.u. Erwartungswertparameter von Klasse 0
mu_1 = c(2,2) # w.a.u. Erwartungswertparameter von Klasse 1
Sigma_0 = matrix(c( 0.5, -0.3,
                  -0.3, 0.5),
                byrow = TRUE,
                nrow = m) # w.a.u. Kovarianzmatrixparameter von Klasse 0
Sigma_1 = matrix(c( 1.0, 0.5,
                  0.5, 1.0),
                byrow = TRUE,
                nrow = m) # w.a.u. Kovarianzmatrixparameter von Klasse 1

# Datensimulation
y = matrix(rep(NaN,n) , nrow = 1) # Labeldatenarray
x = matrix(rep(NaN,n*m), nrow = m) # Featurevektorarray
for(i in 1:n){ # Iteration über Datenpunkte

  # klassenabhängige Datengeneration
  if(i <= n/2){
    y[i] = 0 # Label
    x[,i] = rmvnorm(1, mu_0, Sigma_0) # Featurevektor
  } else {
    y[i] = 1 # Label
    x[,i] = rmvnorm(1, mu_1, Sigma_1) # Featurevektor
  }
}

D = rbind(x,y) # Datensatz konkatenieren
fname = "Prädiktive_Modellierung.csv" # Dateiname
write.csv(D, file = fname, row.names = FALSE) # Datenspeichern
```

Mit folgendem **R**-Code visualisieren wir den Datensatz zusammen mit den klassenspezifischen Stichprobenmitteln und Stichprobenkovarianzmatrizen visualisieren. Das Ergebnis wird in [Abbildung 1](#) dargestellt.

```

# Laden einer m x n Datenmatrix eines Datensatzes mit binären Klassen {0,1}
D = read.csv("Prädiktive_Modellierung.csv")
X = as.matrix(D[1:2,])
m = nrow(X)
y = as.vector(D[3,])
L = c(0,1)

# Datensatz
# Featurevektoren beider Klassen
# Featurevektordimension
# Label beider Klassen
# Klassenlabels

# klassenspezifische Deskriptivstatistiken
x_bar = matrix(rep(NA,m*2) , nrow = 2)
C = array(rep( NA,m*m*2), dim = c(m,m,2))
for (l in L){
  Xl = X[,y == l]
  n = ncol(Xl)
  I_n = diag(n)
  J_n = matrix(rep(1,n^2), nrow = n)
  x_bar[,l+1] = (1/n)* Xl %*% J_n[,1]
  C[, , l+1] = (1/(n-1))*Xl %*% (I_n-(1/n)*J_n) %*% t(Xl)
}

# m x 2 Array für klassenspezifische Stichprobenmittel
# m x m x 2 Array für klassenspezifische Stichprobenkovarianzmatrizen

# x^{(i)} für Label l
# Anzahl Datenvektorealisierungen
# Einheitsmatrix I_n
# 1_{nn}
# Stichprobenmittel
# Stichprobenkovarianzmatrix

# Abbildungsparameter
library(latex2exp)
par(
  family = "sans",
  mfcol = c(1,1),
  pty = "s",
  bty = "l",
  lwd = 1,
  las = 1,
  mgp = c(2,1,0),
  xaxs = "i",
  yaxs = "i",
  font.main = 1,
  cex = 1,
  cex.main = 1)

# Gaussian Isokontur Paket
library(ellipse)
iso_0 = ellipse(C[,1], level = 0.50, centre = x_bar[,1])
iso_1 = ellipse(C[,2], level = 0.50, centre = x_bar[,2])

# Klasse 0
plot(iso_0,
  type = "l",
  col = "Red",
  xlim = c(0,4),
  ylim = c(0,4),
  xlab = TeX("$x_1$"),
  ylab = TeX("$x_2$"))
points(x_bar[1,1], x_bar[2,1],
  col = "Red",
  pch = 3,
  cex = 1.5)

# Klasse 1
lines(iso_1[,1], iso_1[,2],
  col = "Blue")
points(x_bar[1,2], x_bar[2,2],
  col = "Blue",
  pch = 3,
  cex = 1.5)

# Daten y^{(i)} = 0
points(D[1,D[3,] == 0], D[2,D[3,] == 0],
  col = "White",
  bg = "Red",
  pch = 21)

# Daten y^{(i)} = 1
points(D[1,D[3,] == 1], D[2,D[3,] == 1],
  col = "White",
  bg = "Blue",
  pch = 21)

# Legende
legend("topleft", c("y = 0", "y = 1"),
  pch = 16,
  col = c("Red", "Blue"),
  bty = "n",
  cex = 1,
  x.intersp = 1)

# Speicherung

```

```
dev.copy2pdf(  
  file = "./Abbildungen/datensatz.pdf",  
  width = 4.5,  
  height = 4.5)  
dev.off()
```

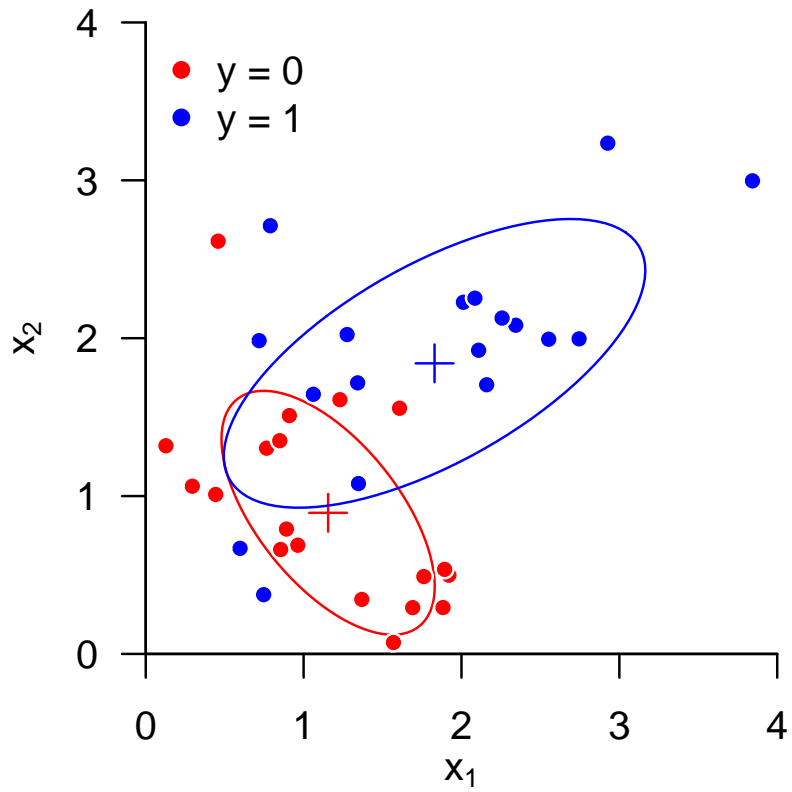


Abbildung 1. Beispieldatensatz mit Datenpunkten aus zwei Klassen  $y = 0$  (rot) und  $y = 1$  (blau).

## *n*-fache leave-one-out Kreuzvalidierung

Mit folgendem **R**-Code führen wir im Sinne des oben beschriebenen Klassifizierungsansatzes “Nearest Neighbour” eine *n*-fache *leave-one-out cross-validation* durch.

```
# Datensatz
D = read.csv("Prädiktive_Modellierung.csv")
L = c(0,1)
x = as.matrix(D[1:2,])
y = as.matrix(D[3,])
n = ncol(y)
y_pred = matrix(rep(NA,n*2), nrow = n)

# Datensatz
# Klassenlabels
# Featurevektoren
# Label
# Anzahl Datenpunkte
# Array wahrer und prädizierter Label

# n-fache leave-one-out cross-validation
for(i in 1:n){

  # Datensatzpartition
  x_train = as.matrix(x[,-i])
  y_train = as.matrix(y[,-i])
  x_test = as.matrix(x[, i])
  y_test = as.matrix(y[, i])
  y_pred[i,1] = y_test

  # iter Featurevektor nicht im Trainingsdatensatz
  # iter Label nicht im Trainingsdatensatz
  # iter Featurevektor als Testdatenpunkt
  # iter Label als Testdatenpunkt
  # wahres Label des iten Labels

  # klassenspezifische Mahalanobisdistanz ("Training")
  D = matrix(rep(NA,n*2), ncol = 2)
  for (l in L){
    Xl = x_train[y_train == l]
    nl = ncol(Xl)
    I_nl = diag(nl)
    J_nl = matrix(rep(1,nl^2), nrow = nl)
    x_bar = (1/nl)* Xl %*% J_nl[,1]
    C = (1/(nl-1))*Xl %*% (I_nl-(1/nl)*J_nl) %*% t(Xl)
    D[l+1] = t(x_test - x_bar) %*% solve(C) %*% (x_test - x_bar)
  }
  # Iteration über Klassen
  # x^{(i)} für Label l
  # Anzahl Datenvektorealisierungen
  # Einheitsmatrix I_n
  # 1_{nn}
  # Stichprobenmittel
  # Stichprobenkovarianzmatrix
  # Mahalanobisdistanz

  # Prädiktion des iten Labels ("Test")
  if (D[i] <= D[i+1]){y_pred[i,2] = 0} else {y_pred[i,2] = 1}
  # prädiziertes Label
}

# Evaluation der leave-one-out cross-validation
rp = sum(y_pred[y_pred[,1] == 1,2] == 1)
rn = sum(y_pred[y_pred[,1] == 0,2] == 0)
fp = sum(y_pred[y_pred[,1] == 0,2] == 1)
fn = sum(y_pred[y_pred[,1] == 1,2] == 0)
ACC = (rp+rn)/(rp+fp+rn+fn)
SEN = rp/(rp+fn)
SPE = rn/(rn+fp)
# Anzahl richtig positiver Prädiktionen (1,1)
# Anzahl richtig negativer Prädiktionen (0,0)
# Anzahl falsch positiver Prädiktionen (0,1)
# Anzahl falsch positiver Prädiktionen (1,0)
# Genauigkeit ("Accuracy")
# Sensitivität
# Spezifität

# Ergebnisausgabe
cat("Accuracy : ", ACC,
    "\nSensitivity : ", SEN,
    "\nSpecificity : ", SPE, "\n")
```

```
Accuracy : 0.75
Sensitivity : 0.85
Specificity : 0.65
```

Wir erhalten für den Beispieldatensatz eine Genauigkeit von 0.75, eine Sensitivität von 0.85 und eine Spezifität von 0.65.