

## Statistisches Data Mining (StDM)

### Woche 7

#### Aufgabe 1 Boston crime data (LDA vs logistische Regression)

Verwenden Sie den Boston Datensatz (boston aus package MASS) um vorherzusagen, ob ein Stadtteil (Zeile im Datensatz) eine grosse Kriminalitätsrate hat (grösser als der Median aller Stadtteile). Verschaffen Sie sich einen Überblick über die Daten und verwenden Sie logistische Regression sowie LDA. Teilen Sie das Datenset in ein Trainings und in ein Testset (je 50% der Daten)

- a) Logistische Regression
- b) LDA

#### Aufgabe 2 Weinerkennung (Leave one out crossvalidation)

Der in dieser Aufgabe betrachtete Datensatz enthält die Resultate der chemischen Analyse von 178 italienischen Weinen. Es wurden die folgenden Grössen gemessen:

<b>Alcohol:</b>	Alkoholgehalt	<b>Malic:</b>	Apfelsäure
<b>Ash:</b>	Menge der Rückstände	<b>Alcalinity:</b>	Basizität von Rückständen
<b>Magnesium:</b>	Magnesiumgehalt	<b>Phenols:</b>	Totaler Phenolgehalt
<b>Flavanoids:</b>	Geschmack tragende Phenole	<b>Nonflavanoid:</b>	neutrale Phenole
<b>Proanthocyanins:</b>	Proanthocyanin	<b>Intensity:</b>	Farbintensität
<b>Hue:</b>	Farbton	<b>OD280:</b>	OD280/OD315
<b>Proline:</b>	Prolingehalt		

- a) Laden Sie die Daten und fertigen einen Boxplot der Daten an. Was fällt auf?
- b) Verwenden Sie LDA, um die Weinsorte vorherzusagen. Beginnen Sie mit einem naiven Ansatz bei dem Sie auf dem gesamten Datensatz trainieren und vorhersagen. Berechnen Sie die Konfusionsmatrix zum Beispiel mit dem Befehl `confusion` aus dem `mda` Paket.
- c) Schreiben Sie eine Leave-One-Out Kreuzvalidierung in dem Sie die i-te Zeile für das Training des LDA Klassifiers weglassen und dann die i-te Zeile vorhersagen. Erzeugen Sie sich so einen Vektor, der jeweils die vorhergesagte Sorte enthält und bestimmen Sie so die Konfusionsmatrix.

Tipp: Die i-te Zeile können Sie mit `x[-i,]` weglassen.

- d) Eine schnelle LOO-Kreuzvalidierung kann man durch die Option `CV=TRUE` im `lda` Befehl erreichen, schauen Sie sich die Hilfe an und berechnen Sie die Konfusionsmatrix analog zu c).

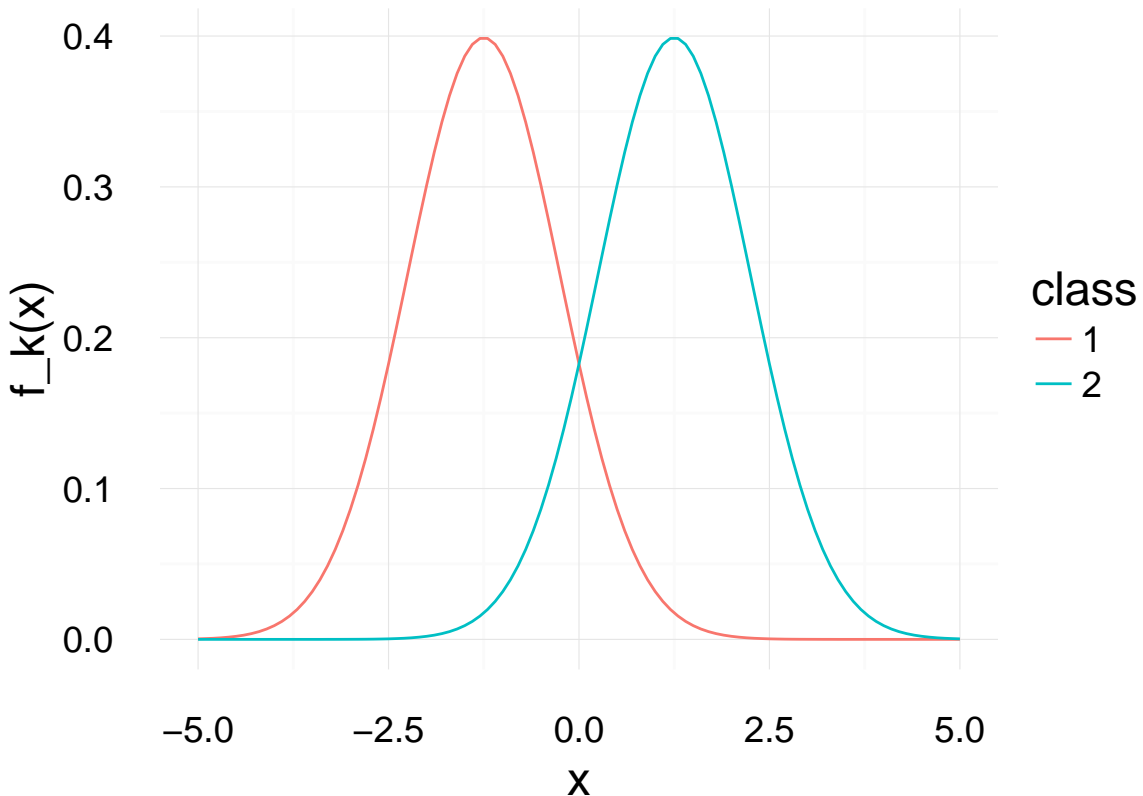
### Aufgabe 3 Simulationsstudie (LDA im Vergleich zu Bayes)

Ziel dieser Aufgabe ist es festzustellen, wie nahe ein Klassifikator an der optimalen Lösung ist. Wie in der Vorlesung besprochen ist die LDA als Bayes Klassifikator optimal, falls man die Verteilungen kennt. Allerdings ist die Situation, bei der man die Verteilung kennt etwas artifiziell. Wir simulieren daher Daten aus einer bekannten Verteilung und versuchen evaluieren die Performance des LDA Klassifikators im Vergleich zum Bayes Klassifikator

Wir betrachten im folgenden 2 normalverteilte Likelihooddichten  $f_k(x)$ , für die Klasse  $k = 0$  um  $\mu_0 = -1.25$  verteilt, für die Klasse 1 um  $\mu_1 = 1.25$  verteilt. Die Standardabweichung ist jedesmal 1.

- a) Was ist die optimale 'Accuracy'? Fragen Sie sich zuerst wo ist hier die Grenze ist, bei der Sie  $x$  zur Klasse 1 bzw. zur Klasse 2 zuordnen. Tragen Sie Konfusionsmatrix die jeweiligen Prozentzahlen ein.

## Warning: package 'ggplot2' was built under R version 3.2.4



- b) Simulationsstudie. Ziehen Sie nun  $N = 10$  Zufallszahlen aus der Klasse 1 und 2, wobei  $\pi_1 = \pi_2 = 0.5$  ist. Sie können dazu die Funktion `makeData` verwenden (versuchen Sie den code zu verstehen). Lernen Sie einen LDA Klassifier mit diesen Daten und wenden ihn auf 100 Testdaten an, die Sie immer neu aus der gleichen Verteilung ziehen. Wiederholen Sie das Experiment 1000 mal. Wie ist die Verteilung der Accuracy, was ist der Mittelwert? Vergleichen Sie ihn mit a)

```
makeData = function(N=10) {  
  z = rbinom(1,N,prob=0.5)  
  #####  
  # Ein kleiner Hack. Die lda Funktion muss Beispiele von beiden Klassen sehen können  
  # Sollte man eigentlich in der Auswertungsfunktion berücksichtigen
```

```

# Allerdings ist für grosses N der Fehler vernachlässigbar.
if (z == 0) {
  z == 1
}
if (z == N) {
  z == N-1
}
y = c(rep(0, z), rep(1, N-z))
x = c(rnorm(z, -1.25, 1), rnorm(N-z, +1.25, 1))
return (data.frame(y = as.factor(y), x = x))
}

```

- c) Wiederholen Sie Aufgabe b) und bestimmen die mittlere Accuracy für  $N = 5, 10, 20, 50, 70, 100, 200$ . Tragen Sie dies geeignet auf. Gegen welchen Wert strebt die mittlere Accuracy?