

## WERKZEUGE MUSTERERKENNUNG & MASCHINELLES LERNEN

# Aufgabenblatt 7

(Ausgabe am Fr 26.5.2017 — Abgabe bis So 11.6.2017)

### Aufgabe 1 10P

In dieser Aufgabe geht es um die Fisher-Diskriminanten (Lineare Diskriminanzanalyse nach Kitano, ME-Skript V.6).

- Laden Sie wieder den Irisdatensatz, lesen Sie die sechs Datensätze aus `fda.rda` ( $\leadsto$  Aufgabenwebseite) ein und verwenden Sie Ihre alte (korrigierte) Funktion `plot_lfd` (Übung 6, Aufgabe 2).
- Schreiben Sie eine Funktion `class_scatter(x, f)`, die für den mit Faktor `f` etikettierten Datensatz `x` den Mittelwertvektor  $\mu$  und die drei Streuungsmatrizen  $S$ ,  $S_W$ ,  $S_B$  (total, Inner- und Außerklassen) berechnet und in einer Liste mit den Einträgen `mean`, `total`, `within`, `between` zurückliefert.
- Erweitern Sie `class_scatter()` um einen Test ('R'-Funktion `stopifnot`) auf die Gültigkeit der Zerlegung  $S = S_W + S_B$  und korrigieren Sie nötigenfalls die Kovarianzberechnung; lesen Sie dazu bitte `?cov` durch.
- Schreiben Sie nun eine Funktion `FDA(x, train=x, n=)` zur FDA-Transformation des Datensatzes `x`. Es sind die `n` ersten Diskriminanten zu berechnen. Wählen Sie Kitanos Kernmatrix  $S_W^{-1}S_B$  und verwenden Sie die Funktionen `class_scatter()` und `eigen()`; alles weitere wie bei `PCA()`.
- Erweitern Sie `FDA()` um ein Argument `method=c('FDA', 'PCA', 'BSA', 'ID')` für die Alternativen `PCA` (gewöhnliche PCA) und `BSA` ('between-scatter' Analyse), welche als Kernmatrix der Transformation  $Q = S$  bzw.  $Q = S_B$  statt  $Q = S_W^{-1}S_B$  (im Fall `FDA`) zu Grunde legen. (Bei `'ID'` entfällt das Transformieren.)
- Starten Sie nun für jeden der sieben Datensätze eine  $(2 \times 2)$ -Leinwand und zeichnen Sie den Scatterplot für die jeweils beiden ersten
  - Originalmerkmale, (2) PCA-Merkmale, (3) BSA-Merkmale, (4) FDA-Merkmale.

- Datensatz `ALDI` enthält Personen zweier Musterklassen ( $\pm$  `jena`) mit ihren Erwerbshäufigkeiten einschlägiger Konsumartikel als Merkmale. Nutzen Sie einen geschickten `FDA()`-Aufruf um herauszubekommen, welche fünf der zwanzig gelisteten Produkte die verlässlichsten Indikatoren für den Wohnsitz Jena sind.

Abzugeben ist die Datei `fda.R` mit Ihrem Programmcode und Ihre schriftliche Antwort zu (g).

### Aufgabe 2 10P

Laden Sie das 'R'-Paket `class` mit dem Kommando `library(class)` und lesen Sie sich die Beschreibung zu den Methoden `knn` und `knn.cv` des Nächste-Nachbarin-Klassifikators (ME-Skript VI.6, Blatt 14,15) durch.

- Schreiben Sie eine 'R'-Funktion `knn.heldout(train, test, k=1)`, die einen Klassifikatortest mit den angegebenen Lern- und Testdaten durchführt und dabei die `k`-Nächster-Nachbar-Regel der Funktion `knn` verwendet; der Ausgabewert sei die Fehlerrate.
- Schreiben Sie eine 'R'-Funktion `knn.leave1out(data, k=1)`, welche einen Kreuzvalidierungstest („leave-one-out“) durchführt; der Ausgabewert sei wiederum die Fehlerrate.
- Laden Sie die `diabetes.rda`-Daten und führen Sie die `k`-NN-Klassifikation durch (Lern- und Testdaten). Verwenden Sie die Befundgrößen `k` in  $\{1, 2, 3, 5, 7, 10, 14, 19, 25\}$  und schreiben Sie alle Fehlerraten (in Prozent!) in eine erste Tabellenzeile.
- Vertauschen Sie nun die Rolle von Lern- und Testdatensatz und füllen Sie die nächste Tabellenzeile. Wie deuten Sie die Ergebnisse?
- Gehen Sie jetzt zur Kreuzvalidierung über (dritte Zeile); verwenden Sie dazu die Vereinigungsmenge von Lern- und Testdatensatz. Was gibt es Auffälliges zu berichten?
- Wiederholen Sie die Kreuzvalidierung zweimal (vierte und fünfte Zeile). Wo und warum unterscheiden sich die Resultate?
- Laden Sie jetzt Lern- und Testdatensatz aus `letter.rda` und evaluieren Sie die 1-NN-Regel, wobei Sie stets den Testdatensatz zur Fehlerwertung verwenden, aber aus den Lern- und Testdaten Anfangspartien wechselnder Längen zwischen 1 und 8192 Mustern (`T ∈ 2^(0:13)`) zum Anlernen selektieren.
- Analysieren Sie jetzt den Gesamtdatensatz (Lern+Test) aus `germany.rda`. Ermitteln Sie bitte für jedes  $x_1, \dots, x_{24}$  die Fehlerrate (1-NN-Regel/Kreuzvalidierung), die sich mit den restlichen 23 Merkmalen — nach Entfernung des aktuellen  $x_n$  — ergibt (Knockoutrate).
- Welches Merkmal ist am unverzichtbarsten? Welches Merkmal ist am überflüssigsten? Sind die Unterschiede dieser minimalen/maximalen Fehlerrate zu denjenigen des vollständigen Merkmalinventars eigentlich statistisch signifikant?

Abzugeben sind der R-Code `knnrule.R`, eine Datei `errors.rda` mit den drei Fehlerraten tabellen zu (c,d,e,f), (g) und (h) als `save`-Matrixobjekte (mit informativen Beschriftungen in `colnames` und `rownames`) und Ihre Kommentare zu (d,e,f,i) als **schriftlicher** Lösungsteil.

## Hinweise zum Übungsablauf

---

- ✦ Die wöchentliche WMM-Vorlesung findet am Mittwoch um 12:15 Uhr statt.  
Das Aufgabenblatt gibt es immer am Freitag (PDF im Netz).  
Der späteste Abgabetermin ist Sonntag 23:59 Uhr.
- ✦ Die Übungsaufgaben dürfen natürlich (und sollten sogar) in Gruppenarbeit (2 Mitglieder) gelöst werden.
- ✦ Schriftliche Lösungen („*Textantworten*“) sind als PDF beizufügen oder direkt im e-Mail-Textkörper unterzubringen.
- ✦ Alle anderen Lösungen (Programmieraufgaben, Daten und Grafiken) sind als elektronischer Anhang der Lösungs-e-Mail abzuliefern.
- ✦ Programmcode (Dateien \*.R) muss auch wirklich in 'R' ausführbar sein.  
(Kommando `Rscript <name.R>` auf einem der Rechner des FRZ-Pools)
- ✦ Ganz wichtig:  
Schriftliche Antworten werden von mir gedruckt, gelesen, kommentiert und korrigiert.  
Deshalb diese Textteile bitte **niemals** im abgegebenen Programmcode verstecken!
- ✦ Je Gruppe und je Aufgabenblatt ist **genau eine** e-Mail zu senden:
  - Vermerk »WMM/*n*« und Gruppenname im **subject**-Feld  
(*n* ∈ ℕ ist die laufende Nummer des Übungsblattes)
  - die Namen der beteiligten Gruppenmitglieder im Texttrumpf
  - Tabellen, Bilder, Programmcode, Sensordaten als Attachments  
(elektronische Anlagen)
  - etwaige schriftliche Antworten im Texttrumpf der Post oder als Attachment  
(Text/PDF)
- ✦ *Pfingstfrieden*: Am Freitag 2.6. gibt es kein Übungsblatt. Die Lösungen für das Übungsblatt vom Freitag 26.5. müssen erst am Sonntag 11.6. abgeliefert werden.
- ✦ Einige Aufgabentexte verweisen Sie zum Nachschlagen von Details auf das Folienskript zur Vorlesung Mustererkennung; Sie finden es unter der URL  
<http://www.minet.uni-jena.de/fakultaet/schukat/ME/Scriptum/>.  
Die Angabe *ME-Skript II.6* bedeutet: Kapitel II, Abschnitt 6

WWW: <http://www.minet.uni-jena.de/www/fakultaet/schukat/WMM/SS17>  
e-Mail: [EG.Schukat-Talamazzini@uni-jena.de](mailto:EG.Schukat-Talamazzini@uni-jena.de)