

# Statistische Methoden der Datenanalyse

Markus Schumacher

## Übung XII

Matthew Beckingham und Henrik Nilsen

28.01.2010

### Anwesenheitsaufgaben

#### Aufgabe 45 Fisherdiskriminante und Likelihoodverhältnis

Betrachten Sie die Fisherdiskriminante für den Fall zweier Wahrscheinlichkeitsdichtefunktionen  $f(\vec{x}|H_0)$  und  $f(\vec{x}|H_1)$ , die jeweils multidimensionale Gaussverteilungen mit identischen Kovarianzmatrizen  $V_0 = V_1 = V$  sein sollen. Die Wahrscheinlichkeitsdichtefunktion ist also gegeben als:

$$f(\vec{x}|H_k) = \frac{1}{(2\pi)^{n/2}|V|^{1/2}} \exp \left[ -\frac{1}{2}(\vec{x} - \vec{\mu}_k)^T V^{-1}(\vec{x} - \vec{\mu}_k) \right] \quad k \in 0,1 \quad (1)$$

(i) Zeigen Sie, dass das Likelihoodverhältnis gegeben ist durch

$$r = \frac{f(\vec{x}|H_0)}{f(\vec{x}|H_1)} = \exp(t), \quad (2)$$

wobei  $t$  die Fisherdiskriminante

$$t(\vec{x}) = a_0 + (\vec{\mu}_0 - \vec{\mu}_1)^T V^{-1} \vec{x} \quad (3)$$

mit einem beliebigem Schwellenwert  $a_0$  ist. Dementsprechend ist eine Optimierung des Likelihoodverhältnisses äquivalent zu einer Optimierung der Fisherdiskriminante.

(ii) Benutzen Sie das Bayes-Theorem mit Aprioriwahrscheinlichkeiten  $\pi_0$  und  $\pi_1$  für  $H_0$  bzw.  $H_1$ , um zu zeigen, dass die bedingte Wahrscheinlichkeit für  $H_0$  bei gegebenen Daten  $\vec{x}$  gegeben ist durch

$$P(H_0|\vec{x}) = \frac{1}{1 + e^{-t}} = s(t), \quad (4)$$

wobei  $s(t)$  die Logistische Funktion ist. Betrachten Sie dazu eine Redefinition des Schwellenwertes von der Form  $a'_0 = a_0 + \log \frac{\pi_0}{\pi_1}$ .

#### Aufgabe 46 Maximale Separation der Fisherdiskriminante

Ein Maß für die Separation zweier Hypothesen  $H_0$  und  $H_1$  mittels einer Fisherdiskriminanten ist gegeben durch

$$J(\vec{a}) = \frac{\vec{a}^T B \vec{a}}{\vec{a}^T W \vec{a}}, \quad (5)$$

wobei die Fisherdiskriminante definiert ist als

$$t(\vec{x}) = \vec{a}^T \vec{x}. \quad (6)$$

Die Separation zwischen  $H_0$  und  $H_1$  wird repräsentiert durch

$$B_{ij} = (\mu_0 - \mu_1)_i (\mu_0 - \mu_1)_j, \quad (7)$$

und die Summe der Kovarianzen zwischen den zwei Klassen ist gegeben durch

$$W_{ij} = (V_0 + V_1)_{ij}. \quad (8)$$

- (i) Bilden Sie die Ableitung  $\partial J(\vec{a})/\partial \vec{a}$  von  $J(\vec{a})$  nach  $\vec{a}$  und zeigen Sie, dass das Maximum von  $J(\vec{a})$  durch die Eigenwertgleichungen

$$W^{-1}B\vec{a} = \lambda\vec{a} \quad (9)$$

gegeben ist.

- (ii) Zeigen Sie, dass für einen beliebigen Vektor  $\vec{a}$  der Vektor  $B\vec{a}$  parallel zu  $(\vec{\mu}_0 - \vec{\mu}_1)$  ist.  
(iii) Zeigen Sie damit, dass

$$\vec{y} \propto W^{-1}(\vec{\mu}_0 - \vec{\mu}_1) \quad (10)$$

eine Lösung der Eigenwertgleichungen aus (i) ist und daher  $J(\vec{a})$  maximiert.

# Hausaufgaben

## Aufgabe 47 Zählexperiment für eine Signal- und Untergrundmessung

12 Punkte

Betrachtet wird ein Experiment, dessen Ziel die Entdeckung eines neuen Teilchens oberhalb eines von momentanen Theorien vorhergesagten Untergrundes ist. Dabei könnte es sich beispielsweise um das Higgs-Boson oder auch um supersymmetrische Teilchen handeln, wobei die Untergrundvorhersage durch das Standardmodell der Teilchenphysik erfolgt.

Die zu betrachtenden Hypothesen sind also die Nullhypothese  $H_0$ , dass nur Ereignisse aus Untergrundprozessen gemessen wurden, sowie die Alternativhypothese  $H_1$ , dass sowohl Signal- als auch Untergrundereignisse beobachtet wurden.

Im Experiment wurde eine Gesamtanzahl  $x$  von Ereignissen aufgezeichnet. Weiterhin wurde eine andere, signalfreie kinematische Region definiert, aus der man die Normierung des Untergrundes bestimmen kann. In dieser Region wurden  $y$  Ereignisse gefunden. Das Verhältnis der Untergrundereignisse in der signalfreien Region zu denjenigen in der Signalregion sei  $\tau$ . Die mittlere Anzahl der Untergrundereignisse in der Signalregion sei  $b$  und die mittlere Anzahl der Signalereignisse im Falle von Hypothese  $H_1$  sei  $s$ .

- (i) Stellen Sie die Likelihoodfunktionen für die Hypothesen  $H_0$  und  $H_1$  auf. Nehmen Sie dabei an, dass die Gesamtanzahlen von Ereignissen in Signal- und Kontrollregion jeweils Poissonverteilt sind.
- (ii) Betrachten Sie jetzt die Schätzer für  $s$  und  $b$  unter der Hypothese  $H_1$  ( $\hat{s}$  bzw.  $\hat{b}$ ) sowie den Schätzer für  $b$  unter der Nullhypothese,  $\hat{\hat{b}}$ .
  - a) Stellen Sie die Profile Likelihood  $\lambda$  auf.
  - b) Bestimmen Sie Ausdrücke für  $\hat{s}$ ,  $\hat{b}$  und  $\hat{\hat{b}}$ . Betrachten Sie dazu

$$\left. \frac{\partial L(H_0)}{\partial b} \right|_{\hat{\hat{b}}}, \quad (11)$$

sowie die beiden gleichzeitigen Einschränkungen

$$\left. \frac{\partial L(H_1)}{\partial s} \right|_{\hat{s}, \hat{b}} \quad \text{und} \quad \left. \frac{\partial L(H_1)}{\partial b} \right|_{\hat{b}, \hat{s}} \quad (12)$$

- c) Berechnen Sie  $q = -2 \log \lambda$ , und zeigen Sie dadurch dass der in der Vorlesung gegebene Gleichung

$$q = 2[x \ln x + y \ln y - (x + y) \ln \left( \frac{x + y}{1 + \tau} \right) - y \ln \tau] \quad (13)$$

korrekt ist.

## Aufgabe 48 Zählexperiment für eine Signal- und Untergrundmessung 2

8 Punkte

In dieser Aufgabe sollen Sie die Ergebnisse von Aufgabe 47 benutzen, um die erwartete Sensitivität des ATLAS-Experimentes auf eine Entdeckung des Higgs-Bosons im Zerfallskanal  $H \rightarrow \gamma\gamma$  zu ermitteln. Monte-Carlo-Studien von Proton-Proton-Kollisionen im ATLAS-Detektor haben gezeigt, dass der Wirkungsquerschnitt für  $pp \rightarrow H + X \rightarrow \gamma\gamma + X$ -Ereignisse, die die Ereignisselektion passieren, gegeben ist durch  $\sigma_S = 25,4$  fb. Der Wirkungsquerschnitt für Untergrundereignisse, die dieselbe Ereignisselektion passieren, beträgt  $\sigma_B = 947$  fb. In einer weiteren Analyse kann ein reiner Untergrunddatensatz mit einem Wirkungsquerschnitt von  $\sigma_T = 10300$  fb ausgewählt werden.

- (i) Benutzen Sie die Relation

$$N = \mathcal{L}\sigma \quad (14)$$

um die Anzahlen von Signal- ( $x$ ), Untergrundereignissen ( $y$ ) sowie die Anzahl von Ereignissen in der Seitenbandregion ( $\tau b$ ) für eine integrierte Luminosität von  $L = 10 \text{ fb}^{-1}$  auszurechnen.

- (ii) Berechnen Sie die Schätzer für die Anzahl der Signalereignisse  $\hat{s}$ , der Untergrundereignisse  $\hat{b}$  unter der Hypothese von Signal plus Untergrund, sowie den Schätzer  $\hat{\hat{b}}$  auf die Anzahl der Untergrundereignisse in der Nur-Untergrund Hypothese.

(iii) Berechnen Sie die Größe

$$q = -2 \log \lambda, \quad (15)$$

wobei  $\lambda$  die in Aufgabe 47 berechnete Profile Likelihood ist.

(iv) Berechnen Sie daraus die Signifikanz des vorhergesagten Signals.

(v) Nehmen Sie nun an, dass die Untergrundrate in der Signalregion mit einer relativen Genauigkeit von  $\Delta b/b = 5\%$  abgeschätzt werden kann. Wenn man annimmt, dass es sich dabei um einen Poissonfehler handelt, kann man eine effektive Seitenbandregion konstruieren mit Ereignisanzahl  $y' = \tau' b$ . Dazu setzt man die relative statistische Unsicherheit in der hypothetischen Seitenbandregion (Poissonfehler) gleich der relativen Genauigkeit der Untergrundvorhersage:

$$\frac{\sqrt{y'}}{y'} = \frac{\Delta b}{b} \iff y' = \frac{1}{(\Delta b/b)^2} \quad (16)$$

Berechnen Sie die Werte für  $y'$  und  $\tau'$  für die effektive Seitenbandregion.

(vi) Benutzen Sie die Werte für  $y'$  und  $\tau'$  sowie die ursprüngliche Anzahl von Ereignissen in der Signalregion  $x$ , um den neuen Wert für  $q$  und daher der Signifikanz für eine Messung mit einem Fehler von 5% auf die Untergrundvorhersage zu bekommen.