

### 3 Generalisierte lineare Modelle (I)

#### Aufgabe 1

Der Datensatz `fakesoep` (Download von der Veranstaltungshomepage) ist dem sozioökonomischen Panel nachempfunden und soll nun mit Hilfe von R analysiert werden. Wir betrachten folgende Variablen (erhoben an  $n = 3000$  befragten Personen):

<code>beink</code>	Bruttoverdienst im letzten Monat
<code>groesse</code>	Körpergröße
<code>alter</code>	Alter
<code>dauer</code>	Dauer der Betriebszugehörigkeit
<code>verh</code>	verheiratet (1: ja, 0: nein)
<code>geschl</code>	Geschlecht (1: weiblich, 0: männlich)
<code>deutsch</code>	Deutsche Staatsangehörigkeit (1: ja, 0: nein)
<code>abitur</code>	Abitur (1: ja, 0: nein)

- Untersuchen Sie die Variable `beink` hinsichtlich der Gestalt ihrer Verteilung. Was fällt Ihnen dabei auf?
- Im folgenden soll die Variable `beink` als Responsevariable betrachtet werden. Begründen Sie, weshalb die Gammaverteilung als Verteilungsannahme zur Modellierung dieser Responsevariablen geeignet sein könnte. Zeigen Sie nun, dass diese Verteilung zu einer Exponentialfamilie gehört. Bestimmen Sie die Größen  $\theta$ ,  $b(\theta)$ ,  $\phi$ , Erwartungswert und Varianz, sowie den natürlichen Link.
- Zeichnen Sie die Dichte einer gammaverteilten Zufallsgröße  $y$  für verschiedene Werte des `shape` Parameters, wobei für den Erwartungswert immer  $E(y) = 1$  gelten soll.
- Fitten Sie ein GLM mit gammaverteiltem Response, allen Kovariablen (Haupteffekte) und natürlichem Link. Verwenden Sie anschließend den log-Link. Interpretieren Sie die Modelle. Welche Strukturannahme würden Sie hier bevorzugen?
- Erstellen Sie einen Plot, in dem die beobachteten und die durch Ihr Modell geschätzten Werte abgebildet werden. Betrachten Sie außerdem das Ergebnis der generischen R-Funktion `plot(glm.object)`, wobei für `glm.object` das von Ihnen soeben erstellte GLM-Objekt einzusetzen ist (und setzen Sie `which = 1:5`).
- Führen Sie für beide Modelle aus (c) eine *Backward*-Variablenselektion bezüglich des AIC-Kriteriums durch. Was passiert? Welches Modell würden Sie nun bevorzugen? (*Hinweis*: Verwenden Sie die Funktion `step()`.)

#### Aufgabe 2

In der vorangegangenen Aufgabe wäre auch ein GLM mit Invers-Gauß-verteilterm Response denkbar gewesen. Eine Invers-Gauß-verteilte Zufallsgröße  $y$  hat folgende Dichte:

$$f(y|\mu, \lambda) = \sqrt{\frac{\lambda}{2\pi y^3}} \exp\left(-\frac{\lambda(y - \mu)^2}{2\mu^2 y}\right)$$

- Zeigen Sie zunächst, dass diese Verteilung zu einer Exponentialfamilie gehört. Bestimmen Sie die Größen  $\theta$ ,  $b(\theta)$ ,  $\phi$ , Erwartungswert und Varianz, sowie den natürlichen Link.
- Versuchen Sie nun, für die SOEP-Daten ein GLM mit natürlichem Link und allen Kovariablen zu fiten (Response: `beink`). Verwenden Sie anschließend den log-Link.