

1 Semi- und Nonparametrische Regression (II)

Aufgabe 1

Vorbereitungen

```
# Einlesen der Daten
fakesoep <- read.table("fakesoep.dat",header=T)
# Package mgcv laden
library(mgcv)
```

a) Datensatz aufteilen

```
fakesoepW <- fakesoep[fakesoep$geschl == 1,]
fakesoepM <- fakesoep[fakesoep$geschl == 0,]
```

In Blatt 11 wurde ein Additives Modell (AM) betrachtet. Nun haben wir einen Gamma-verteilten Response, also brauchen wir Generalisierte Additive Modelle (GAM)!

In der gam()-Funktion können dieselben Link-Funktionen sowie Families spezifiziert werden wie in glm().

```
modelW <- gam(beink ~ s(groesse,bs="cr") + s(alter,bs="cr") + s(dauer,bs="cr")
+ verh + deutsch + abitur, family=Gamma(link="log"), data=fakesoepW)
modelM <- gam(beink ~ s(groesse,bs="cr") + s(alter,bs="cr") + s(dauer,bs="cr")
+ verh + deutsch + abitur, family=Gamma(link="log"), data=fakesoepM)
```

```
# Ergebnis betrachten
summary(modelW)
```

```
##
## Family: Gamma
## Link function: log
##
## Formula:
## beink ~ s(groesse, bs = "cr") + s(alter, bs = "cr") + s(dauer,
##      bs = "cr") + verh + deutsch + abitur
##
## Parametric coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.57156    0.04657  162.593 <2e-16 ***
## verh         -0.22448    0.02653   -8.462 <2e-16 ***
## deutsch      -0.04411    0.04441   -0.993  0.321
## abitur        0.28801    0.02558  11.258 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##              edf Ref.df    F  p-value
## s(groesse)  1.001  1.001 10.91 0.000979 ***
## s(alter)    6.469  7.651  2.31 0.020579 *
```

```

## s(dauer) 6.158 7.282 17.00 < 2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) = 0.208 Deviance explained = 22.3%
## GCV = 0.17816 Scale est. = 0.17406 n = 1280

summary(modelM)

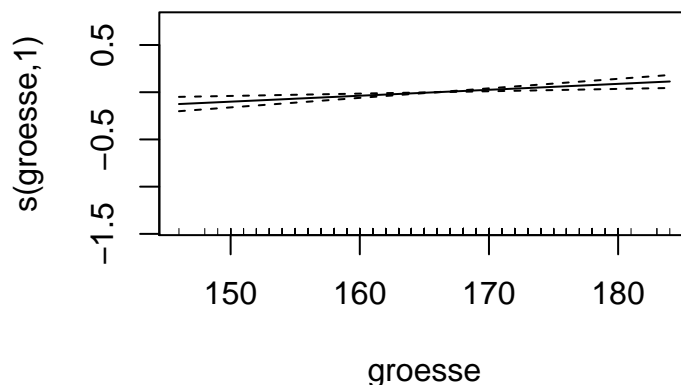
##
## Family: Gamma
## Link function: log
##
## Formula:
## beink ~ s(groesse, bs = "cr") + s(alter, bs = "cr") + s(dauer,
## bs = "cr") + verh + deutsch + abitur
##
## Parametric coefficients:
## Estimate Std. Error t value Pr(>|t|)
## (Intercept) 7.71378 0.03831 201.366 < 2e-16 ***
## verh 0.09797 0.02323 4.216 2.61e-05 ***
## deutsch 0.03077 0.03469 0.887 0.375
## abitur 0.35975 0.02125 16.925 < 2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
## edf Ref.df F p-value
## s(groesse) 1.002 1.005 38.438 6.71e-10 ***
## s(alter) 4.362 5.411 6.314 4.77e-06 ***
## s(dauer) 2.799 3.512 19.489 1.29e-13 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) = 0.25 Deviance explained = 27%
## GCV = 0.14392 Scale est. = 0.15679 n = 1720

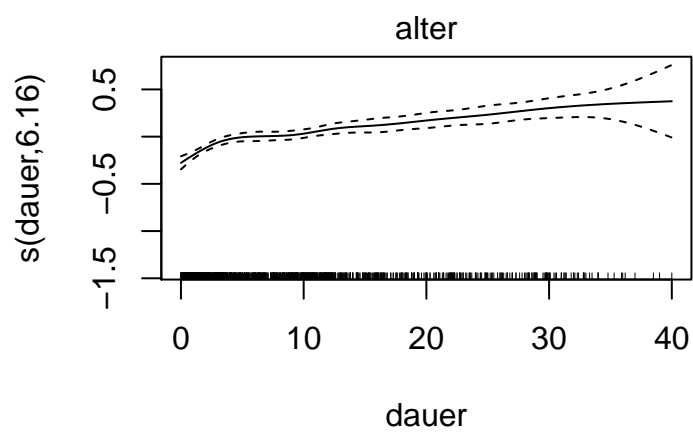
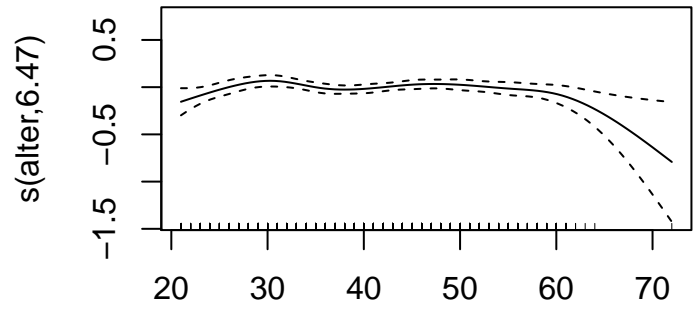
```

Nach den geschätzten Modellen wirkt sich bei Frauen eine Hochzeit im Mittel offenbar negativ auf das Einkommen aus, bei Männern hingegen positiv. Der Einfluss der Variablen Abitur ist wie erwartet.

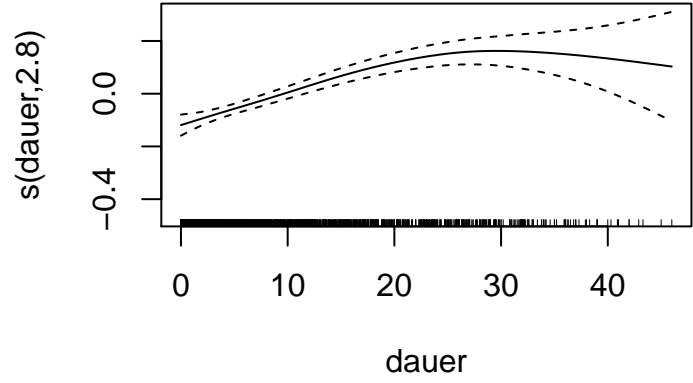
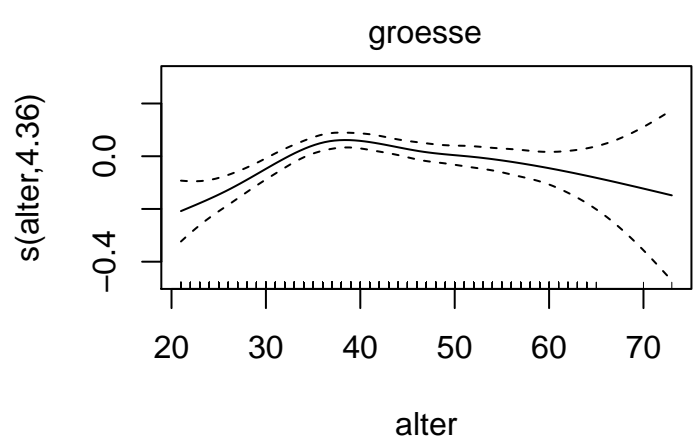
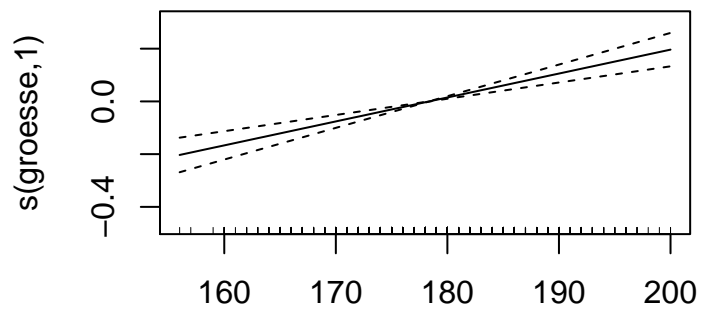
Die glatten Kurven können durch Plots dargestellt werden. Der Einfluss von Körpergröße auf Einkommen erscheint (bei Frauen und Männern) linear.

```
plot(modelW)
```





```
plot(modelM)
```



- b) Um die beiden obigen Modelle in einem Modell zusammenzufassen, muss die Variable Geschlecht sowie Interaktionen aller anderen Variablen mit Geschlecht ins Modell aufgenommen werden. Hierfür müssen zunächst Geschlechts-Dummies erzeugt werden

```
geschl1 <- fakesoep$geschl
geschl0 <- 1-geschl1

modelAll <- gam(beink ~ s(groesse,bs="cr",by=geschl0) +
s(alter,bs="cr",by=geschl0) + s(dauer,bs="cr",by=geschl0) +
s(groesse,bs="cr",by=geschl1) +
s(alter,bs="cr",by=geschl1) + s(dauer,bs="cr",by=geschl1) +
geschl + verh + deutsch + abitur + verh:geschl + deutsch:geschl + abitur:geschl,
family=Gamma(link="log"),data=fakesoep)

# An den Parameter-Schätzern der parametrisch aufgenommenen Variablen erkennt
# man die Äquivalenz der Modelle (abgesehen von kleineren Ungenauigkeiten)
summary(modelAll)

##
## Family: Gamma
## Link function: log
##
## Formula:
## beink ~ s(groesse, bs = "cr", by = geschl0) + s(alter, bs = "cr",
## by = geschl0) + s(dauer, bs = "cr", by = geschl0) + s(groesse,
## bs = "cr", by = geschl1) + s(alter, bs = "cr", by = geschl1) +
## s(dauer, bs = "cr", by = geschl1) + geschl + verh + deutsch +
## abitur + verh:geschl + deutsch:geschl + abitur:geschl
##
## Parametric coefficients:
## Estimate Std. Error t value Pr(>|t|)
## (Intercept) 7.83710 0.06269 125.011 < 2e-16 ***
## geschl -0.15435 0.16899 -0.913 0.3611
## verh 0.09825 0.02376 4.135 3.65e-05 ***
## deutsch 0.03060 0.03548 0.862 0.3885
## abitur 0.35974 0.02174 16.548 < 2e-16 ***
## geschl:verh -0.32221 0.03505 -9.193 < 2e-16 ***
## geschl:deutsch -0.07462 0.05582 -1.337 0.1814
## geschl:abitur -0.07229 0.03301 -2.190 0.0286 *
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
## edf Ref.df F p-value
## s(groesse):geschl0 1.001 1.002 36.933 1.35e-09 ***
## s(alter):geschl0 4.132 5.140 6.012 1.30e-05 ***
## s(dauer):geschl0 2.692 3.381 19.062 4.66e-13 ***
## s(groesse):geschl1 1.001 1.003 11.547 0.000682 ***
## s(alter):geschl1 6.842 7.931 2.734 0.005477 **
## s(dauer):geschl1 6.286 7.398 17.586 < 2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Rank: 62/68
## R-sq.(adj) = 0.356 Deviance explained = 38.5%
## GCV = 0.15839 Scale est. = 0.1641 n = 3000

summary(modelW)

##
## Family: Gamma
## Link function: log
```

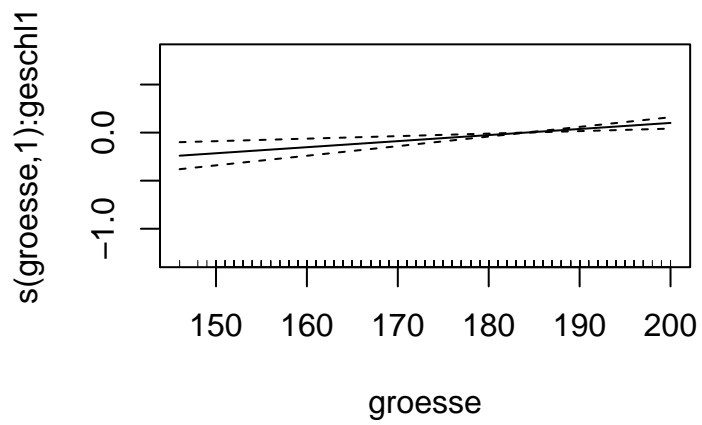
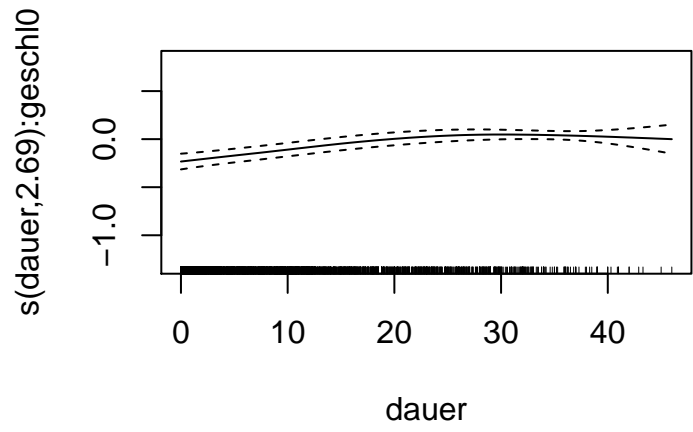
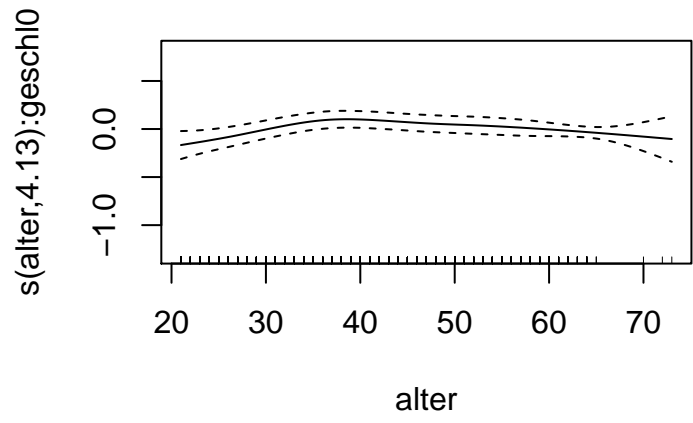
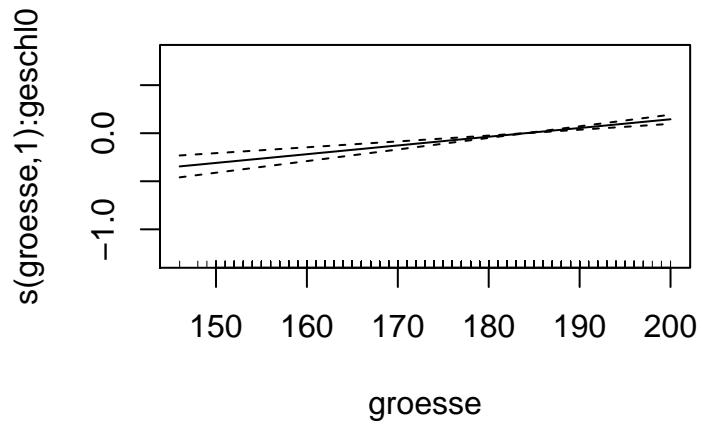
```
##
## Formula:
## beink ~ s(groesse, bs = "cr") + s(alter, bs = "cr") + s(dauer,
##      bs = "cr") + verh + deutsch + abitur
##
## Parametric coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.57156    0.04657 162.593  <2e-16 ***
## verh        -0.22448    0.02653  -8.462  <2e-16 ***
## deutsch     -0.04411    0.04441  -0.993    0.321
## abitur       0.28801    0.02558  11.258  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##           edf Ref.df      F p-value
## s(groesse) 1.001  1.001 10.91 0.000979 ***
## s(alter)   6.469  7.651  2.31 0.020579 *
## s(dauer)   6.158  7.282 17.00 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.208  Deviance explained = 22.3%
## GCV = 0.17816  Scale est. = 0.17406   n = 1280
```

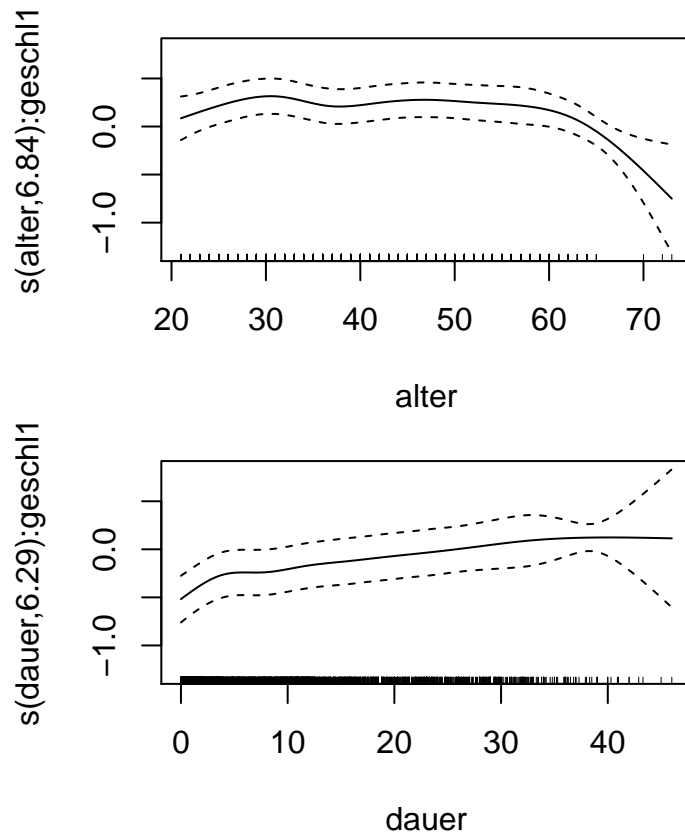
```
summary(modelM)
```

```
##
## Family: Gamma
## Link function: log
##
## Formula:
## beink ~ s(groesse, bs = "cr") + s(alter, bs = "cr") + s(dauer,
##      bs = "cr") + verh + deutsch + abitur
##
## Parametric coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.71378    0.03831 201.366 < 2e-16 ***
## verh         0.09797    0.02323  4.216 2.61e-05 ***
## deutsch     0.03077    0.03469  0.887  0.375
## abitur      0.35975    0.02125 16.925 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##           edf Ref.df      F p-value
## s(groesse) 1.002  1.005 38.438 6.71e-10 ***
## s(alter)   4.362  5.411  6.314 4.77e-06 ***
## s(dauer)   2.799  3.512 19.489 1.29e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.25  Deviance explained = 27%
## GCV = 0.14392  Scale est. = 0.15679   n = 1720
```

Auch das Ergebnis der glatten Modellierung ist ähnlich.

```
plot(modelAll)
```





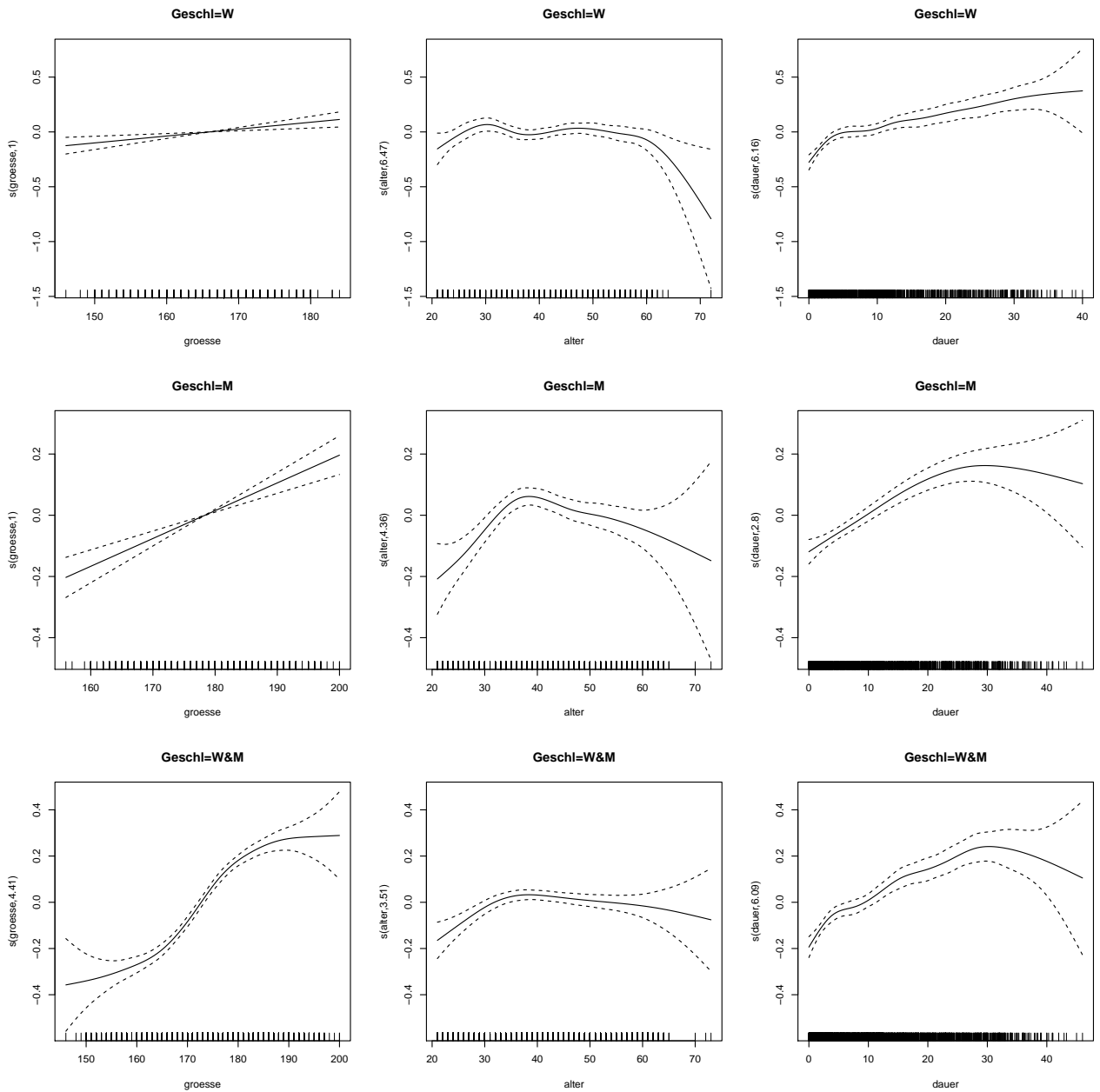
Gemeinsames Modell ohne Trennung nach Geschlecht:

```
modelWM <- gam(beink ~ s(groesse,bs="cr") + s(alter,bs="cr") + s(dauer,bs="cr")
+ verh + deutsch + abitur, family=Gamma(link="log"), data=fakesoep)
```

Vergleich

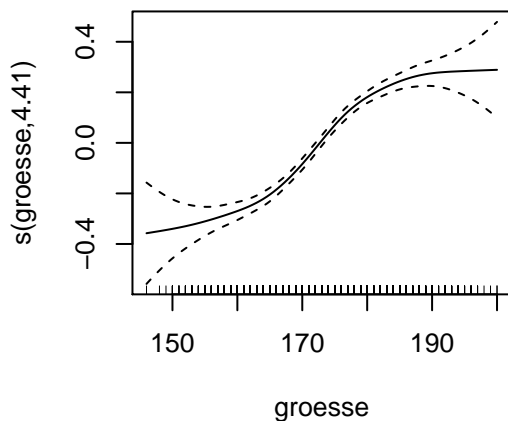
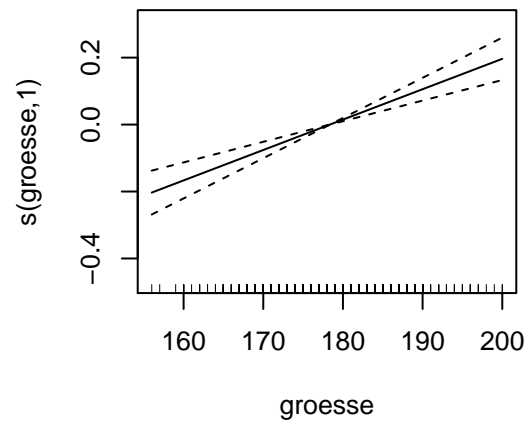
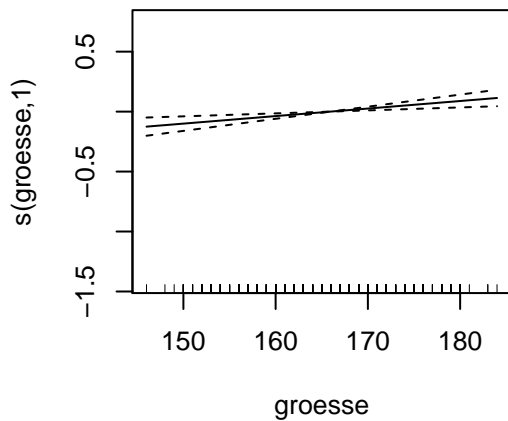
```
par(mfrow=c(3,3))

plot(modelW, main="Geschl=W")
plot(modelM, main="Geschl=M")
# für den gesamten Datensatz (Frauen und Männer zusammen)
plot(modelWM, main="Geschl=W&M")
```



Vergleich des Effekts der Größe zwischen gemeinsamen Modell und den beiden nach Geschlechtern getrennten Modellen:

```
par(mfrow=c(2,2))
plot(modelW, select=1)
plot(modelM, select=1)
plot(modelWM, select=1)
```

Der Effekt der Größe auf das Einkommen ist sowohl für Männer als auch für Frauen eine Gerade, aber das allgemeine Niveau dieses Größeneffekts unterscheidet sich zwischen den Geschlechtern, er ist bei Männern höher. Diese Niveauunterschiede sieht man in den einzelnen Plots wegen der Zentrierung der Kurven um Null herum nicht. Die nonparametrisch geschätzte Kurve in der dritten Grafik zeigt, wie der Schätzer versucht, von der niedrigeren zur höheren Gerade zu springen. In den extremen Größenbereichen, in denen fast nur eines der Geschlechter vorkommt, reproduziert der Schätzer den Verlauf der jeweiligen Gerade. Fazit: geschätzte Kurven in (G)AMs sind empfindlich gegenüber Interaktionen mit übrigen Einflussgrößen. (Hier dem Geschlecht.)

- c) Eine Alternative zur Gamma-Verteilung ist die inverse Gauß-Verteilung (vgl. Blatt 4)

```

modelAllinvG <- gam(beink ~ s(groesse,bs="cr",by=geschl0) +
s(alter,bs="cr",by=geschl0) + s(dauer,bs="cr",by=geschl0) +
s(groesse,bs="cr",by=geschl1) +
s(alter,bs="cr",by=geschl1) + s(dauer,bs="cr",by=geschl1) +
geschl + verh + deutsch + abitur + verh:geschl + deutsch:geschl + abitur:geschl,
family=inverse.gaussian(link="log"),data=fakesoep)

```

Die Ergebnisse sind sehr ähnlich.

```
summary(modelAllinvG)
```

```

##
## Family: inverse.gaussian
## Link function: log
##
## Formula:
## beink ~ s(groesse, bs = "cr", by = geschl0) + s(alter, bs = "cr",
## by = geschl0) + s(dauer, bs = "cr", by = geschl0) + s(groesse,

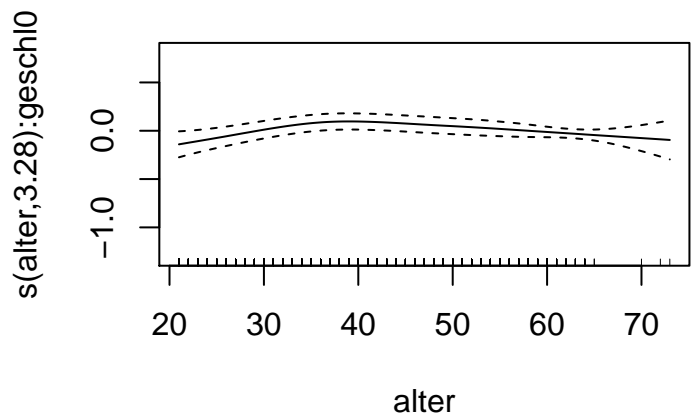
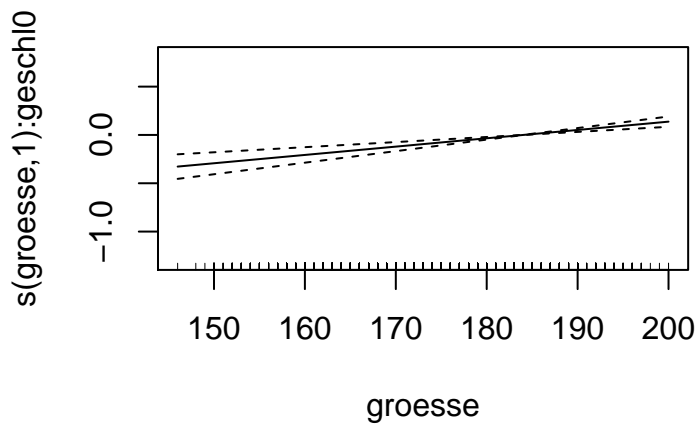
```

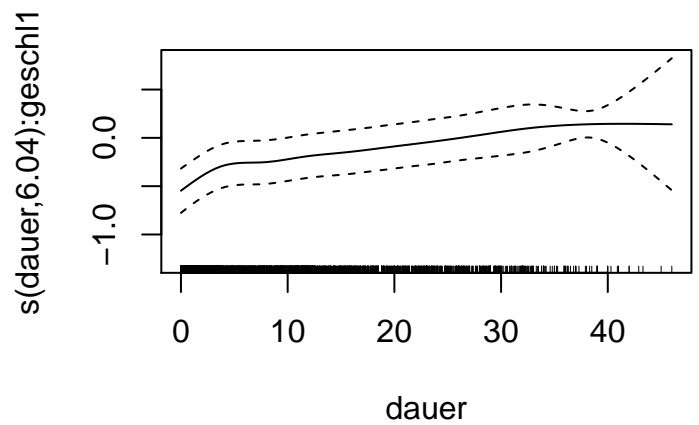
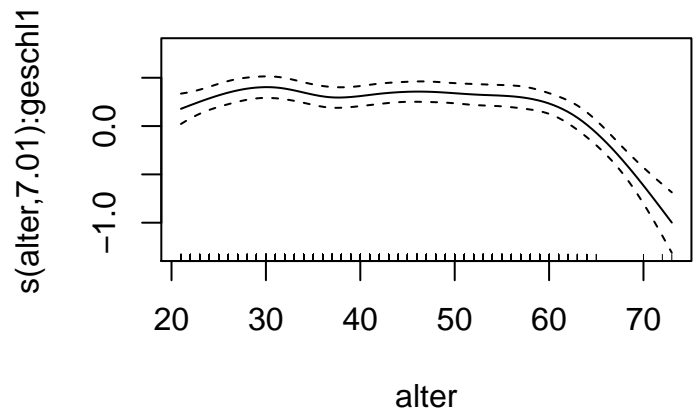
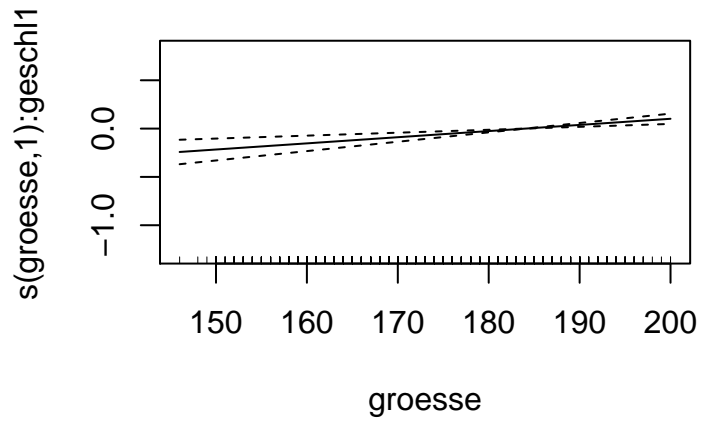
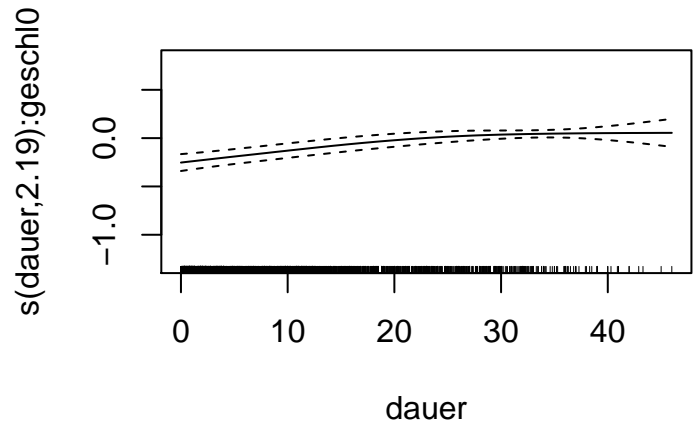
```

##      bs = "cr", by = geschl1) + s(alter, bs = "cr", by = geschl1) +
##      s(dauer, bs = "cr", by = geschl1) + geschl + verh + deutsch +
##      abitur + verh:geschl + deutsch:geschl + abitur:geschl
##
## Parametric coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    7.86131    0.06419 122.470 < 2e-16 ***
## geschl         -0.23273    0.14717  -1.581 0.113911
## verh           0.09664    0.02616   3.695 0.000224 ***
## deutsch        0.02230    0.03811   0.585 0.558429
## abitur         0.35337    0.02634  13.414 < 2e-16 ***
## geschl:verh    -0.32207    0.03542  -9.094 < 2e-16 ***
## geschl:deutsch -0.07519    0.05438  -1.383 0.166851
## geschl:abitur -0.06336    0.03560  -1.780 0.075228 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##              edf Ref.df      F p-value
## s(groesse):geschl0 1.001  1.003 26.614 2.60e-07 ***
## s(alter):geschl0   3.279  4.115  5.803 0.000105 ***
## s(dauer):geschl0   2.187  2.755 20.503 8.22e-12 ***
## s(groesse):geschl1 1.000  1.001 14.558 0.000138 ***
## s(alter):geschl1   7.013  8.046  7.829 1.71e-10 ***
## s(dauer):geschl1   6.039  7.148 25.195 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Rank: 62/68
## R-sq.(adj) = 0.353  Deviance explained = 34.2%
## GCV = 8.3975e-05  Scale est. = 7.6692e-05  n = 3000

```

```
plot(modelAllinvG)
```

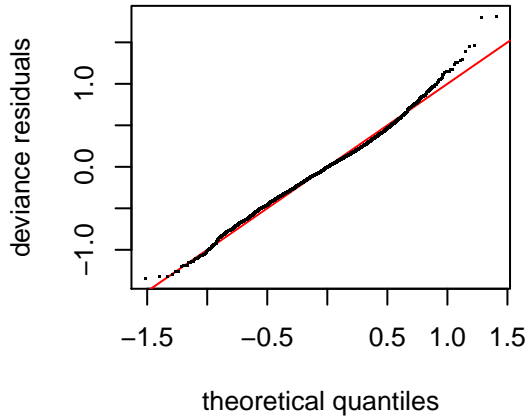
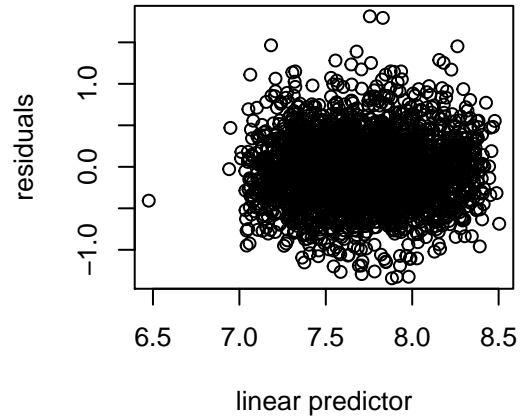




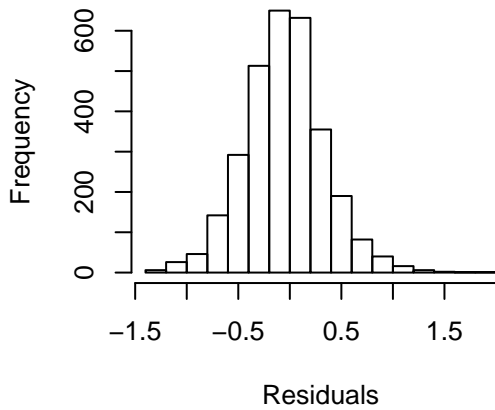
Diagnostik-Plot können mittel `gam.check()` erzeugt werden.

```
gam.check(modelA11)
```

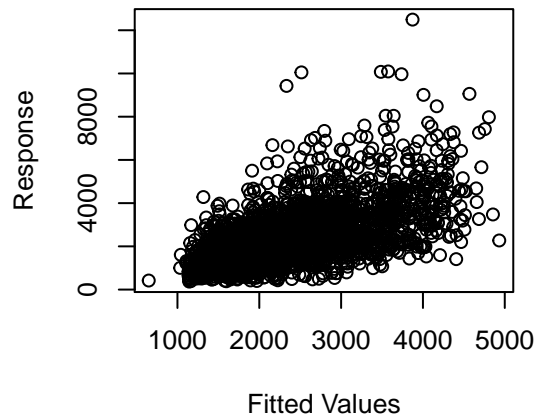
Resids vs. linear pred.



Histogram of residuals

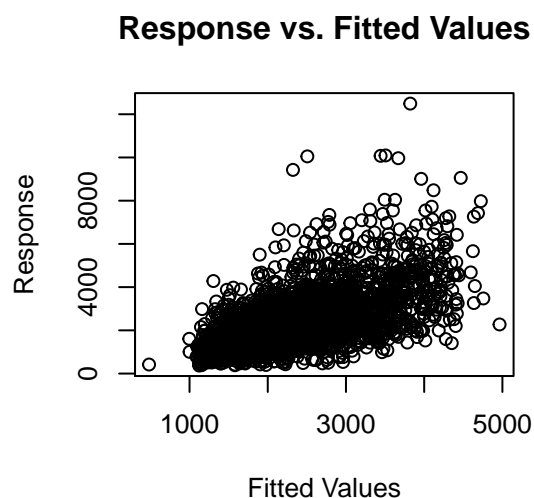
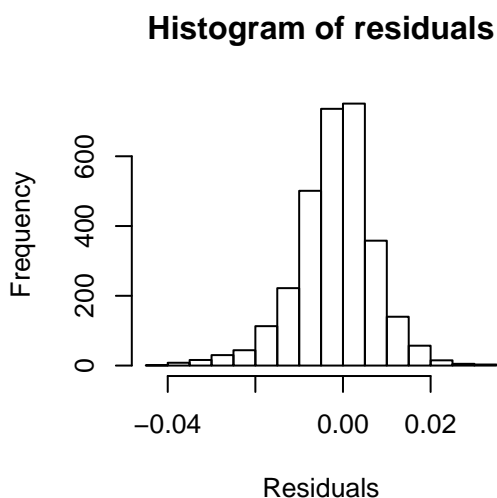
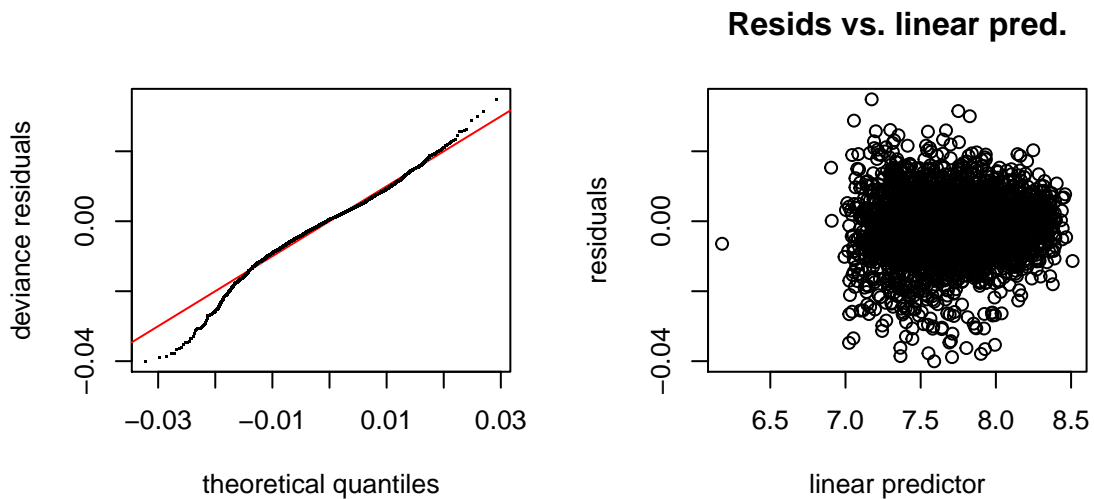


Response vs. Fitted Values



```
##
## Method: GCV   Optimizer: outer newton
## full convergence after 6 iterations.
## Gradient range [-7.537684e-08,1.05339e-09]
## (score 0.1583946 & scale 0.1641047).
## Hessian positive definite, eigenvalue range [2.330142e-08,6.977471e-05].
## Model rank = 62 / 68
##
## Basis dimension (k) checking results. Low p-value (k-index<1) may
## indicate that k is too low, especially if edf is close to k'.
##
##           k'   edf k-index p-value
## s(groesse):geschl0 10.000 1.001 0.975 0.34
## s(alter):geschl0 10.000 4.132 0.991 0.64
## s(dauer):geschl0 10.000 2.692 0.971 0.20
## s(groesse):geschl1 10.000 1.001 0.975 0.30
## s(alter):geschl1 10.000 6.842 0.991 0.71
## s(dauer):geschl1 10.000 6.286 0.971 0.26
```

```
gam.check(modelAllinvG)
```



```
##
## Method: GCV   Optimizer: outer newton
## full convergence after 7 iterations.
## Gradient range [-3.973275e-11,3.074373e-13]
## (score 8.397512e-05 & scale 7.669215e-05).
## Hessian positive definite, eigenvalue range [1.098304e-11,5.639196e-08].
## Model rank = 62 / 68
##
## Basis dimension (k) checking results. Low p-value (k-index<1) may
## indicate that k is too low, especially if edf is close to k'.
##
##           k'   edf k-index p-value
## s(groesse):geschl0 10.000 1.001 0.956 0.26
## s(alter):geschl0    10.000 3.279 0.967 0.62
## s(dauer):geschl0    10.000 2.187 0.952 0.21
## s(groesse):geschl1 10.000 1.000 0.956 0.34
## s(alter):geschl1    10.000 7.013 0.967 0.57
## s(dauer):geschl1    10.000 6.039 0.952 0.29
```

Offenbar ist die Anpassung der (Devianz-)Residuen an die Normalverteilung im Gamma-Modell besser, weshalb jenes dem Invers-Gauß-Modell vorzuziehen ist.

- d) Wir betrachten aus Gründen der Übersichtlichkeit nun Frauen und Männer wieder getrennt und schätzen zum Vergleich GLMs (vgl. Blatt 4).

```
glmW <- gam(beink ~ groesse + alter + dauer
+ verh + deutsch + abitur,family=Gamma(link=log),data=fakesoepW)
```

```
glmM <- gam(beink ~ groesse + alter + dauer
+ verh + deutsch + abitur,family=Gamma(link=log),data=fakesoepM)
```

Der Vergleich mit den GAMs aus (a) kann mittel LR-Tests erfolgen

```
anova(modelW,glmW,test="F")

## Analysis of Deviance Table
##
## Model 1: beink ~ s(groesse, bs = "cr") + s(alter, bs = "cr") + s(dauer,
##   bs = "cr") + verh + deutsch + abitur
## Model 2: beink ~ groesse + alter + dauer + verh + deutsch + abitur
##   Resid. Df Resid. Dev      Df Deviance    F    Pr(>F)
## 1      1262.4      221.81
## 2      1273.0      229.54 -10.627  -7.7315 4.18 5.891e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

anova(modelM,glmM,test="F")

## Analysis of Deviance Table
##
## Model 1: beink ~ s(groesse, bs = "cr") + s(alter, bs = "cr") + s(dauer,
##   bs = "cr") + verh + deutsch + abitur
## Model 2: beink ~ groesse + alter + dauer + verh + deutsch + abitur
##   Resid. Df Resid. Dev      Df Deviance    F    Pr(>F)
## 1      1707.8      244.06
## 2      1713.0      252.62  -5.1637  -8.5644 10.578 2.83e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Offenbar sind GAMs in beiden Fällen zur Modellierung wesentlich besser geeignet.

Insbesondere bei der Kovariablen Körpergröße stellt sich aber die Frage, ob der Einfluss nonparametrisch modelliert werden sollte. Nach den Plots in (a) zu urteilen, ist die glatte Modellierung hier nicht notwendig.

```
modelW_k <- gam(beink ~ groesse + s(alter,bs="cr") + s(dauer,bs="cr")
+ verh + deutsch + abitur,family=Gamma(link=log),data=fakesoepW)
modelM_k <- gam(beink ~ groesse + s(alter,bs="cr") + s(dauer,bs="cr")
+ verh + deutsch + abitur,family=Gamma(link=log),data=fakesoepM)
```

Paradoxerweise wird aber bei Verwendung von Tests die nonparametrische Modellierung der parametrischen vorgezogen, obwohl (bzw. weil) der geschätzte Einfluss (fast) linear ist.

```
anova(modelW,modelW_k,test="F")

## Analysis of Deviance Table
##
## Model 1: beink ~ s(groesse, bs = "cr") + s(alter, bs = "cr") + s(dauer,
##   bs = "cr") + verh + deutsch + abitur
## Model 2: beink ~ groesse + s(alter, bs = "cr") + s(dauer, bs = "cr") +
##   verh + deutsch + abitur
##   Resid. Df Resid. Dev      Df   Deviance    F    Pr(>F)
## 1      1262.4      221.81
## 2      1262.4      221.81 -0.00057356 -8.9877e-05 0.9003 0.002202 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```

anova(modelM,modelM_k,test="F")

## Analysis of Deviance Table
##
## Model 1: beink ~ s(groesse, bs = "cr") + s(alter, bs = "cr") + s(dauer,
##   bs = "cr") + verh + deutsch + abitur
## Model 2: beink ~ groesse + s(alter, bs = "cr") + s(dauer, bs = "cr") +
##   verh + deutsch + abitur
##   Resid. Df Resid. Dev      Df   Deviance      F Pr(>F)
## 1    1707.8    244.06
## 2    1707.8    244.06 -0.0024483 -0.00058527 1.5246 0.006966 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Aufgabe 2

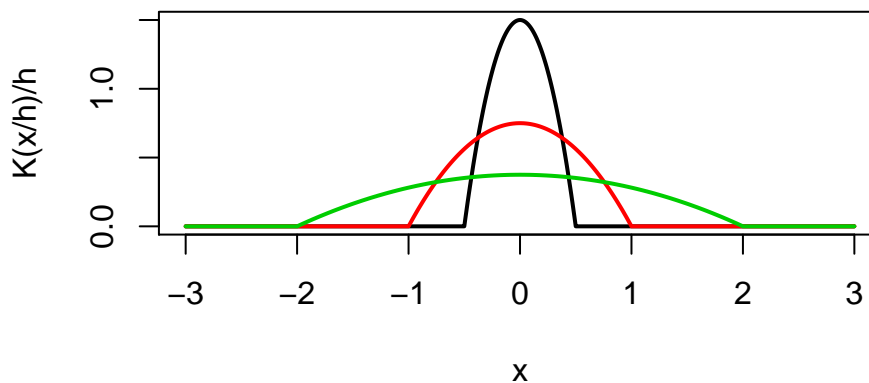
Veranschaulichung von lokaler linearer Regression

```

# Epanechnikov-Kern
epan <- function(x,hvalue)
{
  ifelse(abs(x/hvalue)<=1,(3/4)*(1-(x/hvalue)^2), 0)/hvalue
}

xseq <- seq(-3,3,by=0.01)
plot(xseq,epan(xseq,0.5),type="l",xlab="x",ylab="K(x/h)/h",lwd=2)
lines(xseq,epan(xseq,1),col=2,lwd=2)
lines(xseq,epan(xseq,2),col=3,lwd=2)

```



```

# Daten
x <- 3*runif(100)
fx <- function(x){2 - (x-1.5)^2 + 3*sqrt(x)}
y <- fx(x) + rnorm(100)
plot(x,y,xlab="x",ylab="y")
xseq <- seq(0,3,by=0.01)
lines(xseq,fx(xseq),col=3)

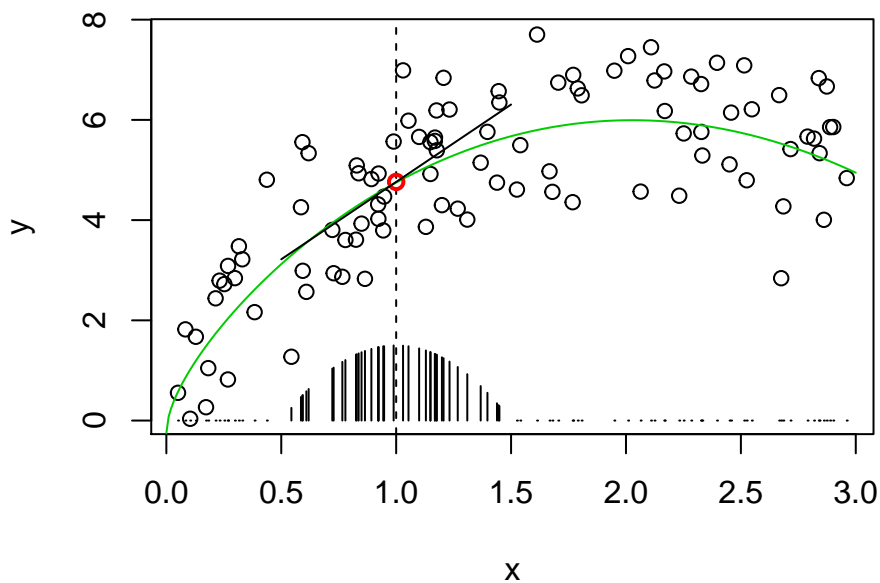
# neues x
xnew <- 1
abline(v=xnew,lty=2)

# Gewichte
w <- epan(x-xnew,0.5)
points(x,w,type="h")

```

```
# Lokale lineare Regression:
```

```
xregr <- x-xnew  
lmlokal <- lm(y ~ xregr, weights=w)  
points(xnew,lmlokal$coef[1],col=2,lwd=2)  
lines(xnew+c(-0.5,0.5),lmlokal$coef[1]+ lmlokal$coef[2]*c(-0.5,0.5))
```



```
#Für eine Kurve Grid von neuen x-Werten
```

```
xnewseq <- seq(0,3,by=0.01)  
fit <- numeric(length(xnewseq))  
i <- 1  
for (xnew in xnewseq)  
{  
  w <- epan(x-xnew,0.5)  
  xregr <- x-xnew  
  lmlokal <- lm(y ~ xregr, weights=w)  
  fit[i] <- lmlokal$coef[1]  
  i <- i+1  
}
```

```
plot(x,y,xlab="x",ylab="y")  
lines(xnewseq,fx(xnewseq),col=3)  
lines(xnewseq,fit,col=2)
```

