

Applied Statistical Regression

AS 2015 – Extending the Linear Model

Marcel Dettling

Institute for Data Analysis and Process Design

Zurich University of Applied Sciences

marcel.dettling@zhaw.ch

<http://stat.ethz.ch/~dettling>

ETH Zürich, 30. November 2015

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Generalized Additive Modelling (GAM)

Motivation:

We require a flexible regression method, similar to 1-dimensional smoothing, that also works in multiple regression setting.

Background:

The generic multiple regression formula is:

$$y_i = f(x_{i1}, x_{i2}, \dots, x_{ip}) + E_i$$

As we have argued before, this is a too challenging problem, as there are just too many functions $f(\cdot)$. While in simple regression, visualization of the function is feasible, this is no longer the case in a multiple regression where $p > 2$ (“*curse of dimensionality*”).

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Solution 1: Linear Modelling with OLS

The canonical approach for solving the multiple regression problem lies in using parametric linear models such as:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + E_i$$

As we know, the predictors x_{ij} may be transformed in any arbitrary way. However, there is no way around exactly specifying these transformations.

Since these models are linear in the parameters $\beta_0, \beta_1, \dots, \beta_p$, there is (under some mild conditions) an analytical and unique solution if the OLS algorithm is used.

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Solution 2: GAM

A ***Generalized Linear Model*** is based on the following:

$$\begin{aligned}y_i &= \beta_0 + f_1(x_{i1}) + f_2(x_{i2}) + \dots + f_p(x_{ip}) + E_i \\ &= \beta_0 + \sum_{j=1}^p f_j(x_{ij}) + E_i\end{aligned}$$

Here, $f_j(\cdot)$ are smooth, flexible, 1-dimensional functions that don't need to be explicitly defined by the user, but can be determined from the data in an explorative fashion.

There are several approaches to determine the $f_j(\cdot)$. Some are better, some are worse. The most popular approach is based on cubic splines, as explained on the next few slides...

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Simple (1-dimensional) GAM

We first explain the concept in 1-dimension, i.e. we only require to fit $f_j(\cdot)$. This is somewhat similar to smooting, but here we actually require a formula and not just visualization.

A very powerful approach is to express $f_j(\cdot)$ using some simple basis functions (i.e. transformations of x_j):

$$f_j(x_j) = \sum_{m=1}^M \gamma_m h_m(x_j)$$

Here, γ_m are some unknown coefficients that are to be estimated from data. Moreover, $h_m(\cdot)$ are arbitrary but explicitly specified basis functions. The choice of M and the complexity of $h_m(\cdot)$ controls the fit to the data.

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Polynomial Basis Functions

A simple, yet intuitive choice for the basis functions $h_m(\cdot)$ is given by powers of x_j , i.e. fitting a polynomial. In particular:

$$h_m(x_j) = x_j^m, \text{ resp. } f_j(x_j) = \sum_{m=1}^M \gamma_m x_j^m = \gamma_1 x_j + \dots + \gamma_M x_j^M$$

Polynomial basis functions have the following properties:

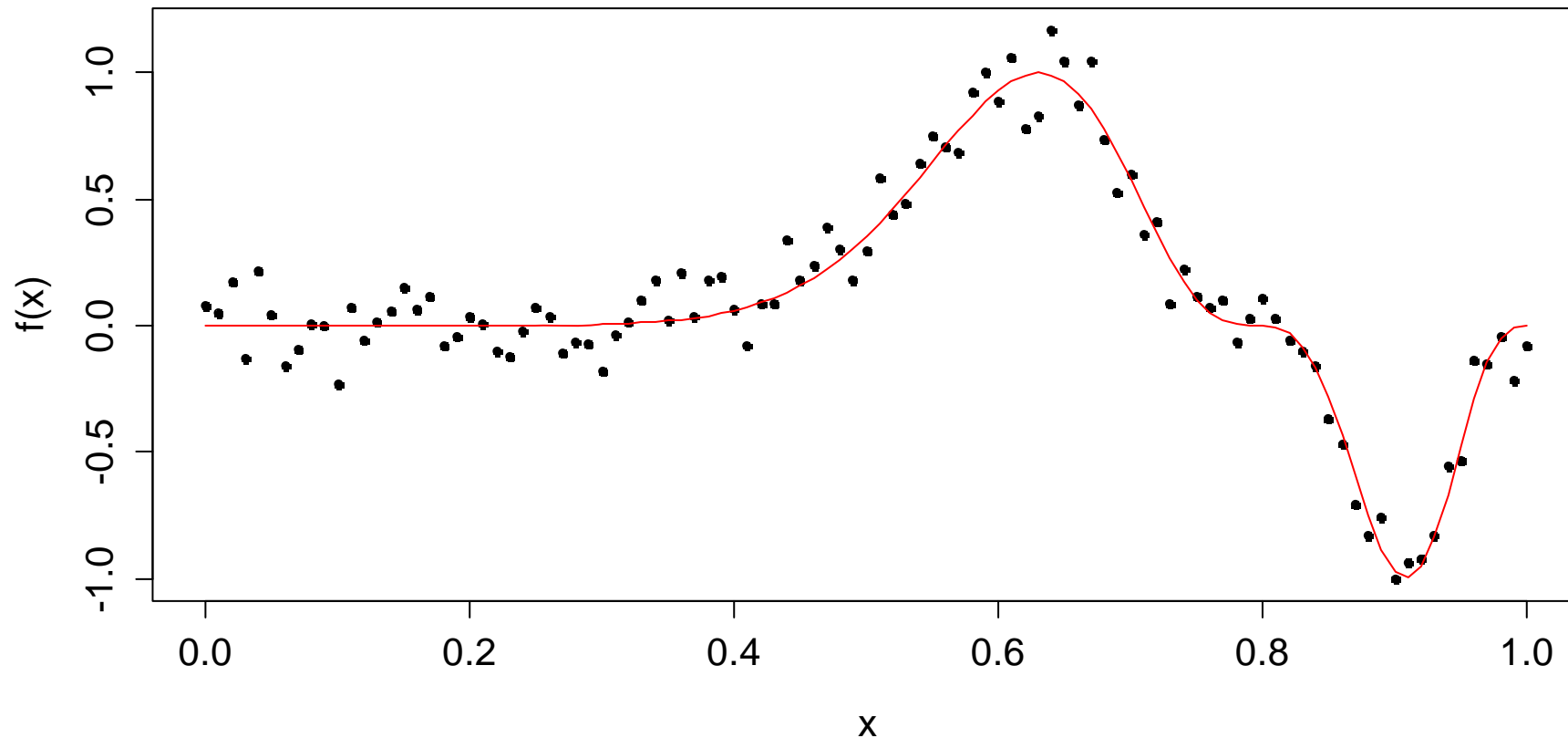
- They allow for a flexible, data-adaptive fit!!!
- Since each of the basis functions $h_m(x_j) = x_j^m$ extends over the entire range of predictor x_j , we may observe some erratic behavior, especially at the boundaries.
- Some simulations results illustrate these drawbacks...

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Example 1

True functional relation: $y = \left(\sin(2\pi x^3)\right)^3 + E$

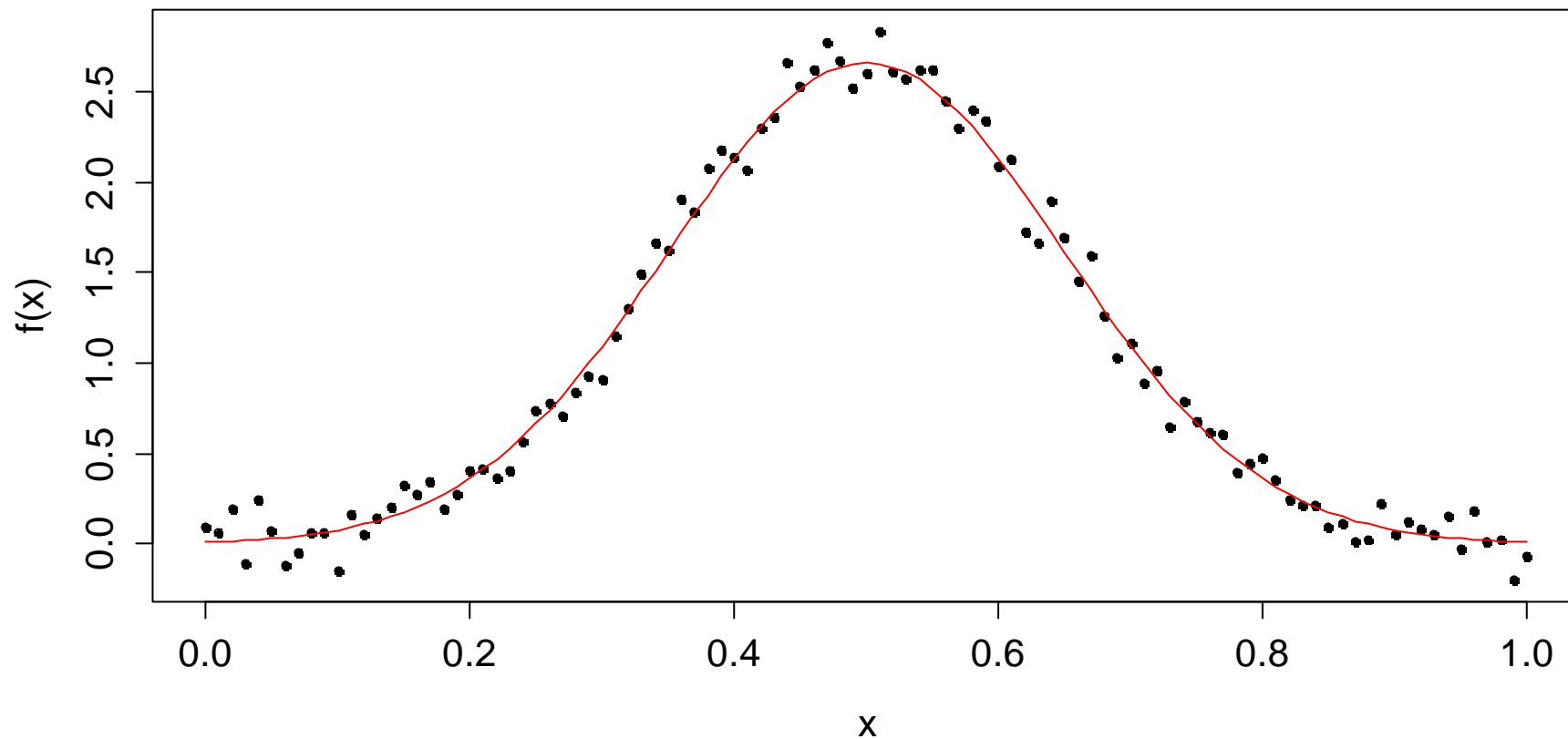


Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Example 2

True function relation: Density function of $N(0.5, 0.15^2)$

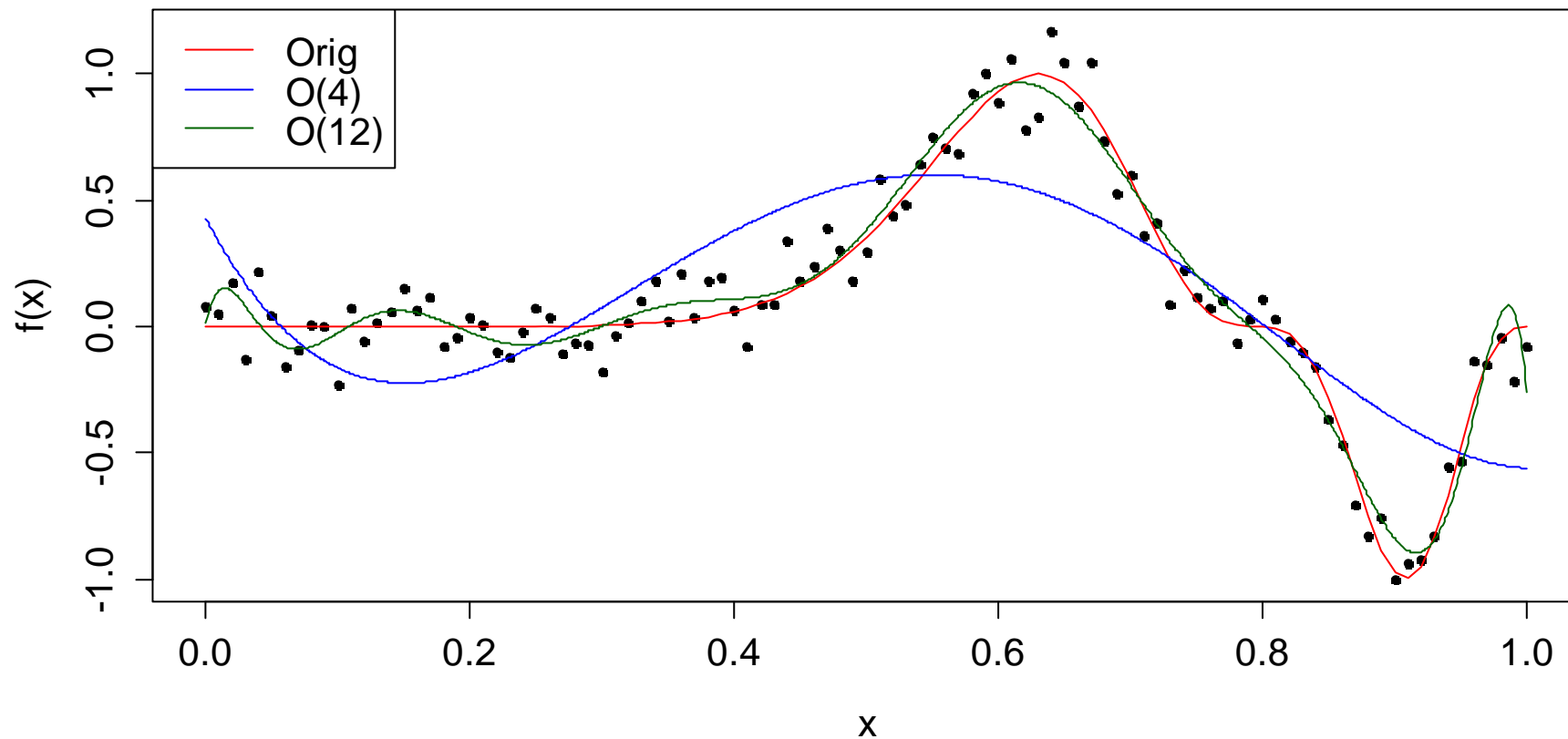


Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Smoothing with Polynomial Basis Functions

Anpassung mit polynomialer Basis

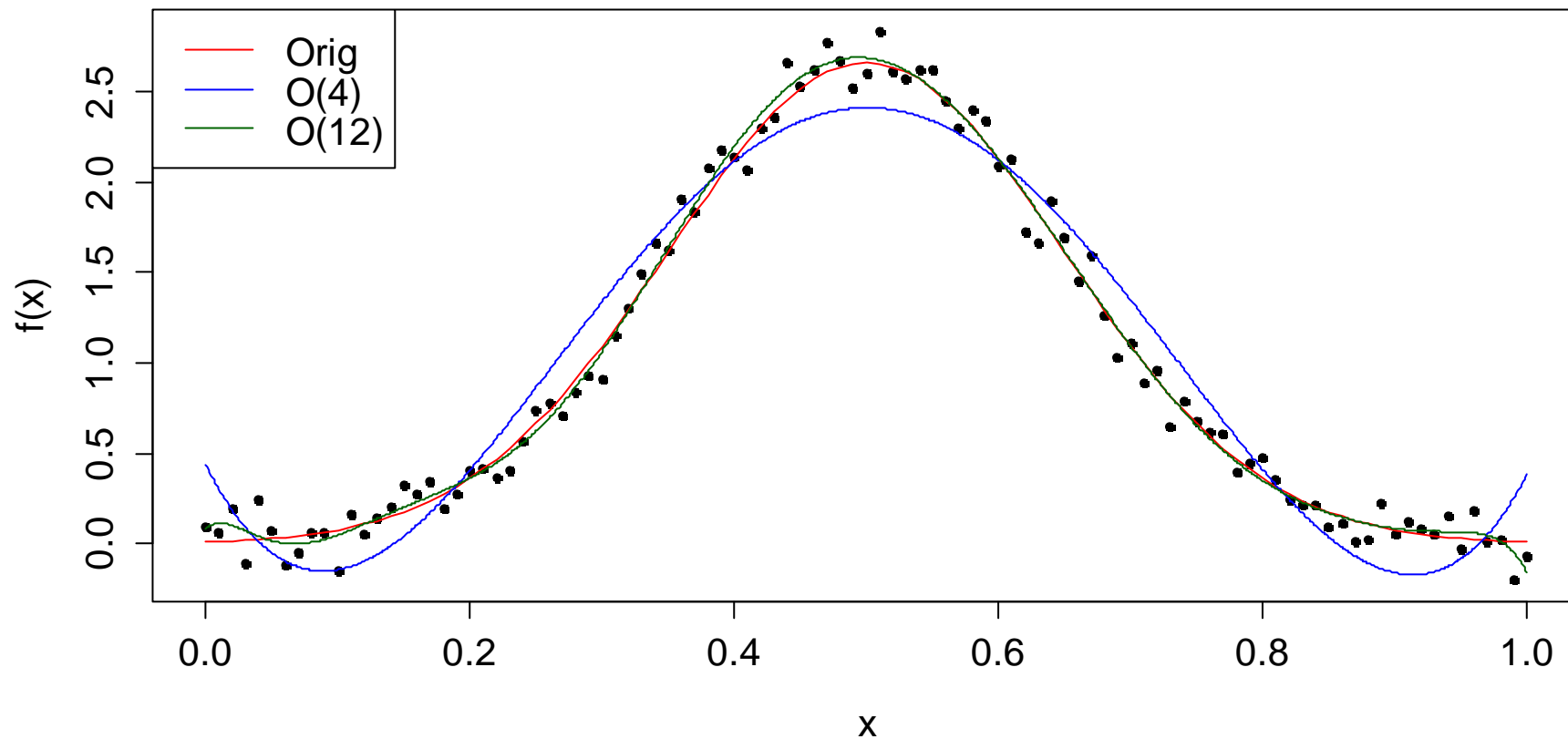


Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Smoothing with Polynomial Basis Functions

Anpassung mit polynomialer Basis

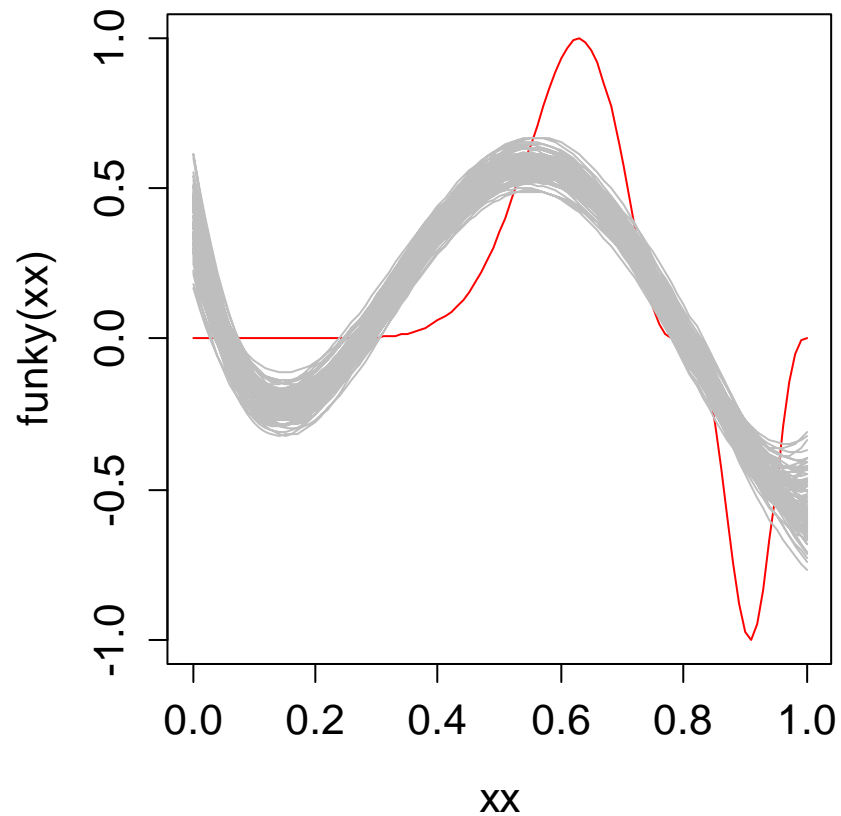


Applied Statistical Regression

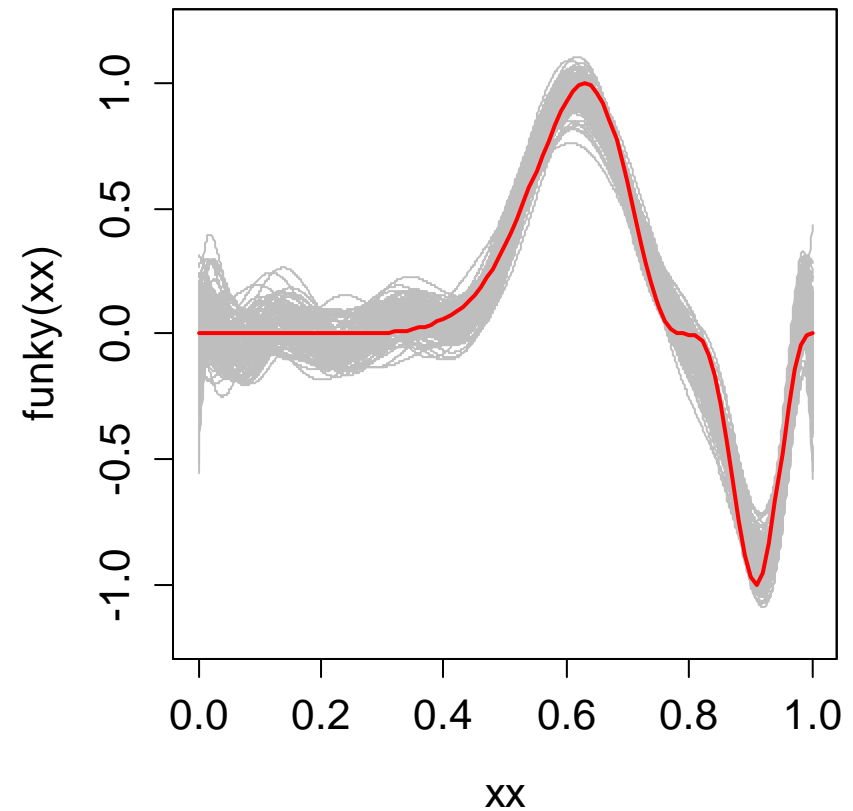
AS 2015 – Generalized Additive Modelling

Resampling on Example 1

Resampling für O(4)



Resampling für O(12)

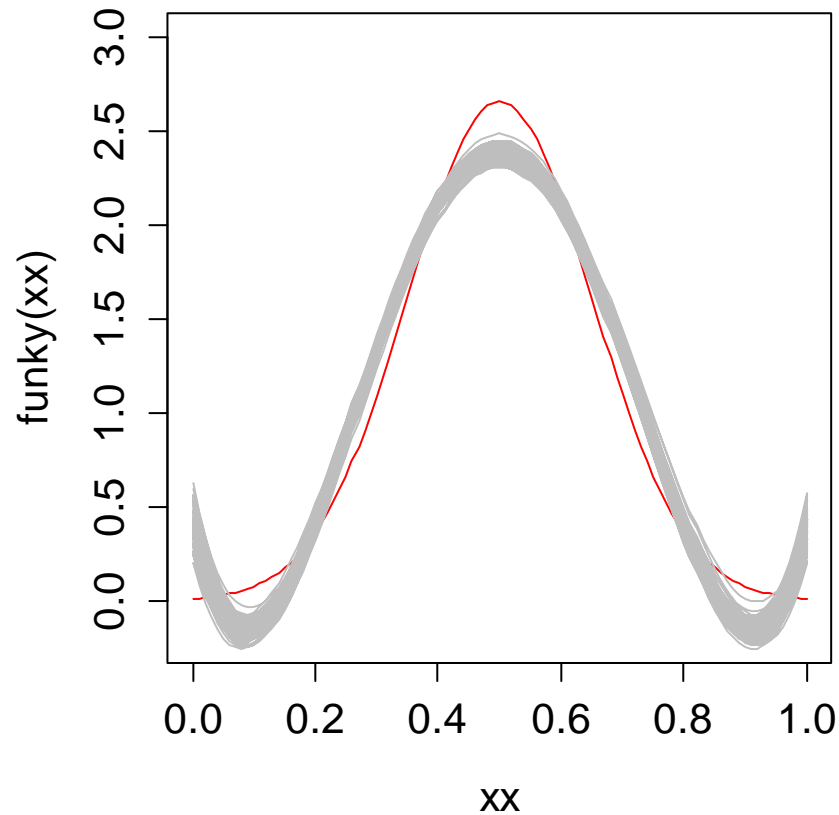


Applied Statistical Regression

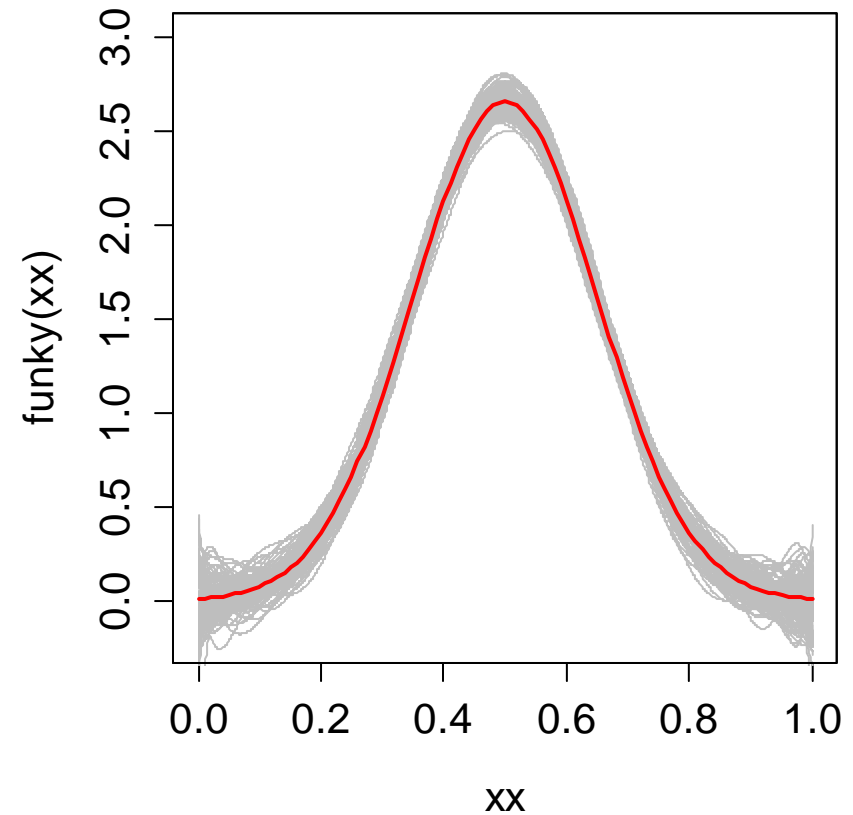
AS 2015 – Generalized Additive Modelling

Resampling on Example 2

Resampling für O(4)



Resampling für O(12)



Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

What is a Better Alternative?

As the simulation results have shown us, using polynomial basis functions has some severe drawbacks and will not result in a fruitful generalized multiple regression approach.

Idea: why not using basis functions that minimize

$$\sum_{i=1}^n (y_i - f_j(x_{ij}))^2 + \lambda \int_{-\infty}^{+\infty} f''(u) du$$

This criterion implements a trade-off between goodness-of-fit and smoothness of the function. Attractive, but how can we find a solution?

→ The solution will always be a cubic B-spline!!!

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Regression Splines

We define a basis consisting of cubic B-splines on the interval $[a, b]$ by imposing the following conditions on the knots, which are fixed at the observations x_{1j}, \dots, x_{nj} :

- 1) Each of the basis functions must be different from zero only over a range of 4 knots, so that its influence remains local.
- 2) The basic form of $h_m(\cdot)$ is a local polynomial of third order.
- 3) These basis functions are twice continuously differentiable at each of the knots. This implies smoothness of the fit consisting of numerous local functions.
- 4) The integral over all basis functions shall be equal to 1.

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Generating a Spline Basis

In R, such a regression spline basis can be generated conveniently:

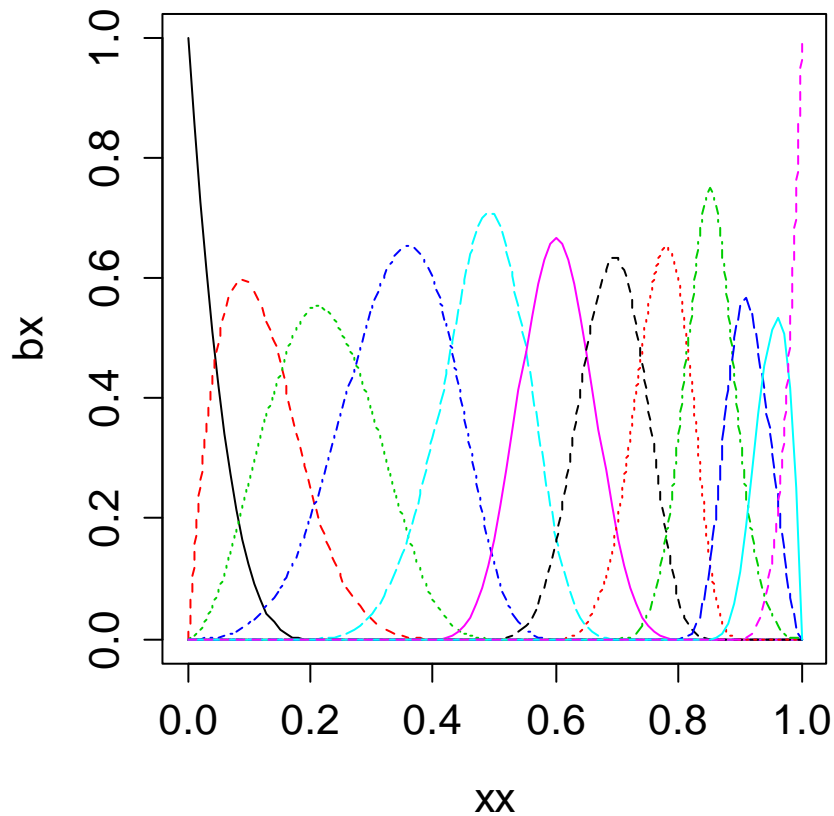
```
> set.seed(21)
> library (splines)
> funky <- function(x) sin(2*pi*x^3)^3
> xx <- seq (0, 1, by=0.01)
> yy <- funky(xx) + 0.1*rnorm (101)
> kn <-c(0,0,0,0,.2,.4,.5,.6,.7,.8,.85,.9,1,1,1,1)
> bx <- splineDesign (kn, xx)
> gs <- lm (yy ~ bx)
> matplot(xx, bx, type="l")
> matplot(xx, cbind (yy, gs$fit), type="pl")
```

Applied Statistical Regression

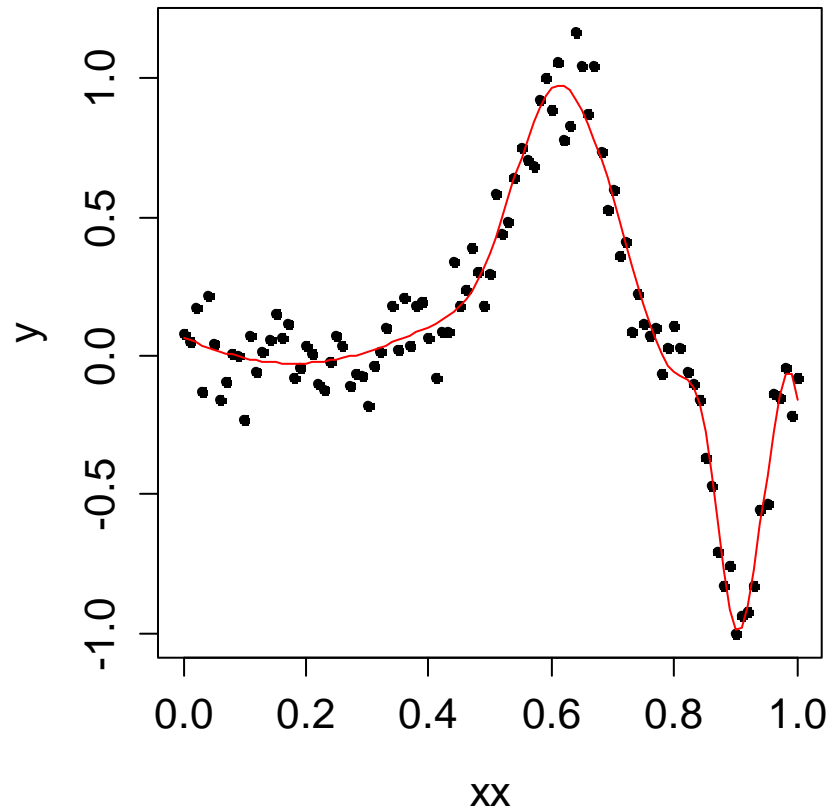
AS 2015 – Generalized Additive Modelling

Spline Basis and Resulting Fit

Kubische B-Splines



Fit mit Spline-Basis



Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

GAM Using a Spline Basis

In practice, we will rarely be satisfied with simple models, but require fitting multiple predictor GAMs. The idea is as follows:

$$y_i = \beta_0 + f_1(x_{i1}) + f_2(x_{i2}) + \dots + f_p(x_{ip}) + E_i$$

The principle is that for each predictor x_j , we will have a flexible and exploratively determined contribution $f_j(\cdot)$ that is rooted on a basis consisting of cubic B-splines with correct complexity. There is an excellent implementation in R...

→ How can this model be estimated?

→ How can one determine the correct smoothness of $f_j(\cdot)$?

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Backfitting-Algorithm

There is no single step solution to a multiple predictor GAM. We pursue an iterative approach that is based on stepwise solution of 1-dimensional problems:

- 1) Initialize $\hat{\beta}_0 = \bar{y}$ and $f_j(\cdot) = 0$ for all $j = 1, \dots, p$
- 2) Repeat for all $j = 1, \dots, p$ until convergence:
 - Compute $\tilde{y}_i = y_i - \hat{\beta}_0 - \sum_{k \neq j} f_k(x_{ik})$
 - Solve the 1-dimensional problem for $\hat{f}_j(\cdot)$ on (x_{ij}, \tilde{y}_i)
 - Center $\hat{f}_j(\cdot) \leftarrow \hat{f}_j(\cdot) - n^{-1} \sum_i \hat{f}_j(x_{ij})$

Note: the solution will only be identifiable if $\sum_i \hat{f}_j(x_{ij}) = 0$

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Implementation in `library(mgcv)`

- The backfitting-algorithm and in particular R function `gam()` also allows for having parametric terms in the model.
- The estimation in R package `mgcv` is not based on the backfitting algorithm specified above, but on the more sophisticated Lanczos approach (w/o details here...).
- **Syntax:** `fit <- gam(resp ~ s(p1) + s(p2) + p3, data=ex)`
- The complexity of the spline basis for each component will be estimated exploratively using cross validation. It may be overruled by typing `s(p1, df=...)`.

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Example: Prestige Data

```
> fit <- gam(prestige ~ s(income) + s(education), data=...)
> summary(fit)
Family: gaussian; Link function: identity
Formula: prestige ~ s(income) + s(education)
---
Parametric coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  46.8333    0.6889    67.98  <2e-16 ***
---
Approximate significance of smooth terms:
              edf Ref.df      F  p-value
s(income)     3.118  3.877 15.29 8.94e-10 ***
s(education)  3.177  3.952 38.34 < 2e-16 ***
---
R-sq.(adj) =  0.836    Deviance explained = 84.7%
GCV = 52.143  Scale est. = 48.414      n = 102
```

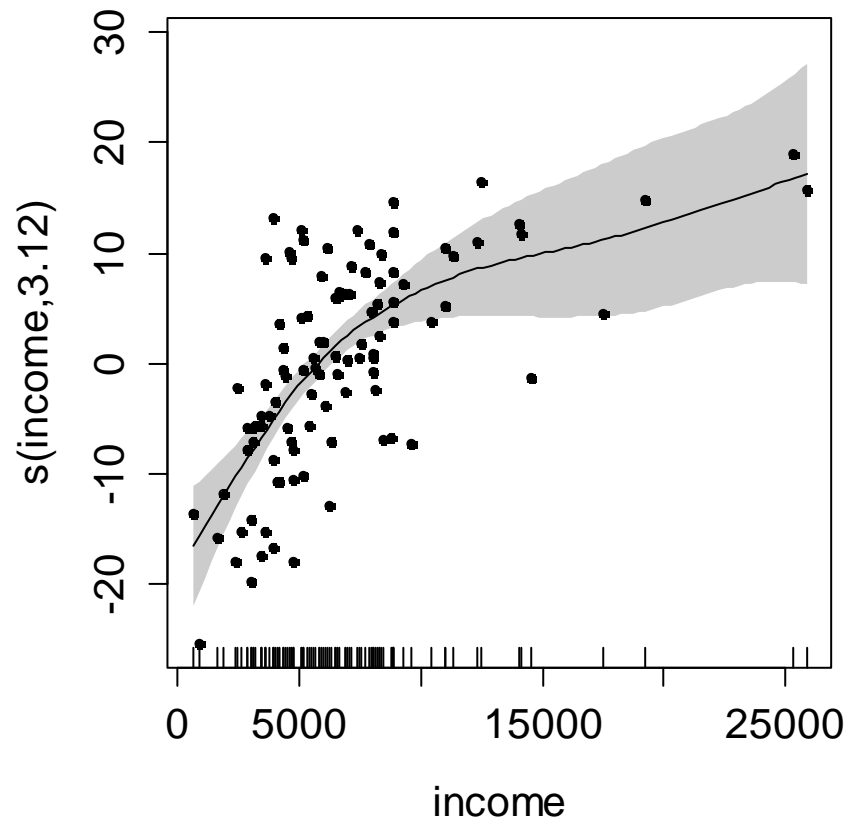
Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

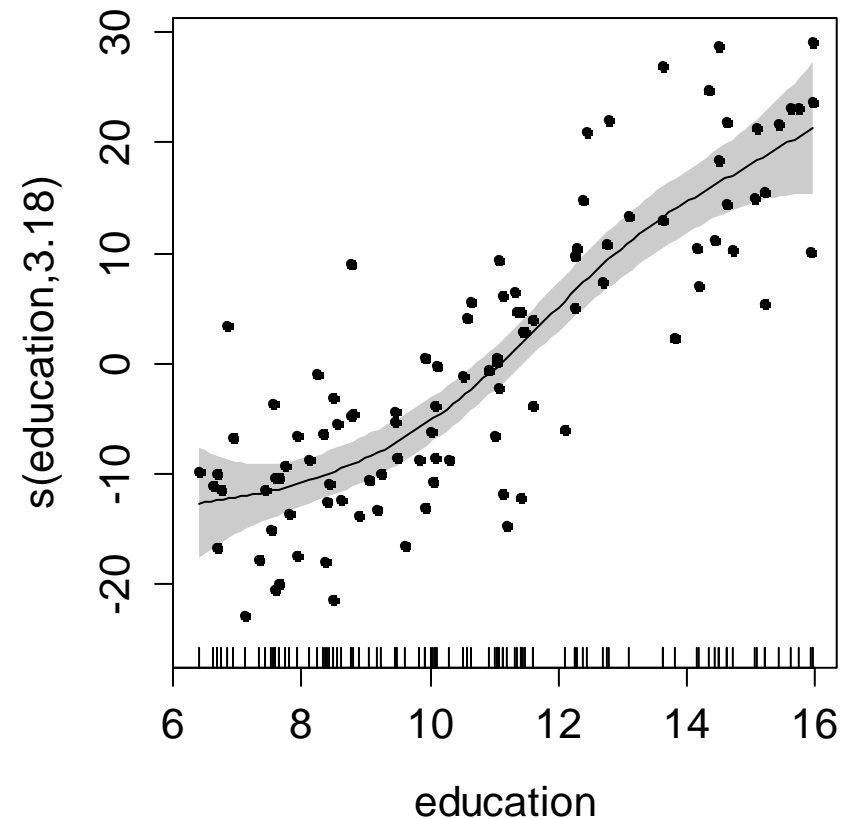
Example: Partial Residual Plots

```
> plot(fit, shade=TRUE, residuals=TRUE, pch=20, main=...)
```

GAM Partial Residual Plot



GAM Partial Residual Plot



Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Example: Residual Analysis

```
> gam.check(fit, pch=20, rep=100)
```

```
Method: GCV      Optimizer: magic
```

```
Smoothing parameter selection converged after 4 iterations  
The RMS GCV score gradient at convergence was 9.783945e-05  
The Hessian was positive definite.
```

```
The estimated model rank was 19 (maximum possible: 19)  
Model rank = 19 / 19
```

```
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'.
```

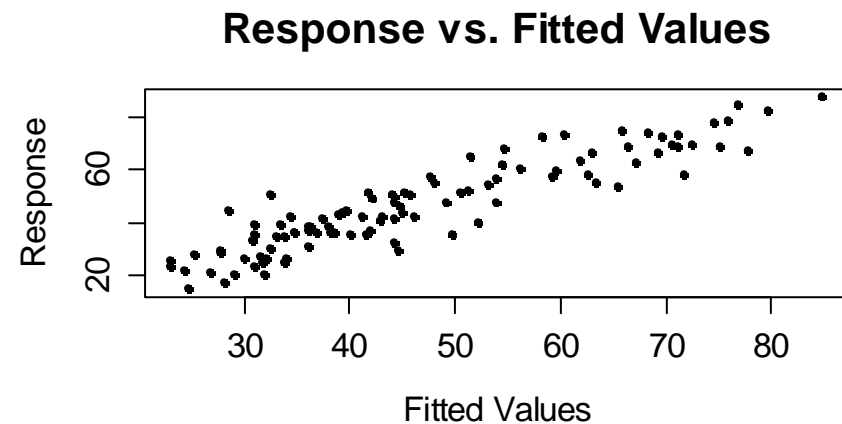
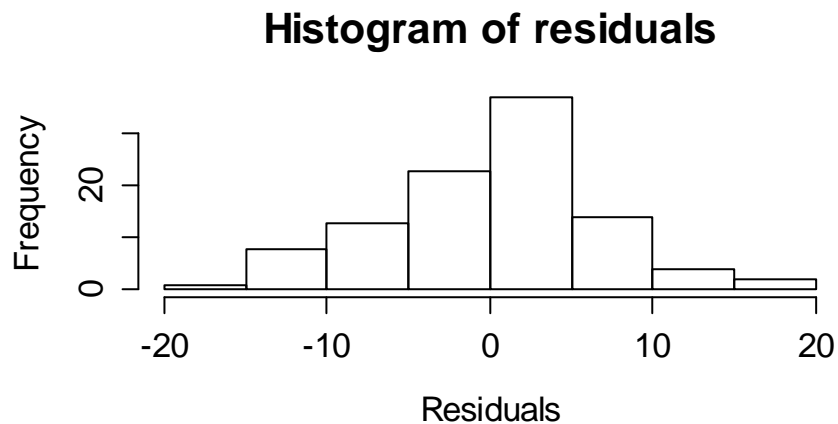
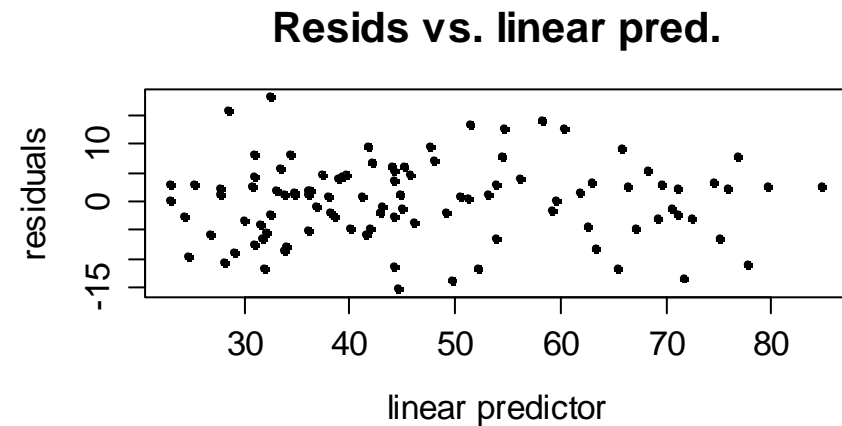
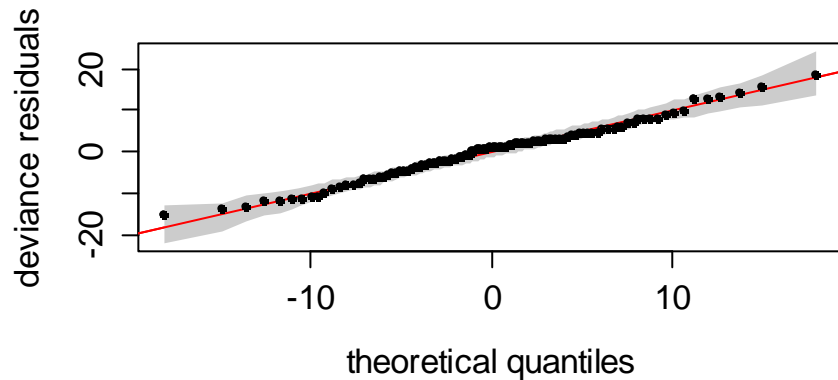
	k'	edf	k-index	p-value
s(income)	9.000	3.118	0.981	0.36
s(education)	9.000	3.177	1.025	0.61

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Example: Residual Analysis

```
> gam.check(fit, pch=20, rep=100)
```



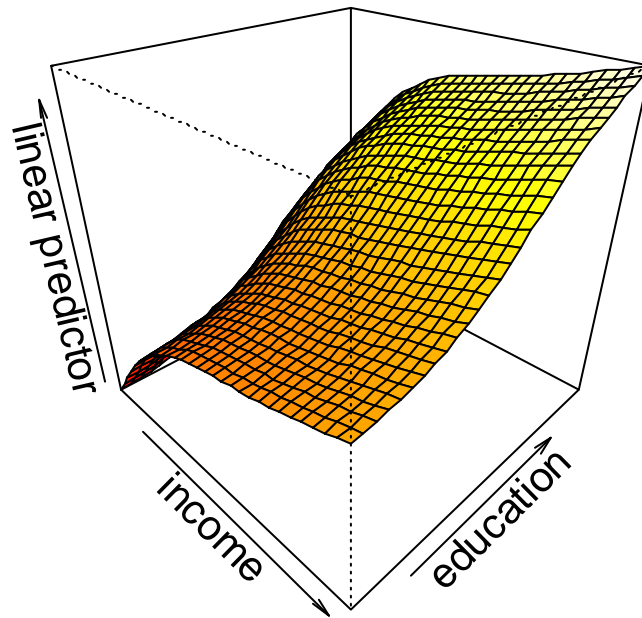
Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Example: Visualizing the Fit

```
> vis.gam(fit, theta=45, phi=30)
```

2-Dimensional Fit Visualization



Note: both predictors contribute in a non-linear fashion, but model is additive!

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Testen for Linearity

Function `gam()` determines the degrees of freedom for each of the predictors data-adaptively. If no flexibility is required, we can obtain $df=1$. In that case, the predictor contributes linearly.

However, in many situations one may be interested in formally testing whether a GAM yields a better fit than using OLS. This can be done on the basis of a test that gauges RSS versus the degrees of freedom of the respective models.

```
> fit
Estimated degrees of freedom:
3.12 3.18 total = 7.3
```

The GAM for the Prestige data spends 7.3 degrees of freedom. The competing OLS model only takes 3 of them!!!

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Testen for Linearity

```
> fit.ols <- gam(prestige ~ income + education, data=...)
```

```
Family: gaussian; Link function: identity
```

```
Formula: prestige ~ income + education
```

```
Total model degrees of freedom 3
```

```
GCV score: 62.84693
```

```
> deviance(fit.ols)
```

```
[1] 6038.851
```

```
> dd <- deviance(fit.ols)-deviance(fit); dd
```

```
[1] 1453.856
```

```
> 1-pchisq(dd, 7.3-3)
```

```
[1] 0
```

The GAM has a highly significant edge on OLS. However, we need to use variable transformations in the OLS model.

Applied Statistical Regression

AS 2015 – Generalized Additive Modelling

Testen for Linearity

There is some alternative (better) functionality that carries out the test for linearity as a one-line-command:

```
> anova(fit.ols, fit, test="Chisq")
Analysis of Deviance Table

Model 1: prestige ~ income + education
Model 2: prestige ~ s(income) + s(education)
  Resid. Df Resid. Dev      Df Deviance Pr(>Chi)
1     99.000     6038.9
2     94.705     4585.0  4.2951   1453.9 6.783e-06 ***
```

As we can see, the computed value for the test statistic is identical to the one on the previous slide. There is some rounding-based difference in the p-value, though.

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Non-Numerical Response Variable

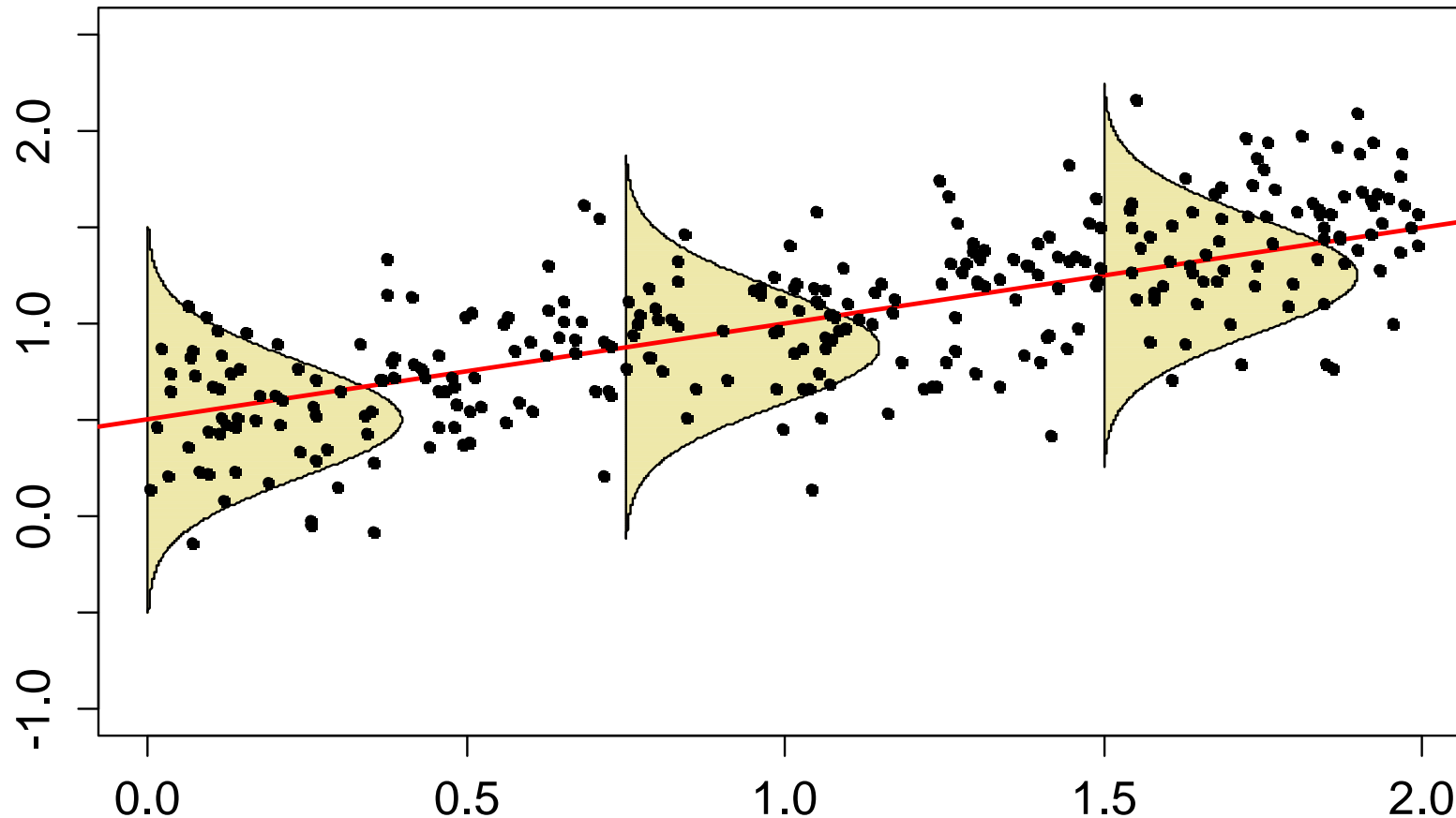
- **So far, the response y_i was a continuous random variable with infinite range, where the conditional distribution was a Gaussian, i.e. $y_i | X_i \sim N(\hat{y}_i, \sigma_E^2)$, see next slide.**
- If the task is modeling binary, binary or multinomial response (i.e. probabilities or proportions) or a count, this is not doable correctly with the methods that were discussed yet.
- We will present some additional techniques which implement linear modeling for these different types of responses. As we will see, there is a generic framework that incorporates all of these, as well as multiple linear regression.

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Conditional Gaussian Distribution

Linear Regression with Gaussian Distributions



Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Binary Response / Logistic Regression

What is the question?

In toxicological studies, one tries to infer whether a lab mouse survives when it is given a particular dose of poisonous matter.

In human medicine, one is often interested in the question how much of a drug is required to see an effect, i.e. pain reduction.

Mathematics:

→ The response variable $y_i \in \{0,1\}$ is binary

→ The conditional distribution $y_i | X_i \sim \text{Bernoulli}(p_i)$

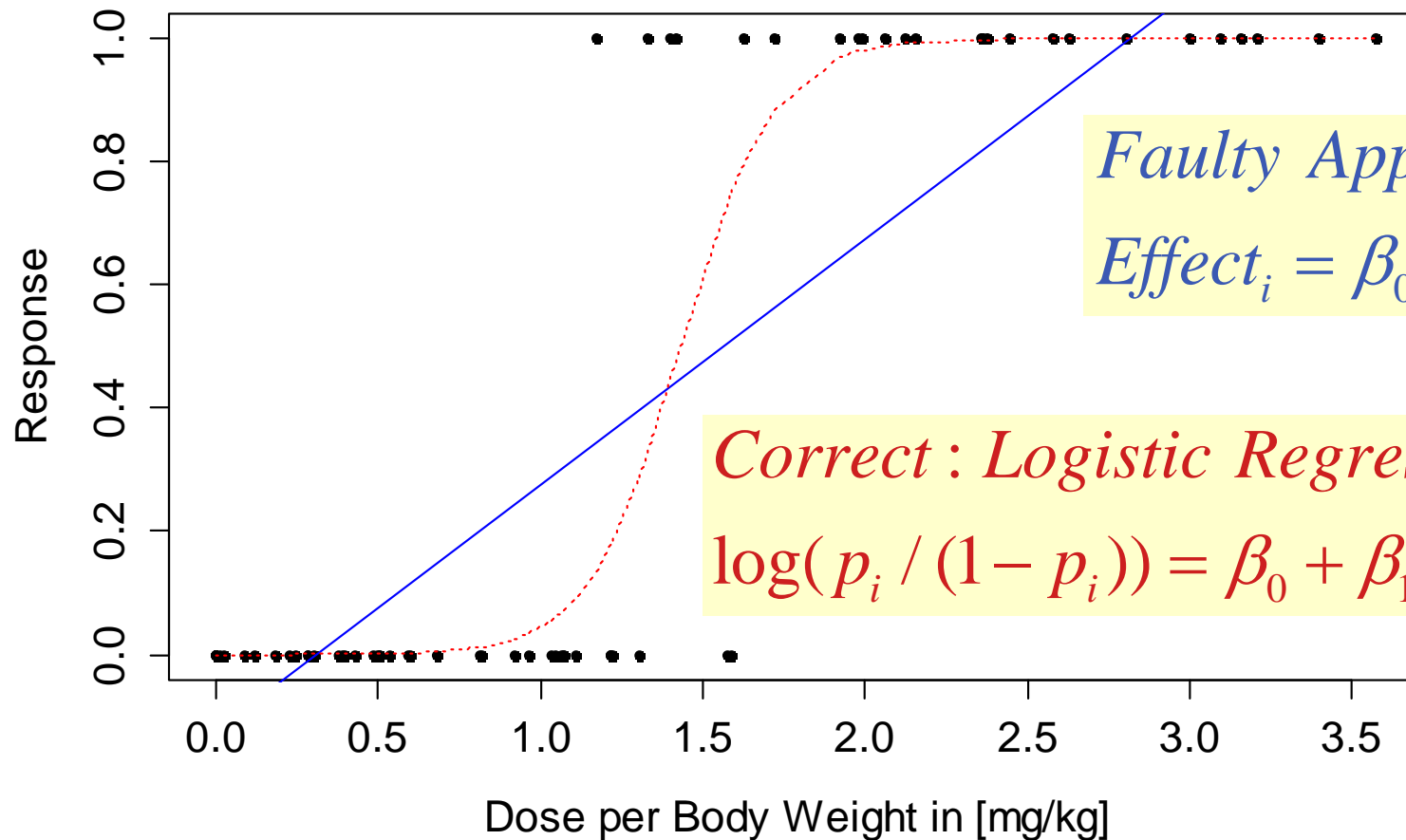
→ The fitted value is the expectation of the above conditional distribution, and hence the probability of death/survival p_i .

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Binary Response / Logistic Regression

Effect of Medication vs. Dose



Faulty Approach:
 $Effect_i = \beta_0 + \beta_1 \cdot Dose_i + E_i$

Correct: Logistic Regression
 $\log(p_i / (1 - p_i)) = \beta_0 + \beta_1 \cdot Dose_i$

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Count Response / Poisson Regression

What are predictors for the locations of starfish?

- analyze the number of starfish at several locations, for which we also have some covariates such as water temperature, ...
- the response variable is a count. The simplest model for this assumes a Poisson as the conditional distribution.

We assume that the logged parameter λ_i at location i depends in a linear way on the covariates:

$$y_i | X_i \sim \text{Pois}(\lambda_i), \text{ where } \log(\lambda_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Generalized Linear Models

What is it?

- General framework for regression type modeling
- Many different response types are allowed
- Notion: the responses' conditional expectation has a monotone relation to a linear combination of the predictors.

$$g(E[y_i | X_i]) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

- Some further requirements on variance and density of y
- **may seem complicated, but is very powerful!**

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Binary Response / Logistic Regression

The essence of a logistic regression model is that the response $y_i \in \{0, 1\}$, the conditional distribution is $y_i | X_i \sim \text{Bernoulli}(p_i)$ and we model the conditional expectation $E[y_i | X_i] = p_i$.

What do we need to take care of?

- Formulation of the model
- Estimation
- Inference
- Model diagnostics
- Model choice

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Example

Premature Birth, by Hubbard (1986)

$y_i \in \{0,1\}$ death (0) or survival (1) after premature birth.

Predictors:

- weight (in grams) at birth
- age at birth (in weeks of pregnancy)
- apgar scores (vital function after 1 and 5 min)
- pH-value of the blood (breathing)

Observations:

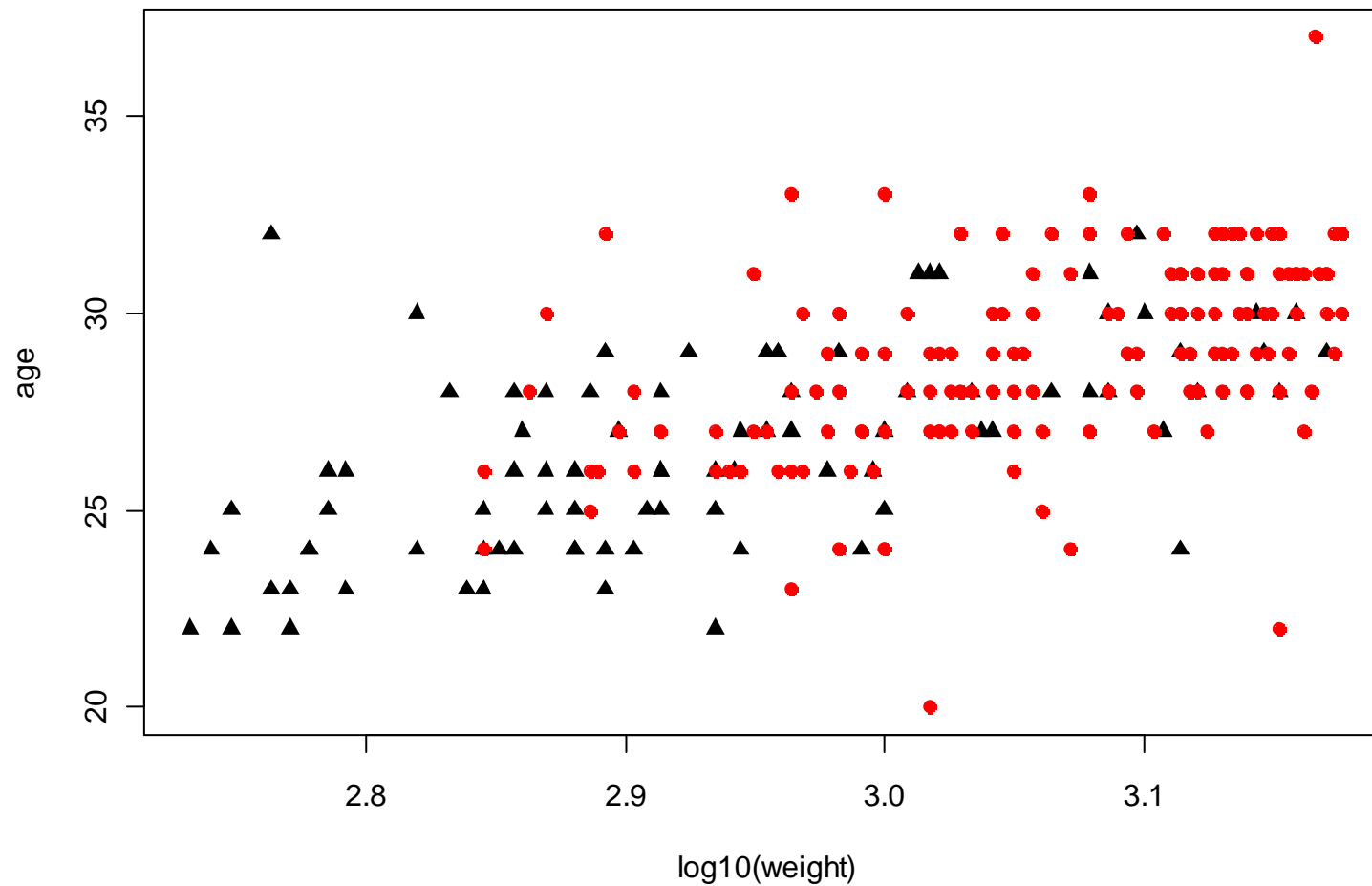
- there are 247 instances

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Example

Survival in Premature Birth



Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Logistic Regression Model

- $y_i \in \{0,1\}$ has a Bernoulli distribution.
- The parameter of this distribution is p_i , the success rate

Now please note that:

$$p_i = P(y_i = 1 | X_i) = E[y_i | X_i]$$

→ the most powerful notion of the logistic regression model is to see it as a model where we try to find a relation between the conditional expected value of y_i and the predictors!

Important: $P(y_i = 1 | X_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + E_i$ is no good here!

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Logit Transformation

Goal: mapping from $[0, 1] \mapsto (-\infty, +\infty)$

Logit transformation: $g(p) = \log\left(\frac{p}{1-p}\right)$

Interpretation: probabilities are mapped to logged odds ("Wettverhältnisse") which can then be modeled linearly.

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} \Leftrightarrow p_i = \frac{\exp(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip})}{1 + \exp(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip})}$$

→ Logistic regression = describing log-odds with a linear model

→ **Can you explain why there is no error term?**

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Some Remarks and Terminology

- For estimating the regression coefficients, we require the observations to be independent.
- There is no restriction for the predictors. They can be continuous, categorical, transformed, interactions, ...
- $\eta_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$ is called the linear predictor
- $g(\cdot)$ is the link function, mapping from $E[y_i | X_i]$ to η_i
- **There are other (less important) link functions:**
 - *probit link*
 - *c-log-log link*

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Estimation

Simple approach: minimize $\sum_{i=1}^n (y_i - p_i)^2$

→ Not a good way to estimate logistic regression parameters

Maximum-Likelihood approach:

General principle: determine the regression coefficients β_j such that the likelihood of the observed data is maximized. If the cases are independent, this amounts to maximizing the log-likelihood:

$$l(\beta) = \sum_{i=1}^n (y_i \log(p_i) + (1 - y_i) \log(1 - p_i)) \quad \text{with} \quad p_i = \frac{\exp(\beta_0 + \beta_1 x_{i1} + \dots)}{1 + \exp(\beta_0 + \beta_1 x_{i1} + \dots)}$$

Under mild conditions, the solution exists, but it cannot be written in closed form. Usually, the IRLS algorithm is employed.

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Estimation in R

```
> glm(survival~log10(weight)+age,  
      family=binomial, data=baby)
```

```
Coefficients: (Intercept)    log10(weight)      age  
            -33.9711          10.1685      0.1474
```

These are the estimates for $\beta_0, \beta_1, \beta_2$. Please note that they come from a numerical optimization, thus don't ignore this lightly:

Warning message:

```
glm.fit: algorithm did not converge
```

In this case, the coefficients are not trustworthy! However, this rarely happens in well posed regression problems.

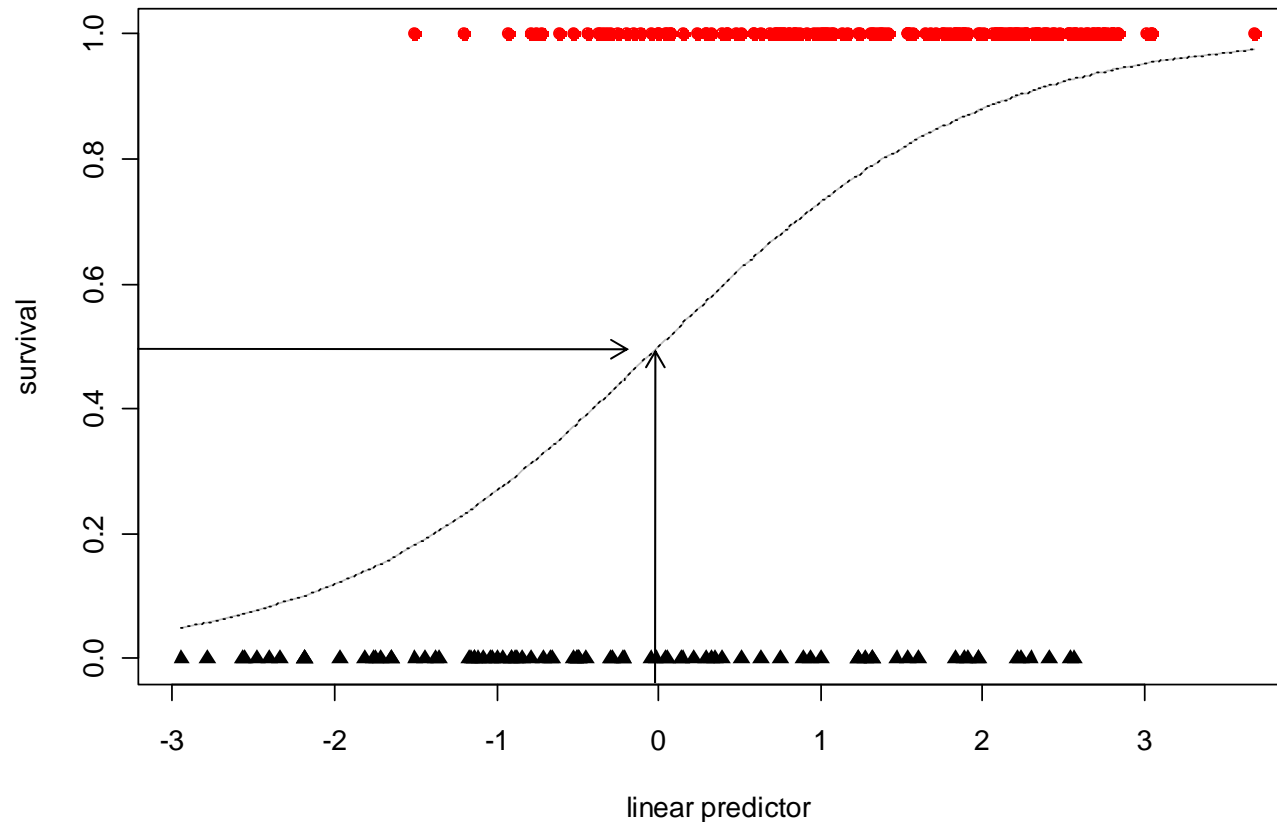
Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Displaying the Fit

$$g\left(P(y = 1 | \log_{10}(\text{weight}), \text{age})\right) = -33.97 + 10.17 \cdot \log_{10}(\text{weight}) + 0.14 \cdot \text{age}$$

Survival vs. Linear Predictor

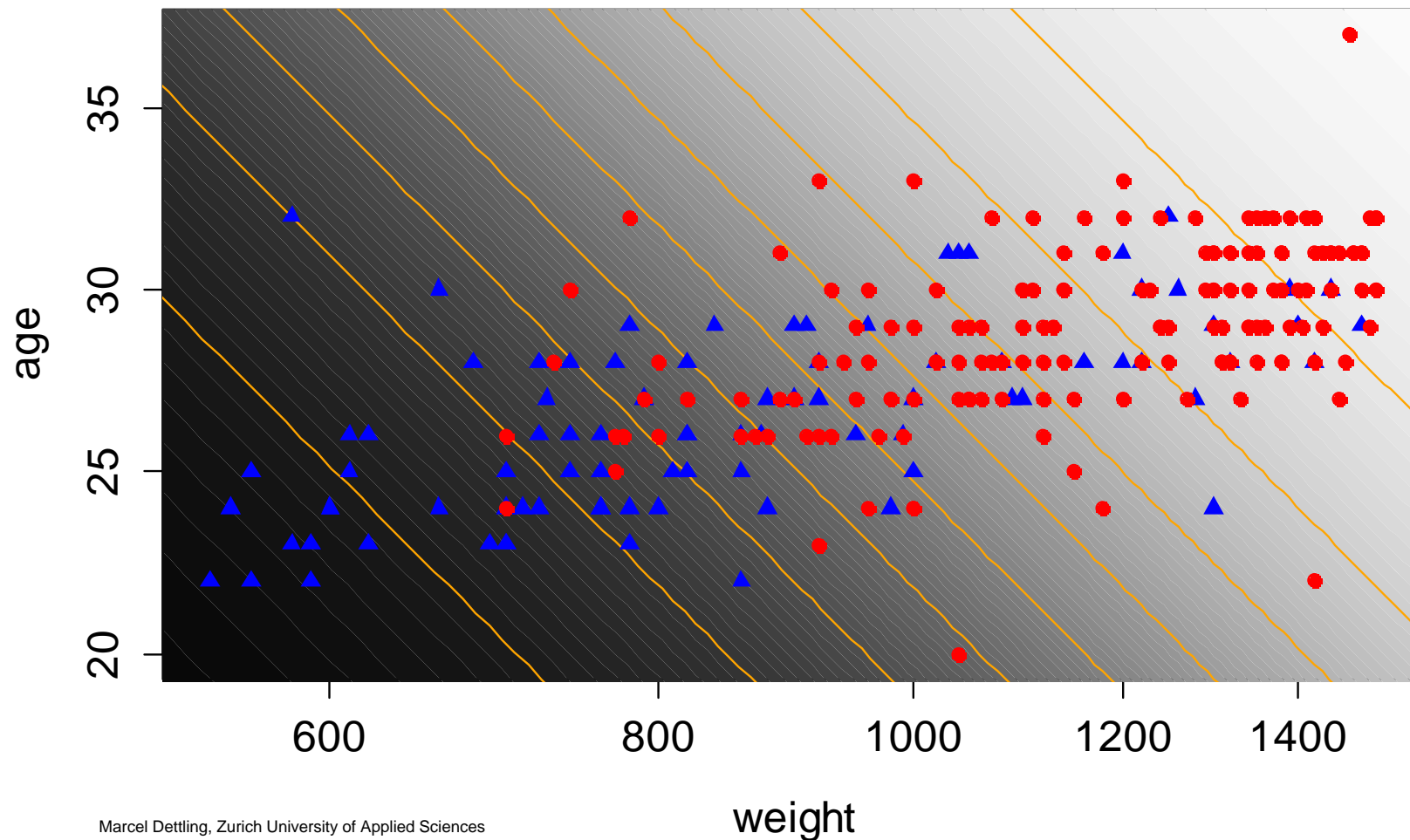


Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Displaying the Fit

Survival after Premature Birth



Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Interpretation of the Coefficients

→ see blackboard...

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Inference for Logistic Regression

While the basic concepts of inference will be familiar from multiple linear regression, various aspects will be markedly different.

Most importantly, the concept for the goodness-of-fit measure needs a second thought. We cannot work with the residuals sum of squares anymore, but employ the so-called **Residual Deviance**:

$$D(y, \hat{p}) = -2 \cdot \sum_{i=1}^n (y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i))$$

Also of importance is the **Null Deviance**, which is the deviance of the simplest possible model that is built from the intercept term only. It is always lower than the Residual Deviance.

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Summary Output from R

```
> summary(glm(survival ~ log10(weight) + age, data=baby,  
              family=binomial(link="logit")))
```

Deviance Residuals: ...

Coefficients:	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	-33.97108	4.98983	-6.808	9.89e-12	***
I(log10(weight))	10.16846	1.88160	5.404	6.51e-08	***
age	0.14742	0.07427	1.985	0.0472	*

Dispersion parameter for binomial family taken to be 1

Null deviance: 319.28 on 246 degrees of freedom

Residual deviance: 235.94 on 244 degrees of freedom

Number of Fisher Scoring iterations: 4

AIC: 241.94

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Coefficient of Determination

There is no direct analogon to the coefficient of determination in logistic regression. Some suggestions for the COD include:

Proportion of deviance explained

```
> 1-fit$dev/fit$null  
[1] 0.2610193
```

A better statistic for measuring the explanatory content

$$R^2 = \frac{1 - \exp((D_{res} - D_{null}) / n)}{1 - \exp(-D_{null} / n)} = 0.395$$

→ There are even more suggestions in the literature.

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Inference: Individual Parameter Tests

Multiple Linear Regression:

Gaussian errors $\rightarrow \hat{\beta}_j$ are normally distributed

Logistic Regression:

There are no errors, variability arises from Bernoulli distribution

MLE-theory tells us that under mild conditions, the coefficients $\hat{\beta}_j$ are approximately normally distributed with a covariance matrix V that can be derived from the coefficients.

Hence: If $H_0 : \beta_j = b$ then use $Z = \frac{\hat{\beta}_j - b}{\hat{\sigma}_{\hat{\beta}_j}} \sim N(0,1)$

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Inference: 95%-CI for the Coefficients β_j

It is straightforward to generate a hand-constructed 95%-CI:

$$\hat{\beta}_j \pm 2 \cdot se(\hat{\beta}_j) = \hat{\beta}_j \pm 2 \cdot \hat{\sigma}_{\hat{\beta}_j}$$

Rather than using the approximate 2 for the 97.5%-quantile of the Gaussian distribution, we can use the exact values. For $\log_{10}(\text{weight})$, we so obtain:

```
> 10.16846 + qnorm(c(0.025, 0.975)) * 1.88160  
[1] 6.480592 13.856328
```

However, in R it is more convenient to use `confint()` which here uses a slightly different, more sophisticated and exact computation by interpolation of the likelihood profile traces.

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Inference: Comparing Hierarchical Models

Analogon to the partial F-test in multiple linear regression

Big Model: has $(p + 1)$ coefficients $\beta_0, \beta_1, \dots, \beta_q, \dots, \beta_p$

Small Model: has $(q + 1)$ coefficients $\beta_0, \beta_1, \dots, \beta_q$

Null hypothesis: $H_0 : \beta_{q+1} = \beta_{q+2} = \dots = \beta_p = 0$

MLE theory suggests to use the likelihood ratio or log-likelihood difference as a test statistic. This amounts to taking the difference of the residual deviances. It asymptotically follows a Chisquare distribution with $(p - q)$ degrees of freedom:

$$2\left(\ell^{Big} - \ell^{Small}\right) = D\left(y, \hat{p}_{Small}\right) - D\left(y, \hat{p}_{Big}\right) \sim \chi_{(p-q)}^2$$

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Example: Global Test

Idea: Compare the actual model against the simplest possible model with only the intercept. The latter fits the overall success ratio $\hat{p}_{Null} = \sum y_i / n$ to all observations.

Since our actual model and the null model are nested, we can perform a hierarchical model comparison. In the baby survival example, there are two predictors and hence:

$$D(y, \hat{p}_{Null}) - D(y, \hat{p}_{Big}) \sim \chi_2^2$$

The two deviances are reported in the summary output.

Null deviance: 319.28 on 246 degrees of freedom

Residual deviance: 235.94 on 244 degrees of freedom

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Example: Global Test

A quick and simple check for rejection of $H_0 : \beta_1 = \beta_2 = 0$ is to compare the difference in deviance vs. the difference in degrees of freedom.

If $D(y, \hat{p}_{Null}) - D(y, \hat{p}_{Big}) \gg (p - q)$ then reject H_0

The exact p-valued can be computed in R by:

```
> 1-pchisq(fit$null-fit$dev, df=(fit$df.null-fit$df.res))  
[1] 0
```

The p-value is (numerically) zero, hence the null hypothesis is very clearly rejected. Conjecture: there is a strongly significant contribution of $\log_{10}(\text{weight})$ and age to the odds for survival.

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Using `drop1()` *for Testing*

The R function `drop1()` performs hierarchical model comparison for exclusion of one model term at a time.

```
> drop1(fit, test="Chisq")
Single term deletions
Model: survival ~ I(log10(weight)) + age
      Df Deviance    AIC    LRT   Pr(Chi)
<none>      235.94 241.94
log10(weight)  1   270.19 274.19 34.247 4.855e-09 ***
age           1   239.89 243.89  3.948  0.04694  *
```

Question:

- where is the difference to the summary output?
- it exists, though it's not obvious and asymptotically vanishes

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Model Diagnostics

Diagnostics are:

- in principle as important with logistic regression as they are with multiple linear regression models, but more difficult.
 - again based on differences between fitted & observed values
- we have to take into account that the variance of the response residuals $r_i = y_i - \hat{p}_i$ is non-constant.
- we have to come up with novel types of residuals:

Pearson and Deviance residuals

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Pearson Residuals

Take response residual (difference between observed and fitted value) and divide by an estimate of its standard deviation:

$$R_i = \frac{y_i - \hat{p}_i}{\sqrt{\hat{p}_i(1 - \hat{p}_i)}}$$

- R_i^2 is the contribution of the i th observation to the Pearson statistic for model comparison (that we did not discuss).
- It is important to note that Pearson residuals exceeding a value of two in absolute value warrant a closer look. They appear if an observation with $\hat{p}_i = 0.8$ resp. $\hat{p}_i = 0.2$ in reality has class label 0 resp. 1.

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Deviance Residuals

Take the contribution of the i th observation to the log-likelihood, i.e. the chi-square statistic for model comparison.

$$d_i = (y_i \cdot \log(\hat{p}_i) + (1 - y_i) \cdot \log(1 - \hat{p}_i))$$

For obtaining a well interpretable residual, we take the square root and the sign of the difference between true and fitted value:

$$D_i = \text{sign}(y_i - \hat{p}_i) \cdot \sqrt{d_i}$$

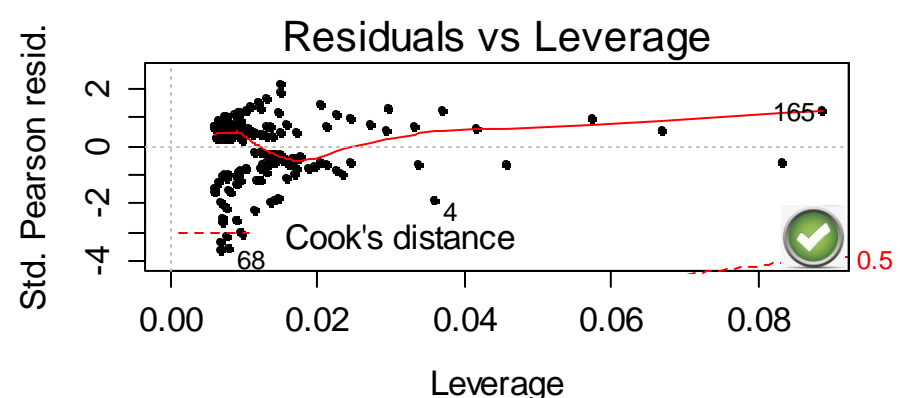
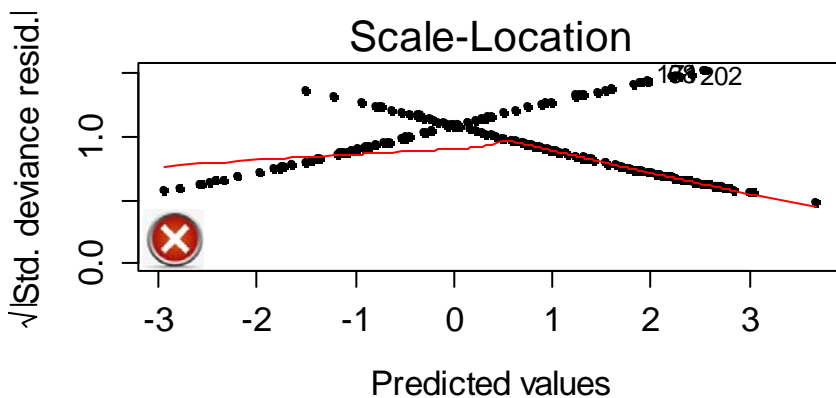
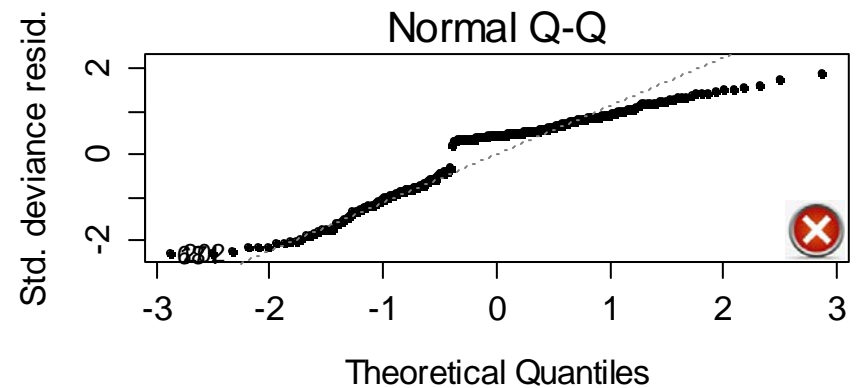
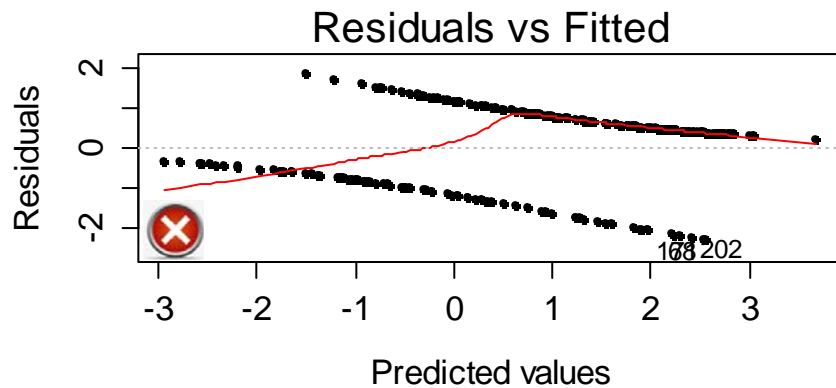
- - *deviance residuals > 2 warrant a closer look.*
- *the distribution of the deviance residuals is not known.*

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Model Diagnostics in R

The 4 standard plots are not well suited for logistic regression!!!

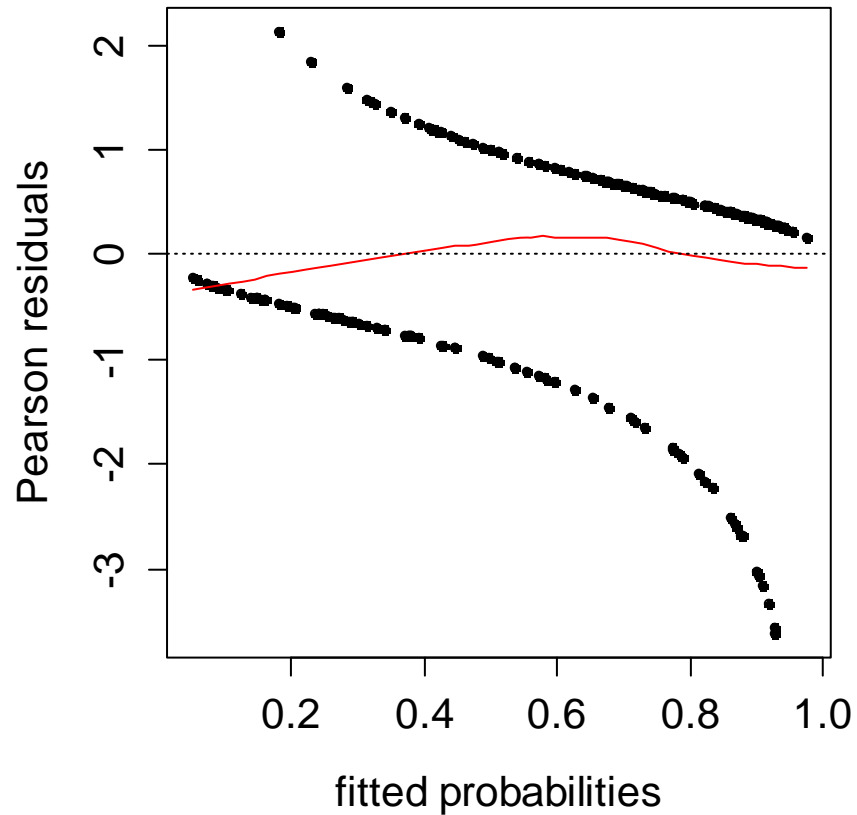


Applied Statistical Regression

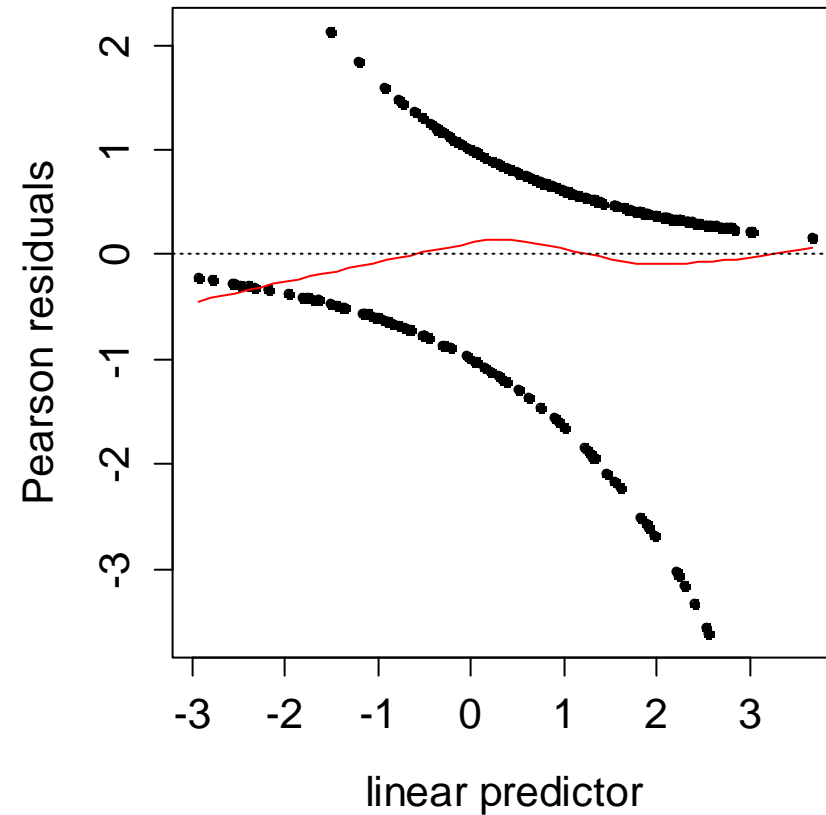
AS 2015 – Generalized Linear Modeling

Improved Tukey-Anscombe Plots

Tukey-Anscombe 1



Tukey-Anscombe 2



Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Improved Tukey-Anscombe Plot

It is mandatory to use a non-robust smoother in the TA plot!

```
xx <- predict(fit, type="response")
yy <- residuals(fit, type="pearson")
loess.smooth(xx, yy, family="gaussian", pch=20)
abline(h=0, lty=3)
```

Remarks:

- On the y-axis, use *Pearson* or *Deviance* residuals
- On the x-axis, use the *linear predictor* or *probabilities*
- One can, but does not have to use studentized residuals
- The LogReg residuals do not follow a Gaussian distribution
- The LogReg residuals always lie on two curves
- Residual analysis is easier with grouped data!

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

AIC and Variable Selection

General remark:

All comparison between models of different size can also be done using the AIC criterion. Not only in logistic regression, but also here.

The criterion:

$$AIC = D(y_i, \hat{p}) + 2p$$

Variable selection:

- stepwise approaches as with multiple linear regression
- factor variables need to be treated the right way!

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Binomial Regression Models

Concentration in log of mg/l	Number of insects n_i	Number of killed insects y_i
0.96	50	6
1.33	48	16
1.63	46	24
2.04	49	42
2.32	50	44

- For the number of killed insects, we have $y_i \sim \text{Bin}(n_i, p_i)$
- We are mainly interested in the proportion of insects surviving
- These are grouped data for which we do binomial regression.
We could run a logistic regression with 243 observations instead, but the grouped data approach is more powerful!

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Model and Estimation

The goal is to find a relation:

$$p_i = P(y_i = 1 | X_i) \sim \eta_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

We will again use the logit link function such that $\eta_i = g(p_i)$

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

Here, p_i is the expected value $E[y_i / n_i]$, and thus, also this model here fits within the GLM framework. The log-likelihood is:

$$l(\beta) = \sum_{i=1}^k \left[\log\binom{n_i}{y_i} + n_i y_i \log(p_i) + n_i (1 - y_i) \log(1 - p_i) \right]$$

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Fitting with R

We need to generate a two-column matrix where the first contains the “successes” and the second contains the “failures”

```
> killsurv
```

```
      killed surviv
[1,]      6     44
[2,]     16     32
[3,]     24     22
[4,]     42      7
[5,]     44      6
```

```
> fit <- glm(killsurv~conc, family="binomial")
```

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Summary Output

The result for the insecticide example is:

```
> summary(glm(killsurv ~ conc, family = "binomial"))
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	-4.8923	0.6426	-7.613	2.67e-14	***
conc	3.1088	0.3879	8.015	1.11e-15	***

Null deviance: 96.6881 on 4 degrees of freedom

Residual deviance: 1.4542 on 3 degrees of freedom

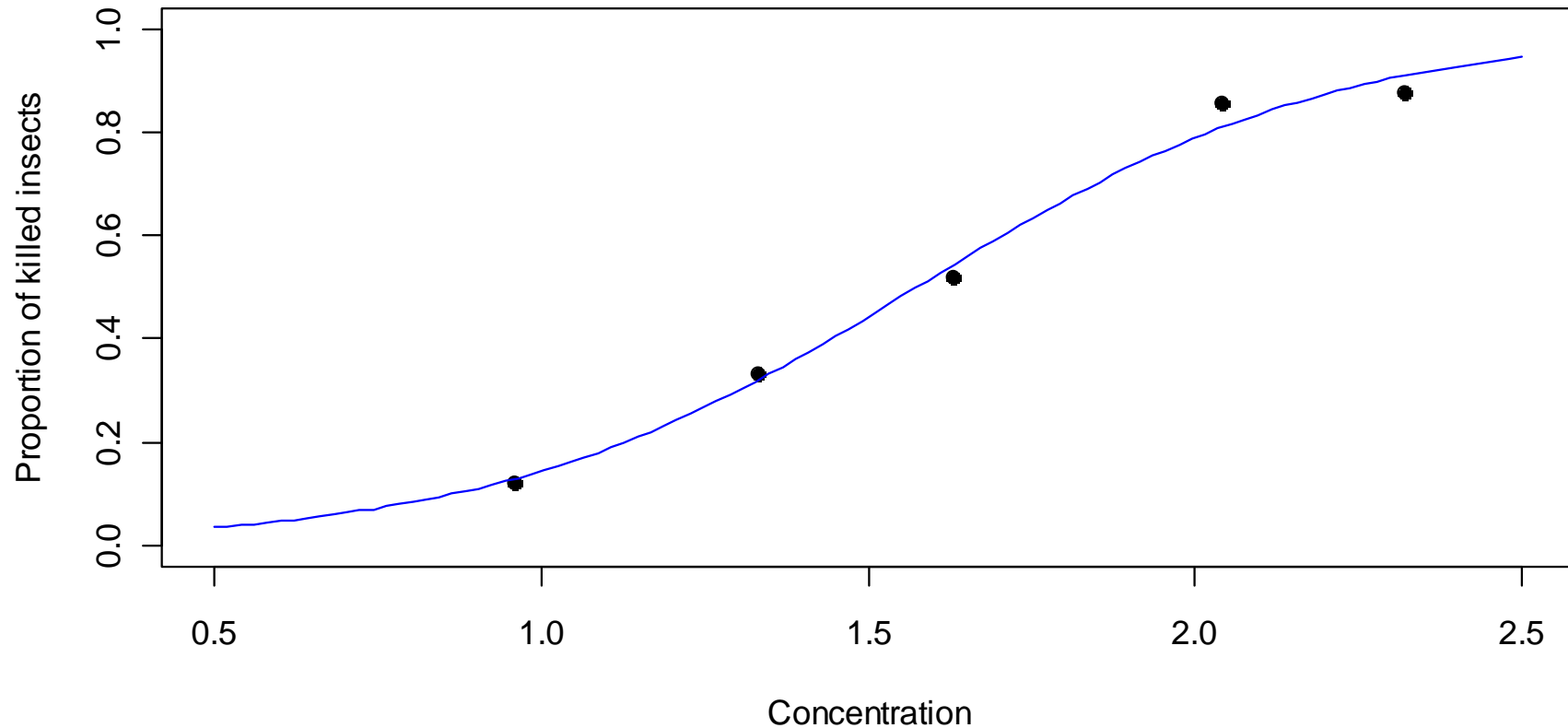
AIC: 24.675

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Proportion of Killed Insects

Insecticide: Proportion of Killed Insects



Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Global Tests for Binomial Regression

For GLMs there are three tests that can be done:

- **Goodness-of-fit test**
 - based on comparing against the saturated model
 - not suitable for non-grouped, binary data
- **Comparing two hierarchical models**
 - likelihood ratio test leads to deviance differences
 - test statistics has an asymptotic Chi-Square distribution
- **Global test**
 - comparing versus an empty model with only an intercept
 - this is a nested model, take the null deviance

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Goodness-of-Fit Test

→ the residual deviance will be our goodness-of-fit measure!

Paradigm: take twice the difference between the log-likelihood for our current model and the saturated one, which fits the proportions perfectly, i.e. $\hat{p}_i = y_i / n_i$

$$D(y, \hat{p}) = 2 \sum_{i=1}^k \left[y_i \log \left(\frac{y_i}{\hat{y}_i} \right) + (n_i - y_i) \log \left(\frac{(n_i - y_i)}{(n_i - \hat{y}_i)} \right) \right]$$

Because the saturated model fits as well as any model can fit, the deviance measures how close our model comes to perfection.

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Evaluation of the Test

Asymptotics:

If Y_i is truly binomial and the n_i are large, the residual deviance is approximately χ^2 distributed. The degrees of freedom is:

$$k - (\# \text{ of predictors}) - 1$$

```
> 1 - pchisq(deviance(fit), df.residual(fit))  
[1] 0.69287
```

Quick and dirty:

Deviance \gg *df* : \rightarrow model is not worth much.
More exactly: check $df \pm 2\sqrt{df}$

\rightarrow only apply this test if at least all $n_i \geq 5$

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Overdispersion

What if *Deviance* \gg *df* ???

1) Check the structural form of the model

- model diagnostics
- predictor transformations, interactions, ...

2) Outliers

- should be apparent from the diagnostic plots

3) IID assumption for p_i within a group

- unrecorded predictors or inhomogeneous population
- subjects influence other subjects under study

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Overdispersion: a Remedy

We can deal with overdispersion by estimating:

$$\hat{\phi} = \frac{X^2}{n-p} = \frac{1}{n-p} \cdot \sum_{i=1}^n \frac{(y_i - n_i \hat{p}_i)^2}{n_i \hat{p}_i (1 - \hat{p}_i)}$$

This is the sum of squared Pearson residuals divided with the df

Implications:

- regression coefficients remain unchanged
- standard errors will be different: inference!
- need to use an F-test for comparing nested models

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Results when Correcting Overdispersion

```
> phi <- sum(resid(fit)^2)/df.residual(fit)
> phi
[1] 0.4847485
> summary(fit, dispersion=phi)
              Estimate Std. Error z value Pr(>|z|)
(Intercept)  -4.8923      0.4474  -10.94  <2e-16 ***
conc          3.1088      0.2701   11.51  <2e-16 ***
---
(Dispersion parameter taken to be 0.4847485)
Null deviance: 96.6881  on 4  degrees of freedom
Residual deviance:  1.4542  on 3  degrees of freedom
AIC: 24.675
```

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Global Tests for Binomial Regression

For GLMs there are three tests that can be done:

- **Goodness-of-fit test**
 - based on comparing against the saturated model
 - not suitable for non-grouped, binary data
- **Comparing two hierarchical models**
 - likelihood ratio test leads to deviance differences
 - test statistics has an asymptotic Chi-Square distribution
- **Global test**
 - comparing versus an empty model with only an intercept
 - this is a nested model, take the null deviance

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Testing Hierarchical Models / Global Test

For binomial regression, these two tests are conceptually equal to the ones we already discussed in binary logistic regression.

→ *We refer to our discussion there and do not go into further detail here at this place!*

Null hypothesis and test statistic:

$$H_0 : \beta_{q+1} = \beta_{q+2} = \dots = \beta_p = 0$$

$$2\left(\ell^{(B)} - \ell^{(S)}\right) = D\left(y, \hat{p}^{(S)}\right) - D\left(y, \hat{p}^{(B)}\right)$$

Distribution of the test statistic:

$$D^{(S)} - D^{(B)} \sim \chi_{p-q}^2$$

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Practical Example

With this example taken from the lecturer's research, we illustrate the pro's and con's of working with logistic vs. binomial regression, i.e. grouped vs. non-grouped data

CHURN	REGION	GENDER	AGE	TENURE	PRODUCT
1	D-CH	male	65	84	PH + INET + TV
1	F-CH	female	45	34	INET + TV
1	F-CH	female	68	52	INET + TV
1	D-CH	female		102	INET
1	D-CH	male	45	21	TV
1	D-CH	male	43	63	PH + INET + TV
1	I-CH	male	28	47	TV

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Practical Example

Goal: understanding *churn*, i.e. end of contract

Model: $churn \sim region + gender + age + tenure + product$

The data per se are non-grouped, with millions of observations. But in this problem, it **pays off to work with grouped data**.

The main advantages when doing so are:

- Dealing with missing values in *age* and *tenure*: we do not lose any observations when factorizing these two variables.
- Instead of millions of rows, the design matrix is reduced to just 885 rows. This speeds up the computing tremendously.
- Much better inference and residual analysis is possible!

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Aggregating the Data in R

```
## Aggregating the data
> gdat <- aggregate(dat$churn,by=list(dat$region, dat$sex,
                                     dat$age.group, dat$dauer.group,
                                     dat$produkt),table)

## Excerpt of the data
> gdat[c(34, 92, 122, 588),]
  region sex age dauer produkt churn.no churn.yes
34   F-CH male Missing [0,24]   PHON         53         8
92   F-CH male (45,60] (72,180] PHON         50         6
122  F-CH female (30,45] [0,24]   TV          826        194
588  F-CH female (45,60] (72,180] INET+TV      103         14
```

→ Now, there are $3 \cdot 3 \cdot 6 \cdot 3 \cdot 7 = 1134$ groups, of which only 885 are populated. We will now fit a binomial regression model using only the main effects (i.e. without any interaction terms).

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Summary Output

```
> drop1(fit, test="Chisq")
```

```
Model: churn ~ region + sex + age + dauer + produkt
```

	Df	Deviance	AIC	LRT	Pr(>Chi)	
<none>		2866.6	6254.7			
region	2	3212.0	6596.1	345.4	< 2.2e-16	***
sex	2	3344.4	6728.5	477.8	< 2.2e-16	***
age	5	6745.2	10123.3	3878.6	< 2.2e-16	***
dauer	2	4172.9	7557.0	1306.3	< 2.2e-16	***
produkt	6	10718.3	14094.4	7851.7	< 2.2e-16	***

```
---
```

```
Null deviance: 19369.7 on 884 degrees of freedom
```

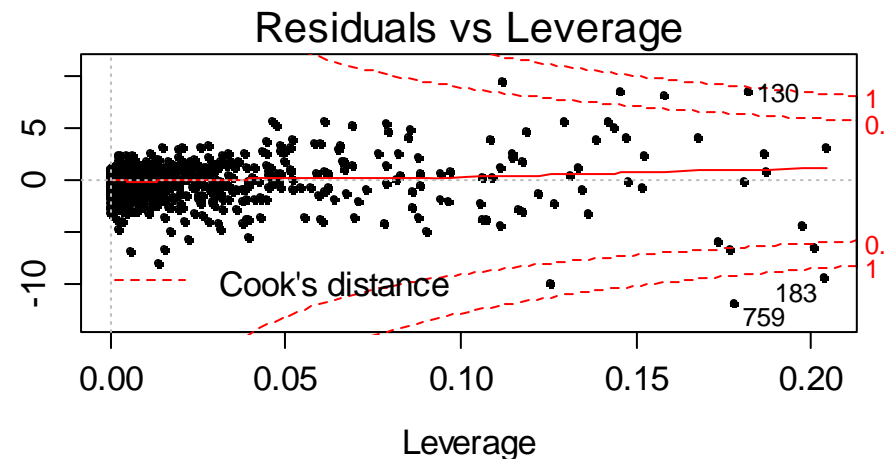
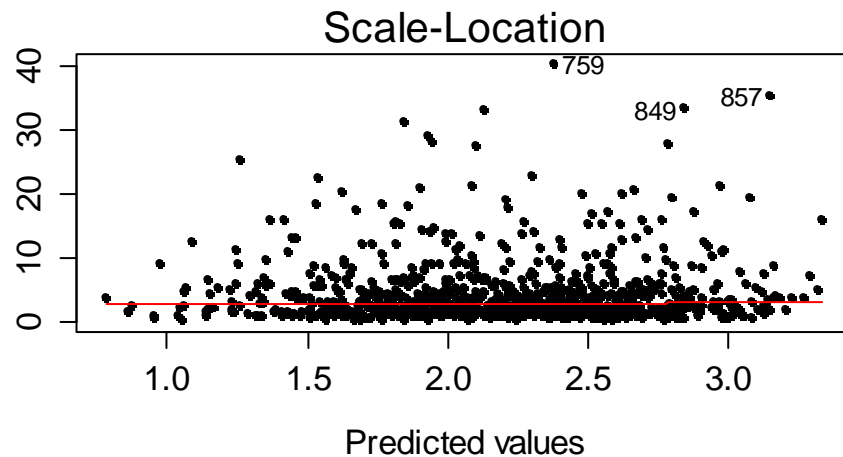
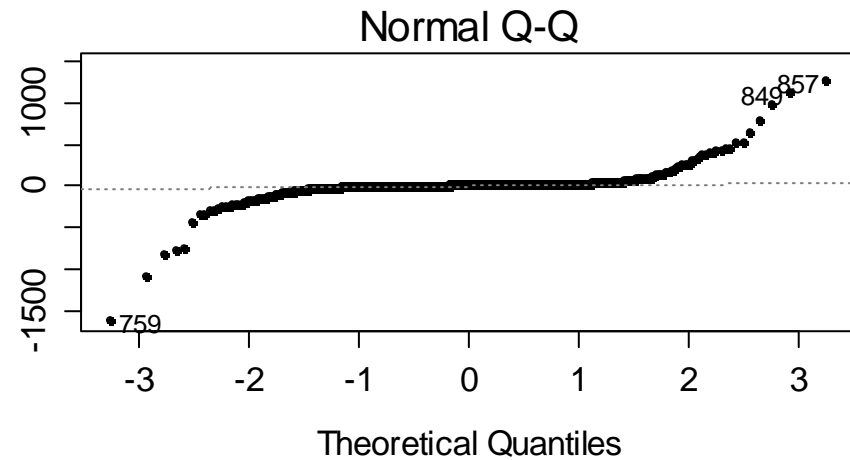
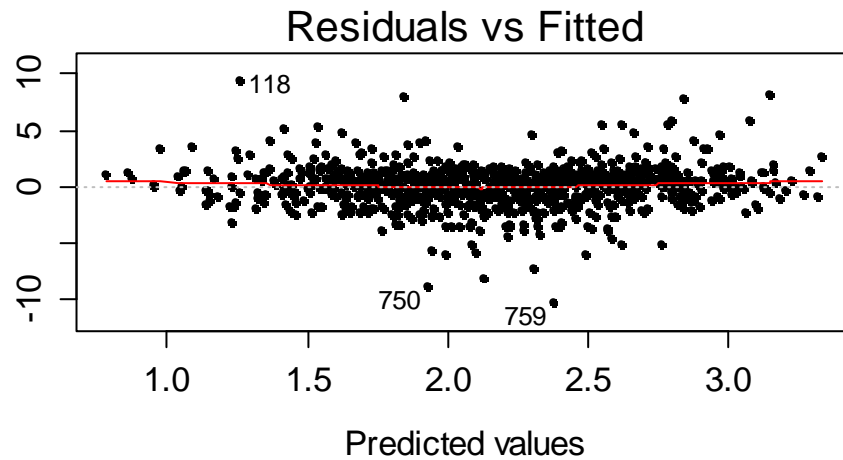
```
Residual deviance: 2866.6 on 867 degrees of freedom
```

→ Very strong overdispersion, the model does not fit well!

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

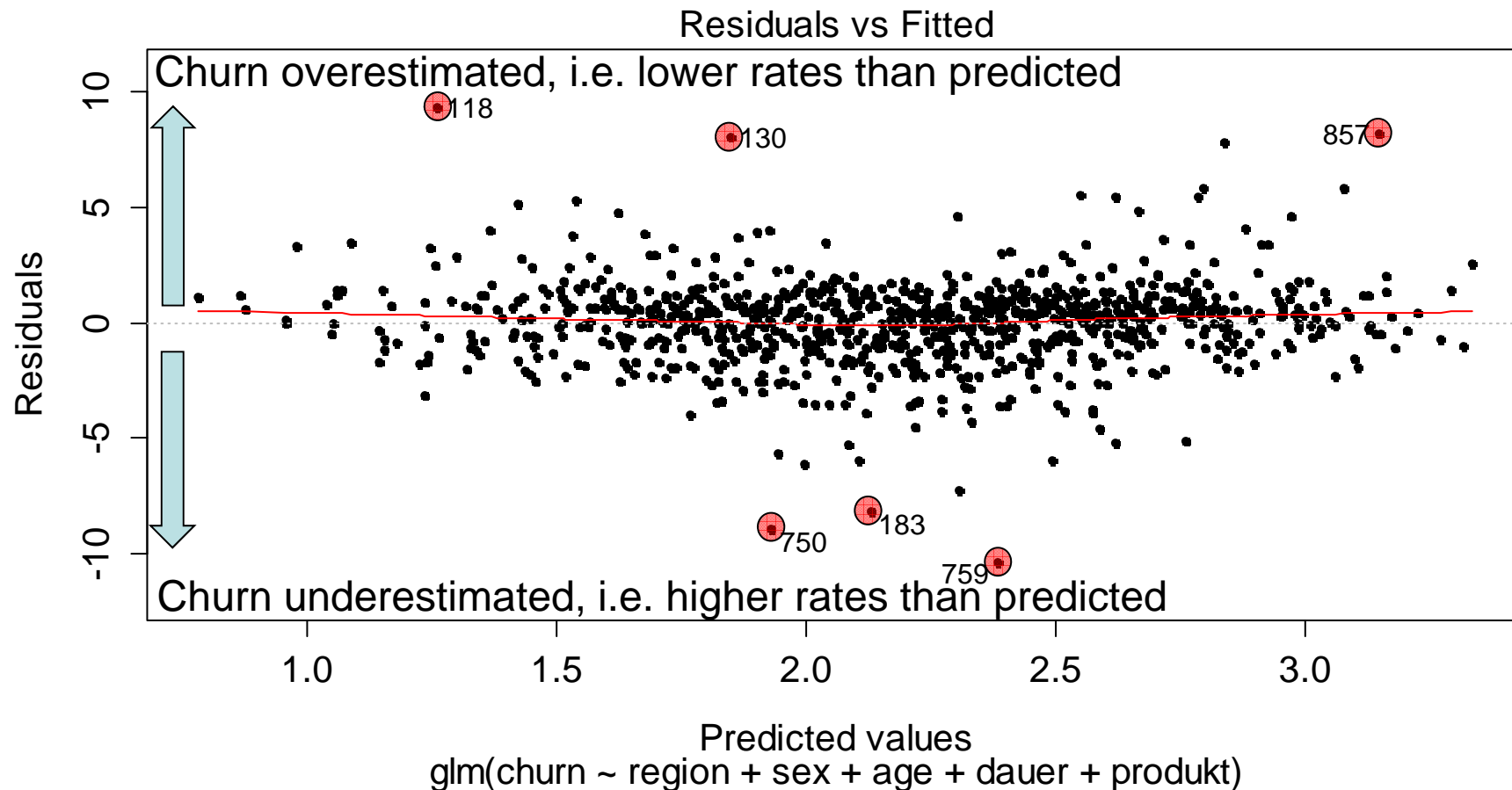
Model Diagnostics



Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Detail: Residuals vs. Predicted



Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Discussion of the Practical Example

The analysis of grouped data shows that we have a very incomplete understanding of the churn mechanics. There are groups for which the churn probability is very strongly over- or underestimated. All-in-all, the goodness-of-fit test for our binomial model is rejected.

What to do?

- Use more and/or better predictors for *churn*.
- If not available, try to work with interaction terms.
- Using a dispersion parameter doesn't help for prediction!
- Models can/should also be evaluated using cross validation.

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Poisson-Regression

When to apply?

- *Generally, if the **response variable is a count**. However:*
 - for bounded counts, the binomial model can be useful
 - for large numbers the normal approximation can serve
- *The use of Poisson regression is a must if:*
 - the counts are small and/or population size unknown
 - the population size is big and hard to come by, and the probability of an event, resp. the counts are small.

Model, Estimation, Inference:

Poisson Regression fits within the GLM framework!

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Example: Tortoise Species on Galapagos

The data are as follows:

```
> library(faraway); data(gala); head(gala[,-2])
```

	Species	Area	Elevation	Nearest	Scruz	Adjacent
Baltra	58	25.09	346	0.6	0.6	1.84
Bartolome	31	1.24	109	0.6	26.3	572.33
Caldwell	3	0.21	114	2.8	58.7	0.78
Champion	25	0.10	46	1.9	47.4	0.18
Coamano	2	0.05	77	1.9	1.9	903.82
Daphne.Major	18	0.34	119	8.0	8.0	1.84

Because the predictors all take positive values only and are skewed to the right, we urgently need transformations, namely:

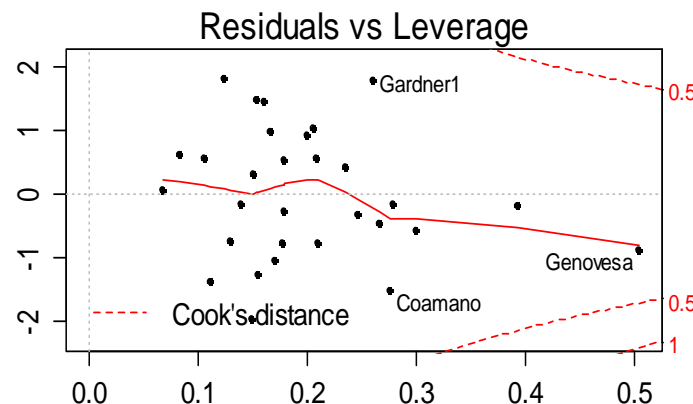
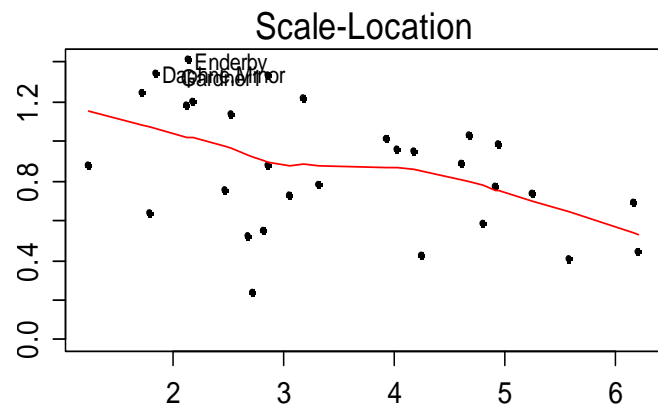
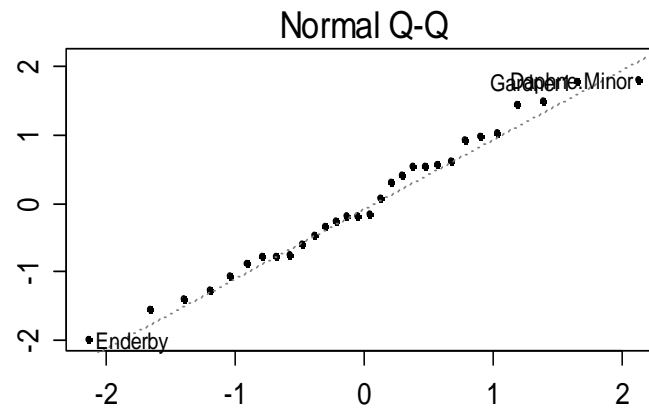
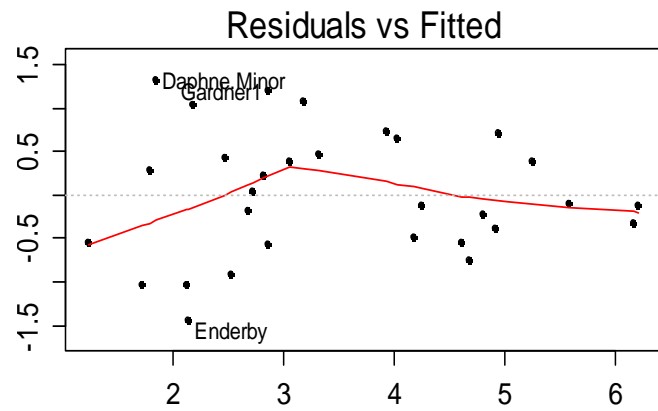
→ a log-transformation for all variables!

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Trying Multiple Linear Regression

```
> fit03 <- lm(log(Species) ~ log(Area) + ..., data=gala[, -2])
```



The normal plot is fine and there are no outliers.

But it seems that the relation has a bias. The variance is ↘.

Model needs to be improved!

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Poisson Regression: Theory

We have count response: $Y_i | X \sim Pois(\lambda_i)$

→ The goal is to relate the parameter λ_i , which is also the conditional expectation $\lambda_i = E[Y_i | X]$ linearly to the predictors. Since it takes positive values only, we require a log-trsf:

$$\log(\lambda_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

This is a GLM. The coefficients can be estimated by MLE. Assuming independence, the likelihood function is:

$$P(Y_1 = y_1, \dots, Y_n = y_n | X) = \prod_{i=1}^n P(Y_i = y_i | X) = \prod_{i=1}^n \frac{\lambda_i^{y_i} \cdot e^{-\lambda_i}}{y_i!}$$

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Poisson Regression: R Commands

```
> fit <- glm(Species ~ log(Area)+..., family=poisson, data=...)
> summary(fit)
```

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	3.323245	0.286430	11.602	< 2e-16	***
log(Area)	0.350370	0.018005	19.459	< 2e-16	***
log(Elevation)	0.033108	0.057034	0.580	0.56158	
log(Nearest)	-0.040153	0.014071	-2.854	0.00432	**
I(log(Scruz + 0.4))	-0.035848	0.013207	-2.714	0.00664	**
log(Adjacent)	-0.089452	0.006944	-12.882	< 2e-16	***

Null deviance: 3510.73 on 29 degrees of freedom
Residual deviance: 359.94 on 24 degrees of freedom
AIC: 532.77

→ These results are based on numerical optimization.
Thus, watch the convergence of the IRLS algorithm.

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Does the Model Fit?

Quick check: *residual deviance* \gg *df* ???

More precisely:
$$D = 2 \sum_{i=1}^n \left[y_i \log \left(\frac{y_i}{\hat{\lambda}_i} \right) - (y_i - \hat{\lambda}_i) \right] \sim \chi_{n-(p+1)}^2$$

Thus, when testing H_0 : "*Model is correct*", we obtain:

```
> pchisq(359.94, 24, lower=FALSE)
[1] 1.185031e-61
```

- The ***model does not fit well***. There is (much) more variation in the response than the Poisson distribution alone suggests. *Why is this and where does it come from?*
- Diagnostic plots / visualization is key!

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Residual Analysis

Analyze deviance or (as in R) Pearson residuals:

$$P_i = \frac{(y_i - \hat{\lambda}_i)}{\sqrt{\hat{\lambda}_i}} \quad \text{approx. } \sim N(0,1)$$

Thus, residuals $|P_i| > 2$ are bigger than the Poisson distribution suggests. And even larger residuals $|P_i| > 4$ would not exist if the Poisson model was correct.

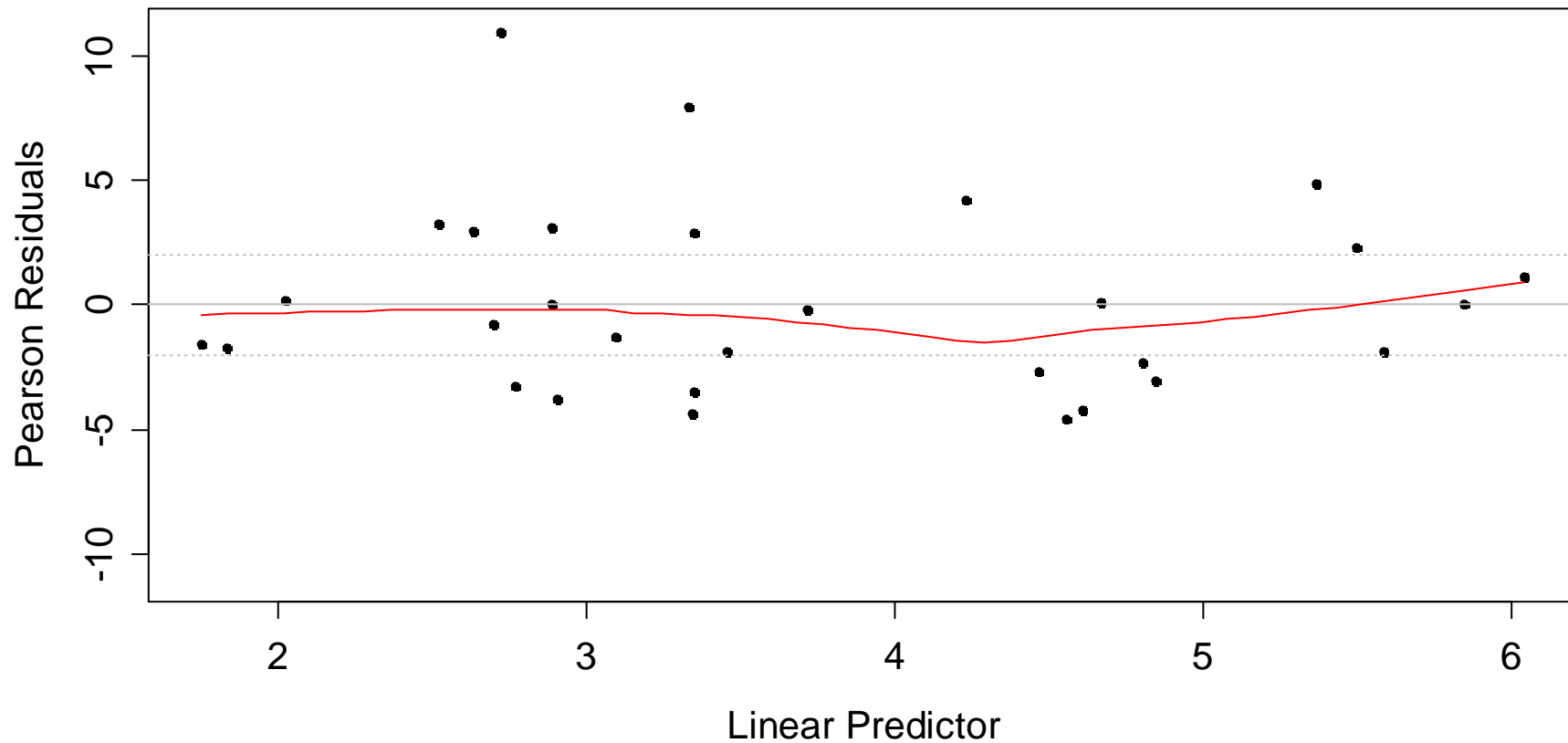
```
> xx <- predict(fit, type="link")
> yy <- resid(fit, type="pearson")
> plot(xx, yy, main="Tukey-Anscombe Plot...")
> lines(loess.smooth(xx, yy), col="red")
```

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Tukey-Anscombe Plot

Tukey-Anscombe Plot for Galapagos Tortoise



Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Dealing with Overdispersion

If the predictor-response relation is correct, but the variation is observed to be bigger than the distribution model suggests:

$\hat{\beta}_0, \dots, \hat{\beta}_p$ and $\hat{\lambda}_i$ unbiased

Standard errors $se(\hat{\beta}_0), \dots, se(\hat{\beta}_p)$ are wrong

Standard errors are corrected using a dispersion parameter:

$$\hat{\phi} = \frac{\sum_i (y_i - \hat{\lambda}_i)^2 / \hat{\lambda}_i}{n - (p + 1)}$$

In R:

```
> sum(resid(fit, type="pearson")^2) / fit$df.res  
[1] 16.64651
```

Applied Statistical Regression

AS 2015 – Generalized Linear Modeling

Final Result

```
> summary(fit, dispersion=16.64651)
```

```
Coefficients:
```

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	3.32325	1.16864	2.844	0.00446	**
log(Area)	0.35037	0.07346	4.769	1.85e-06	***
log(Elevation)	0.03311	0.23270	0.142	0.88686	
log(Nearest)	-0.04015	0.05741	-0.699	0.48430	
I(log(Scruz + 0.4))	-0.03585	0.05389	-0.665	0.50589	
log(Adjacent)	-0.08945	0.02833	-3.157	0.00159	**

```
---
```

```
Dispersion parameter for poisson family: 16.647
```

```
Null deviance: 3510.73 on 29 degrees of freedom
```

```
Residual deviance: 359.94 on 24 degrees of freedom
```

```
AIC: 532.77
```

→ Inference result is similar to the one from multiple linear regression. **Mathematics says: this is not a surprise!**