Flexible Regression and Smoothing: Using GAMLSS in R

Mikis Stasinopoulos, Rob Rigby, Gillian Heller, Vlasios Voudouris and Fernanda De Bastiani

Graz University of Technology, Austria, November 2016



Stasinopoulos et al.

Flexible Regression and Smoothing:

2016 1 / 43



2 Motivating examples

- The Dutch boys data
- The Munich rent data
- The fish species data
- A stylometric application

3 What is GAMLSS?





About the course

- Day 1 (Morning)
 - Why GAMLSS?
 - Introduction to the R packages, Diagnostics and Algorithms
 - Practical
- 2 Day 1 (Afternoon)
 - The gamlss.family distributions (Continuous, Discrete and Mixed distributions)
 - Practical
- Oay 2 (Morning)
 - Additive terms (linear, smoothing and random effects)
 - Practical
- Oay 2 (Afternoon)
 - Model selection
 - Centile estimation
 - Practical



3 / 43

2016

Information

The R Series **Flexible Regression** and Smoothing Using GAMLSS in R Mikis D. Stasinopoulos Robert A. Rigby Gillian Heller **Vlasios Voudouris** Fernanda De Bastiani CRC CRC Press



Stasinopoulos et al.

Flexible Regression and Smoothing:

A CHAPMAN & HALL BOOK

2016 4 / 43

The Dutch boys data

BMI : the BMI of 7294 boys age : the age in years Source: van Buuren and Fredriks (2001)



The Dutch boys data: statistical challenges



Agamlss

Flexible Regression and Smoothing:

The Dutch boys data: Histograms by age



The Dutch boys data: centile estimation





age

The Dutch boys data: centiles





The Munich rent data

- ${\tt R}\,$: the monthly net rent for flats in the city of Munich.
- $\ensuremath{\texttt{Fl}}$: the floor space area in square meters
 - A : year of construction
- $\verb"loc"$: whether the location is below, 1, average, 2, or above average 3
 - ${\tt H}\,$: two level factor indicating whether there is central heating, (0), or not, (1).

Source: Munich rental guide 1993



The Munich rent data



Flexible Regression and Smoothing:

2016 11 / 43

The fish species data

The fish species data



The fish species data





2016 13 / 43

A stylometric application

64 observations

- word : is the number of times a word appears in a single text
- freq : the number of different words which occur exactly
 word times in the text

Source: Prof. Mario Cortina-Borja



The stylometric data



What we need for modelling the above data?

We need

- flexible distributions for the response variable
- to be able to deal with heterogeneity in the data
- to be able to model skewness and kurtosis
- to be able to model overdispersion in count data
- We need modelling all the parameters of the distributions
- flexible functions to model the relationship between the parameter of the distribution and the explanatory variables

The Munich rent data



2016 17 / 43

A stylometric application

The Munich rent data: Linear model

Model assumptions

$$\mathbf{y} \stackrel{\mathrm{ind}}{\sim} N(\mu, \sigma^2),$$

 $\mu = \mathbf{X} \boldsymbol{\beta}$

Estimation

 $\hat{oldsymbol{eta}} = (\mathbf{X}^{ op}\mathbf{X})^{-1}\mathbf{X}^{ op}\mathbf{y}$

$$\hat{\sigma}^2 = \frac{\hat{\epsilon}^{\top}\hat{\epsilon}}{n}$$



2016 18 / 43

The linear model assumptions



2016 19 / 43

The Munich rent data: Linear model

```
r1 <- gamlss(R ~ Fl+A+H+loc, family=NO, data=rent)
## GAMLSS-RS iteration 1: Global Deviance = 28159
## GAMLSS-RS iteration 2: Global Deviance = 28159
11 <- lm(R ~ Fl+A+H+loc,data=rent)</pre>
coef(r1)
  (Intercept) Fl
                                   H1
                                                    loc2
##
                                 A
## -2775.038803 8.839445 1.480755 -204.759562 134.052349
        loc3
##
## 209.581472
coef(11)
             Fl
                                   H1
   (Intercept)
                                 A
                                                    loc2
##
## -2775.038803
             8.839445 1.480755 -204.759562 134.052349
        loc3
##
## 209.581472
                                                         amlss
```

2016 20 / 43

The Munich rent data: LM residuals plot



2016 21 / 43

The Munich rent data: Generalised Linear Model

The model

$$\mathbf{y} \stackrel{\text{ind}}{\sim} ExpFamily(\mu, \phi)$$

 $g(\mu) = \mathbf{X} \boldsymbol{\beta}.$

The exponential family

$$f_Y(y;\mu,\sigma) = \exp\left\{rac{y heta - b(heta)}{\phi} + c(y,\phi)
ight\}$$



The Munich rent data: GLM fit

```
12 <- glm(R ~ Fl+A+H+loc, family=Gamma(link="log"), data=rent)
r2 <- gamlss(R ~ Fl+A+H+loc, family=GA, data=rent)
## GAMLSS-RS iteration 1: Global Deviance = 27764.59
## GAMLSS-RS iteration 2: Global Deviance = 27764.59</pre>
```

The Munich rent data: GLM and GAIC

```
r22 <- gamlss(R ~ Fl+A+H+loc, family=IG, data=rent)
## GAMLSS-RS iteration 1: Global Deviance = 27991.56
## GAMLSS-RS iteration 2: Global Deviance = 27991.56
GAIC(r1, r2, r22) # AIC
      df
##
              AIC
## r2 7 27778.59
## r22 7 28005.56
## r1 7 28173.00
GAIC(r1, r2, r22, k=log(length(rent$R))) # SBC or BIC
      df
              ATC
##
## r2
       7 27817,69
## r22 7 28044.66
## r1 7 28212.10
```

The Munich rent data: GLM residuals plot



2016 25 / 43

The Munich rent data: Generalised Additive Model

The model

$$\begin{array}{rcl} \mathbf{y} & \stackrel{\mathrm{ind}}{\sim} & ExpFamily(\boldsymbol{\mu}, \boldsymbol{\phi}) \\ g\left(\boldsymbol{\mu}\right) & = & \mathbf{X}\boldsymbol{\beta} + s_1(\mathbf{x}_1) + \ldots + s_J(\mathbf{x}_J) \end{array}$$



Stasinopoulos et al.

Flexible Regression and Smoothing:

2016 26 / 43

The Munich rent data: GAM commands

```
r3 <- gamlss(R ~ pb(Fl)+pb(A)+H+loc, family=GA, data=rent)
## GAMLSS-RS iteration 1: Global Deviance = 27683.22
## GAMLSS-RS iteration 2: Global Deviance = 27683.22
## GAMLSS-RS iteration 3: Global Deviance = 27683.22
AIC(r2,r3)
## df AIC
## r3 11.21547 27705.65
## r2 7.00000 27778.59</pre>
```

The Munich rent data: GAM term plot



Flexible Regression and Smoothing:

2016 28 / 43

The Munich rent data: GAM summary

```
summary(r3)
## Family: c("GA", "Gamma")
##
## Call:
## gamlss(formula = R ~ pb(Fl) + pb(A) + H + loc, family = GA, data = rent)
##
## Fitting method: RS()
##
## ------
## Mu link function: log
## Mu Coefficients:
             Estimate Std. Error t value Pr(>|t|)
##
## (Intercept) 3.0851197 0.5666171 5.445 5.84e-08 ***
## pb(Fl) 0.0103084 0.0004030 25.578 < 2e-16 ***
## pb(A) 0.0014062 0.0002879 4.884 1.12e-06 ***
## H1
         -0.3008111 0.0225705 -13.328 < 2e-16 ***
## loc2
        0.1886692 0.0299153 6.307 3.51e-10 ***
## loc3
          0.2719856 0.0322699
                                8.428 < 2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The Munich rent data: GAM drop one term

```
drop1(r3)
## Single term deletions for
## mu
##
## Model:
## R ~ pb(Fl) + pb(A) + H + loc
             Df AIC LRT Pr(Chi)
##
## <none>
                27706
## pb(Fl) 1.4680 28261 558.59 < 2.2e-16 ***
## pb(A) 4.3149 27798 101.14 < 2.2e-16 ***
## H 1.8445 27862 160.39 < 2.2e-16 ***
## loc 2.0346 27770 68.02 1.825e-15 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



The Munich rent data: GAM worm plot



2016 31 / 43

Motivating examples A stylometr

A stylometric application

The Munich rent data: Mean and dispersion additive models

The model

$$\begin{array}{rcl} \mathbf{y} & \stackrel{\text{ind}}{\sim} & D(\boldsymbol{\mu}, \boldsymbol{\sigma}) \\ g_{1}(\boldsymbol{\mu}) & = & \mathbf{X}_{1} \beta_{1} + s_{11}(\mathbf{x}_{11}) + \ldots + s_{1J_{1}}(\mathbf{x}_{1J_{1}}) \\ g_{2}(\boldsymbol{\sigma}) & = & \mathbf{X}_{2} \beta_{2} + s_{21}(\mathbf{x}_{21}) + \ldots + s_{2J_{2}}(\mathbf{x}_{2J_{2}}) \end{array}$$



2016

32 / 43

The Munich rent data: MADAM commands

```
r4 <- gamlss(R ~ pb(Fl)+pb(A)+H+loc, sigma.fo=~pb(Fl)+pb(A)+H+loc, family=GA,
              data=rent)
## GAMLSS-RS iteration 1: Global Deviance = 27572.14
## GAMLSS-RS iteration 2: Global Deviance = 27570.29
## GAMLSS-RS iteration 3: Global Deviance = 27570.28
## GAMLSS-RS iteration 4: Global Deviance = 27570.28
r5 <- gamlss(R ~ pb(Fl)+pb(A)+H+loc, sigma.fo=~pb(Fl)+pb(A)+H+loc, family=IG,
              data=rent)
## GAMLSS-RS iteration 1: Global Deviance = 27675.74
## GAMLSS-RS iteration 2: Global Deviance = 27672.97
## GAMLSS-RS iteration 3: Global Deviance = 27673
## GAMLSS-RS iteration 4: Global Deviance = 27673.01
## GAMLSS-RS iteration 5: Global Deviance = 27673.01
## GAMLSS-RS iteration 6: Global Deviance = 27673.02
AIC(r3, r4, r5)
           df
##
                    AIC
## r4 22,25035 27614,78
## r3 11.21547 27705.65
```

Agamlss

r5 21.82318 27716.66

The Munich rent data: MADAM term plot for σ



Flexible Regression and Smoothing:

2016 34 / 43

The Munich rent data: MADAM worm plot



Flexible Regression and Smoothing:

2016 35 / 43

The Munich rent data: GAMLSS

The model

$$\begin{array}{lll} \mathbf{y} & \stackrel{\text{ind}}{\sim} & D(\mu, \sigma, \nu, \tau) \\ g_1(\mu) & = & \mathbf{X}_1 \beta_1 + s_{11}(\mathbf{x}_{11}) + \ldots + s_{1J_1}(\mathbf{x}_{1J_1}) \\ g_2(\sigma) & = & \mathbf{X}_2 \beta_2 + s_{21}(\mathbf{x}_{21}) + \ldots + s_{2J_2}(\mathbf{x}_{2J_2}) \\ g_3(\nu) & = & \mathbf{X}_3 \beta_3 + s_{31}(\mathbf{x}_{31}) + \ldots + s_{3J_3}(\mathbf{x}_{3J_3}) \\ g_4(\tau) & = & \mathbf{X}_4 \beta_4 + s_{41}(\mathbf{x}_{41}) + \ldots + s_{4J_4}(\mathbf{x}_{4J_4}) \end{array}$$



A stylometric application

The Munich rent data: GAMLSS assumptions



2016 37 / 43

The Munich rent data: GAMLSS commands

r7 <- gamlss(R ~ pb(F1)+pb(A)+H+loc,sigma.fo=~pb(F1)+pb(A)+H+loc, nu.fo=~pb(F1)+pb(A)+H+loc, family=BCCGo, data=rent)

AIC(r4, r6, r7) ## df AIC ## r7 28.41391 27608.15 ## r6 22.48092 27611.02 ## r4 22.25035 27614.78

The Munich rent data: GAMLSS worm plot



Flexible Regression and Smoothing:

2016 39 / 43

What is GAMLSS?

GAMLSS: are semi-parametric regression type models.

- regression type: we have many explanatory variables X and one response variable y and we believe that $X \to y$
- parametric: a parametric distribution assumption for the response variable,
- semi: the parameters of the distribution, as functions of explanatory variables, may involve non-parametric smoothing functions
- GAMLSS philosophy: try different models

GAMLSS is a generalisation of GLM and GAM models.



Conclusions

- GAMLSS is a very flexible statistical model
- It is a unified framework for univariate regression type of models
- Allows any distribution for the response variable Y
- Models all the parameters of the distribution of Y
- Allows a variety of penalised additive terms in the models for the distribution parameters
- The fitted algorithm is modular, where different components can be added easily
- it can easily introduced to students since it relies on known concepts
- It deals with overdispersion, skewness and kurtosis

This is a collaborative work: A list of people who help with the R implementation

present	past
Vlasios Voudouris	Popi Akantziliotou
Paul Eilers	Nicoleta Mortan
Gillian Heller	Fiona McElduff
Marco Enea	Raydonal Ospina
Majid Djennad	Konstantinos Pateras
Fernanda De Bastiani	
Luiz Nakamura	
Daniil Kiose	
Andreas Mayr	
Thomas Kneib	
Nadja Klein	
Abu Hossain	

Conclusions

There are lot more to be done

the END

for more information see

www.gamlss.org



Stasinopoulos et al.

Flexible Regression and Smoothing:

2016 43 / 43