Models for Additive Response Components RESEARCH



ISO 9001 certified

JOANNEUM



JOANNEUM RESEARCH Forschungsgesellschaft mbH

> Institute of Applied Statistics

Motivation

Regression with Gaussian errors can be seen as additive Figure 2: Binominal Model for Rates with linear $f(x_i)$ and decomposition of the expected response $\beta_1 = \beta_2 = 1$

$$E(Y) = \mu = \sum_{j} \mu_{j} \text{ with } \mu_{j} = f_{j}(X_{j})$$

or $\mu_{j} = \sum_{j} X_{j} \beta_{j}$

Components μ_i refer to the predictor variables X_i .

Does not hold for GLM/GAM



Risk Estimation

Comparing population segments, e.g. regions, s = 1, ..., S

Crude Rate

 $cr_s = O_s / N_s$ $O_s = \sum_c O_{cs}$ the observed disease frequency c = 1, ..., C age classes

Controlling for age effects

Standardization is usually applied

and System Analysis

Gerhard Neubauer

Steyrergasse 25a A-8010 Graz, Austria

Phone: +43 316 876-15 57 Fax: +43 316 876-915 57

gerhard.neubauer@joanneum.at http://www.joanneum.at/sta



Poisson Model with Additive Response Components

 $Y \sim P(\lambda)$ with canonical link function $g(.) = \ln(.)$ and $h(.) = \exp(.)$

have

Relies on the choice of a standard population

Results depend on this choice

Breslow & Day, 1980, 1987

Standardized Rate

 $sr_s = E_s / N_s$ $E_s = \sum_c E_{cs}$ the expected disease frequency $E_{cs} = W_{cs}O_{cs}, \quad 0 \le W_{cs} \le \infty$

The effect of age is not removed but equalized

Adjusted Rate

Neubauer, 2001

 $ar_s = A_s / N_s$ $A_s = \sum_c \mu_0(X_{0cs})$ with $\mu_0(X_{0cs})$ estimated from the additive response component model

 $E(O_{cs}) = \mu_0(X_{0cs}) + \mu_1(age_{cs})$

The effect of age is estimated and removed from the prediction

Application to Hospital Discharge Data

The Data

Hospital discharges with a ICD-9 diagnosis from the group Cardiovascular diseases in Styria 1995–2000

The Problem

Identify high-risk regions on the NUTS-3 level. The six NUTS-3 regions in Styria are: GRAZ, LIEZEN, NORTH WEST, NORTH EAST, SOUTH WEST, SOUTH EAST

GLM/GAM The model for the mean

 $\mu = \exp(\eta) = \exp[\sum_{i} f_{i}(X_{i})] \ge 0$

Additive Response Components The model for the mean

 $\mu = \sum_{i} \mu_{i} = \sum_{i} \exp[f_{i}(X_{i})] \ge 0$

The additive components

 $\mu_{i} = \exp[f_{i}(X_{i})] \ge 0$

The nonlinear predictor

$$\eta = \ln(\mu) = \ln\{\sum_{j} \exp[f_{j}(X_{j})]\}$$

Figure 1: Poisson Model with cubic $f(x_i)$



Binomial Model



ML Estimation

For exponential family the log likelihood as a function of the canonical parameter θ_i is given by

$$l(\theta_i) = \frac{y_i \theta_i - b(\theta_i)}{\phi} + \text{constant.}$$

Let $f_{ij} = f_j(X_{ij})$ be a polynominal of degree d_j , i. e. $f_{ij} = \sum_{k=1}^{a_j} x_{ij}^k \beta_{jk}$, then $\underline{y_i - \mu_i} \quad \frac{\partial \mu_i}{\partial \mu_i} \quad \frac{\partial \eta_i}{\partial \mu_i} \quad x_{ir}^a$ $\partial l_i \quad \partial \eta_i$ ∂l_i $\partial \eta_i \quad \partial f_{ir}$ $\partial \beta_{ra}$ $\partial \eta_i$ $\partial \beta_{ra}$ V_i¢

For the Poisson and the Binominal model we have $\eta_i = \ln[\sum_i \exp(f_{ii})]$ and further using canonical links we obtain

$$\frac{\partial l_i}{\partial \beta_{ra}} = \frac{y_i - \mu_i}{V_i \phi} \mu_{ir} x^a_{ir}, \text{ and }$$

$$\frac{\partial l_i}{\partial \beta_{ra}} = \frac{y_i - \mu_i}{V_i \phi} (1 - \mu_i) \mu_{ir} x_{ir}^a$$

for the Poisson and the Binominal Model respectively. In matrix notation we have

Figure 3: The six NUTS-3 regions of Styria



The Model

 Y_{ir} the number of discharges in region r

 $E(Y_{ir}) = \mu_{ir} = \exp[g_r(sex_{ir}, size_{ir})] + \exp[g_r(age_{ir})]$ where size_{ir} is the size of the population, i.e. the number of people at risk and g(.) is a third order polynomial.

The Adjusted Rate

$$ar_{s} = n_{s}^{-1} \sum_{i} \exp[g_{r}(sex_{ir}, size_{ir})]$$

Results

Adjusted rates are much smaller than standardized rates, as for

- sr_s age effects are equalized, and for
- **a**r_s age effects are estimated and removed.

References

Breslow, N. E. & Day, N. E. (1987). Statistical Methods in Cancer Research. Vol. II—The Design and Analysis of Cohort Studies. Lyon: IARC.

Breslow, N. E. & Day, N. E. (1980).

Research. Vol. I—The Analysis of

Case-Control Studies. Lyon: IARC.

Statistical Methods in Cancer

Hastie, T. J. & Tibshirani, R. J. (1990). Generalized Additive Models. London: Chapman and Hall.

> McCullagh, P. and Nelder, J.A. (1983). Generalized Linear Models. London: Chapman and Hall.

Neubauer, G. (2001). Krankenhaus-Entlassungsdaten, <u>Krankheitshäufigkeiten</u> und statistische Modelle (Hospital discharge data, disease frequency and statistical models). Presentation at ROeS Seminar 2001, 24.9.-27.9.2001 Mayerhofen i. Z., Austria.

Acknowledgements

I would like to thank Herwig FriedI for helpful discussions and comments throughout the work on this paper.

oefpos02043

with Additive Response Components

 $Y \sim B(n, p)$ with canonical link function $g(\mu) = logit(\mu)$ and $h(\eta) = \exp(\eta) / [1 + \exp(\eta)]$

GLM/GAM

The model for the mean

$$0 \le \frac{\mu}{n} = \frac{\exp(\eta)}{1 + \exp(\eta)} = \frac{\exp[\sum_{j} f_{j}(X_{j})]}{1 + \exp[\sum_{j} f_{j}(X_{j})]} \le 1$$

Additive Response Components

The model for the mean

$$0 \le \frac{\mu}{n} = \frac{1}{n} \sum \mu_j = \frac{\sum_j \exp[f_j(X_j)]}{1 + \sum_j \exp[f_j(X_j)]} \le 1$$

The additive components

$$\mu_{k} = \frac{\exp[f_{k}(X_{k})]}{1 + \sum_{j} \exp[f_{j}(X_{j})]}$$

The nonlinear predictor

$$\eta = \ln[\mu / (1 - \mu)] = \ln\{\sum_{j} \exp[f_{j}(X_{j})]\}$$



- **n** x d_r-matrix X_r = $(\partial f_r / \partial \beta_{ra})$
- $\square D = diag(\partial \mu_i / \partial \eta_i)$
- W = diag $[1 / var(Y_i)]$
- **Q**r = diag $(\partial \eta_i / \partial f_{ir})$

Solving (1) by the Newton-Raphson algorithm requires the second derivates

$$\mathbf{Y}_{s} = \mathbf{X}_{r}^{T} \quad \left(\frac{\partial \mathbf{Q}_{r} \mathbf{D} \mathbf{W}}{\partial f_{s}} (\mathbf{y} - \boldsymbol{\mu}) - \mathbf{Q}_{r} \mathbf{D} \mathbf{W} \mathbf{Q}_{s} \frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}\right) \mathbf{X}_{s}. (2)$$

For canonical links (1) and (2) simplify as now DW = I and further using Fisher scoring l_{rs}'' is replaced by its expectation

$$E(l_{rs}'') = -X_r^T Q_r Q_s \frac{\partial \mu}{\partial \eta} X_s.$$

in the estimation algorithm.

Rank Order of Regions

Crude and standardized rates: North of Styria has highest risk,

Adjusted rates:

(1)

LIEZEN and SOUTH WEST have highest risk.

A recommendation for a health care intervention would differ depending on the chosen method.

Table 1: Rate estimates for the six NUTS-3 regions of Styria

	Crude rate		Standardized rate		Adjusted rate	
NUTS 3 region	Estimate	Rank	Estimate	Rank	Estimate	Rank
GRAZ	0.0312	6	0.0231	6	0.0029	6
LIEZEN	0.0322	4	0.0242	5	0.0097	2
NORTH WEST	0.0419	1	0.0292	1	0.0030	5
NORTH EAST	0.0357	2	0.0277	2	0.0062	3
SOUTH WEST	0.0313	5	0.0262	4	0.0101	1
SOUTH EAST	0.0324	3	0.0269	3	0.0053	4

a tradition of innovation