2 General Models slide 89

Smoking data

Table 2: Lung cancer deaths in British male physicians (Doll and Hill, 1952). The table gives man-years at risk T/number of cases y of lung cancer, cross-classified by years of smoking t, taken to be age minus 20 years, and number of cigarettes smoked per day, d.

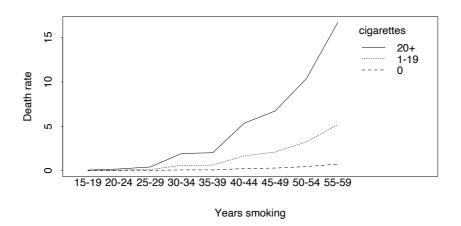
Years of	Daily cigarette consumption \emph{d}									
smoking t										
	Nonsmokers	1–9	10–14	15–19	20–24	25–34	35+			
15–19	10366/1	3121	3577	4317	5683	3042	670			
20-24	8162	2937	3286/1	4214	6385/1	4050/1	1166			
25-29	5969	2288	2546/1	3185	5483/1	4290/4	1482			
30-34	4496	2015	2219/2	2560/4	4687/6	4268/9	1580/4			
35–39	3512	1648/1	1826	1893	3646/5	3529/9	1336/6			
40–44	2201	1310/2	1386/1	1334/2	2411/12	2424/11	924/10			
45-49	1421	927	988/2	849/2	1567/9	1409/10	556/7			
50-54	1121	710/3	684/4	470/2	857/7	663/5	255/4			
55–59	826/2	606	449/3	280/5	416/7	284/3	104/1			

Regression Methods

Autumn 2024 - slide 90

Smoking data

Lung cancer deaths in British male physicians. The figure shows the rate of deaths per 1000 man-years at risk, for each of three levels of daily cigarette consumption.



Regression Methods

C	احدادا	4-+-
Smo	KING	uata

- \square Suppose number of deaths y has Poisson distribution, mean $T\lambda(d,t)$, where T is man-years at risk, d is number of cigarettes smoked daily and t is time smoking (years).
- □ Take

$$\lambda(d,t) = \beta_0 t^{\beta_1} \left(1 + \beta_2 d^{\beta_3} \right) :$$

- background rate of lung cancer is $\beta_0 t^{\beta_1}$ for non-smoker,
- additional risk due to smoking d cigarettes/day is $\beta_2 d^{\beta_3}$.
- \square With $x_i = (T_i, d_i, t_i)$, can write this as

$$y_j \sim \text{Poiss}\{\mu(\beta; x_j)\},$$

 $\mu(\beta; x) = T\beta_0 t^{\beta_1} \left(1 + \beta_2 d^{\beta_3}\right), \quad j = 1, \dots, n:$

a nonlinear model with Poisson-distributed response.

Regression Methods

Autumn 2024 - slide 92

Comments

- \Box Linear model $y \sim (X\beta, \sigma^2 I_n)$
 - applicable for continuous response $y \in \mathbb{R}$
 - assumes linear dependence of mean response $\mathrm{E}(y)$ on covariates X
 - sometimes assumes y normal
- ☐ Lots of data not like this
- □ Need extensions for
 - nonlinear dependence on covariates
 - arbitrary response distribution (binomial, Poisson, exponential, ...)
 - dependent responses
 - variance non-constant (and related to mean?)
 - censoring, truncation, . . .
 - _

Regression Methods

Autumn 2024 - slide 93

Simple fixes

- \square Just fit a linear model anyway
 - Might work as an approximation, but usually extrapolates really badly.
- ☐ Fit a linear model to transformed responses
 - E.g., take variance-stabilising transformation for y, such as $2\sqrt{y}$ when y is Poisson
 - Can be helpful, but usually the obvious transformation can't give linearity.
- ☐ Instead we attempt to fit the model using likelihood estimation.

Regression Methods

2.1 Inference slide 95

Revision: Likelihood

Definition 15 Let y be a data set, assumed to be the realisation of a random variable $Y \sim f(y; \theta)$, where the unknown parameter θ lies in the parameter space $\Omega_{\theta} \subset \mathbb{R}^p$. Then the likelihood (for θ based on y) and the corresponding \log likelihood are

$$L(\theta) = L(\theta; y) = f_Y(y; \theta), \quad \ell(\theta) = \log L(\theta), \quad \theta \in \Omega_{\theta}.$$

The maximum likelihood estimate (MLE) $\widehat{\theta}$ satisfies $\ell(\widehat{\theta}) \geq \ell(\theta)$, for all $\theta \in \Omega_{\theta}$. Often $\widehat{\theta}$ is unique and in many cases it satisfies the score (or likelihood) equation

$$\frac{\partial \ell(\theta)}{\partial \theta} = 0,$$

which is interpreted as a vector equation of dimension $p \times 1$ if θ is a $p \times 1$ vector. The observed information and expected (Fisher) information are defined as

$$J(\theta) = -\frac{\partial^2 \ell(\theta)}{\partial \theta \partial \theta^{\mathrm{T}}}, \quad I(\theta) = \mathrm{E} \left\{ J(\theta) \right\};$$

these are $p \times p$ matrices if θ has dimension p.

Regression Methods Autumn 2024 – slide 96

Revision: Maximum likelihood estimator

 \square In large samples from a **regular model** in which the true parameter is $\theta_{p\times 1}^0$, the maximum likelihood estimator $\widehat{\theta}$ has an approximate normal distribution,

$$\widehat{\theta} \stackrel{.}{\sim} \mathcal{N}_p \left\{ \theta^0, J(\widehat{\theta})^{-1} \right\},$$

so we can compute an approximate $(1-2\alpha)$ confidence interval for the rth parameter θ^0_r as

$$\widehat{\theta}_r \pm z_{\alpha} v_{rr}^{1/2},$$

where v_{rr} is the rth diagonal element of the matrix $J(\widehat{\theta})^{-1}.$

- ☐ This is easily implemented:
 - we code the negative log likelihood $-\ell(\theta)$ (and check the code carefully!);
 - we minimise $-\ell(\theta)$ numerically, ensuring that the minimisation routine returns $\widehat{\theta}$ and the Hessian matrix $J(\widehat{\theta}) = -\partial^2 \ell(\theta)/\partial \theta \partial \theta^{\mathrm{T}}|_{\theta=\widehat{\theta}}$
 - we compute $J(\widehat{\theta})^{-1}$, and use the square roots of its diagonal elements, $v_{11}^{1/2},\ldots,v_{dd}^{1/2}$, as standard errors for the corresponding elements of $\widehat{\theta}$.

Regression Methods

Revision: Regular model

We say that a statistical model $f(y;\theta)$ is regular (for likelihood inference) if

- 1. the true value θ^0 of θ is interior to the parameter space $\Omega_{\theta} \subset \mathbb{R}^p$;
- 2. the densities defined by any two different values of θ are distinct;
- 3. there is an open set $\mathcal{I} \subset \Omega_{\theta}$ containing θ^0 within which the first three derivatives of the log likelihood with respect to elements of θ exist almost surely, and

$$|\partial^3 \log f(Y_i; \theta)/\partial \theta_r \partial \theta_s \partial \theta_t| \le g(Y_i)$$

uniformly for $\theta \in \mathcal{I}$, where $0 < E_0\{g(Y_i)\} = K < \infty$; and

4. for $\theta \in \mathcal{I}$ we can interchange differentation with respect to θ and integration, that is,

$$\frac{\partial}{\partial \theta} \int f(y;\theta) \ dy = \int \frac{\partial f(y;\theta)}{\partial \theta} \ dy, \quad \frac{\partial^2}{\partial \theta \partial \theta^{\mathrm{T}}} \int f(y;\theta) \ dy = \int \frac{\partial^2 f(y;\theta)}{\partial \theta \partial \theta^{\mathrm{T}}} \ dy.$$

The results are also true under weaker conditions, for non-identically distributed and dependent data.

Regression Methods

Autumn 2024 - slide 98

Revision: Comments on regular models

Condition

- 1. is needed so that $\widehat{\theta}$ can lie 'on both sides' of θ^0 and hence can have a limiting normal distribution, once standardized—fails, for example, if θ has a discrete component (e.g. changepoint $\gamma \in \{1, \ldots, n\}$);
- 2. is needed to be able to identify the model on the basis of the data;
- 3. ensures the validity of Taylor series expansions of $\ell(\theta)$ —not usually a problem;
- 4. ensures that the score statistic has a limiting normal distribution—can fail in some models sometimes good news, leading to faster convergence than $n^{-1/2}$.

All the above assumes the postulated model is correct! — there is a literature on what happens when we fit the wrong model, or if the parameter dimension increases with n, or ... usually there are no generic results for such cases.

Regression Methods

Revision: Likelihood ratio statistic

- \square Model $f_B(y)$ is **nested** within model $f_A(y)$ if A reduces to B on restricting some parameters:
 - for example, the model $Y_1,\ldots,Y_n\stackrel{\mathrm{iid}}{\sim}\mathcal{N}(0,\sigma^2)$ is nested within the model $Y_1,\ldots,Y_n\stackrel{\mathrm{iid}}{\sim}\mathcal{N}(\mu,\sigma^2)$, because the first is obtained from the second by setting $\mu=0$;
 - the maximised log likelihoods satisfy $\widehat{\ell}_A \geq \widehat{\ell}_B$, because the more comprehensive model A contains the simpler model B.
- ☐ The likelihood ratio statistic for comparing them is

$$W = 2(\widehat{\ell}_A - \widehat{\ell}_B).$$

 \Box If the model is regular, the simpler model is true, and A has q more parameters than B, then

$$W \stackrel{.}{\sim} \chi_q^2$$
.

This implicitly assumes that ML inference for model A is OK, so that the approximation $\widehat{\theta}_A \stackrel{.}{\sim} \mathcal{N}\{\theta_A, J_A(\widehat{\theta}_A)^{-1}\}$ is adequate.

Regression Methods

Autumn 2024 - slide 100

Revision: Profile log likelihood

- Consider a regular log likelihood $\ell(\psi, \lambda)$, where the **parameter of interest** ψ is variation independent of the **nuisance parameter** λ , i.e., $(\psi, \lambda) \in \Omega_{\psi} \times \Omega_{\lambda}$, and the overall MLE is $(\widehat{\psi}, \widehat{\lambda})$.
- \square For a confidence set for ψ , without reference to λ , we use the **profile log likelihood**

$$\ell_{\mathrm{p}}(\psi) = \max_{\lambda \in \Omega_{\lambda}} \ell(\psi, \lambda) = \ell(\psi, \widehat{\lambda}_{\psi}),$$

say, and, based on the limiting distribution of the likelihood ratio statistic, take as $(1-2\alpha)$ confidence region the set

$$\left\{ \psi \in \Omega_{\psi} : 2\{\ell(\widehat{\psi}, \widehat{\lambda}) - \ell(\psi, \widehat{\lambda}_{\psi})\} \le \chi^{2}_{\dim \psi}(1 - 2\alpha) \right\}.$$

 \square When ψ is scalar, this yields

$$\left\{ \psi \in \Omega_{\psi} : \ell(\psi, \widehat{\lambda}_{\psi}) \right\} \ge \ell(\widehat{\psi}, \widehat{\lambda}) - \frac{1}{2}\chi_1^2 (1 - 2\alpha) \right\},$$

and $\frac{1}{2}\chi_1^2(0.95) = 1.92$.

 \square Such intervals are generally better than the standard interval $\widehat{\psi}\pm z_{\alpha}\mathrm{SE}$, particularly when the distribution of $\widehat{\psi}$ is asymmetric, but require more computation, since they involve many maximisations of ℓ .

Regression Methods

Model setup

- Independent random variables Y_1, \ldots, Y_n , with observed values y_1, \ldots, y_n , and covariates x_1, \ldots, x_n .
- \square Suppose that probability density of Y_j is $f(y_j; \eta_j, \phi)$, where $\eta_j = \eta(\beta, x_j)$, and ϕ is common to all models.
- ☐ Log likelihood is

$$\ell(\beta, \phi) = \sum_{j=1}^{n} \ell_j(\beta, \phi) = \sum_{j=1}^{n} \log f\{y_j; \eta(\beta, x_j), \phi\}.$$

- \square More generally, just let $\ell_j(\beta,\phi)$ denote the log likelihood contribution from the $j{
 m th}$ observation.
- \square Suppose ϕ known (for now), suppress it, and estimate β .

Example 16 (Normal regression model) Express the normal regression model in the terms above.

Regression Methods

Autumn 2024 - slide 102

Note to Example 16

Here $Y_j \stackrel{\mathrm{ind}}{\sim} \mathcal{N}(\mu_j, \sigma^2)$ with $\mu_j = \eta_j = \eta(x_j; \beta)$, so obviously

$$\eta_j = \eta(x_j; \beta), \quad \phi = \sigma^2, \quad \ell_j \equiv -\frac{1}{2} \{ (y_j - \eta_j)^2 / \phi + \log \phi \}.$$

Regression Methods

Autumn 2024 - note 1 of slide 102

Iterative weighted least squares (IWLS)

- $\hfill \Box$ General approach for estimation in regression models, based on Newton–Raphson iteration
- \square Assume that ϕ is fixed, and write

$$\ell(\beta) = \sum_{j=1}^{n} \ell_j \{ \eta_j(\beta) \}.$$

 \square MLEs $\widehat{\beta}$ usually satisfy

$$\frac{\partial \ell(\widehat{\beta})}{\partial \beta_r} = 0, \qquad r = 1, \dots, p,$$

or equivalently

$$\frac{\partial \ell(\widehat{\beta})}{\partial \beta} = \frac{\partial \eta^{\mathrm{T}}}{\partial \beta} \frac{\partial \ell}{\partial \eta} = \frac{\partial \eta^{\mathrm{T}}}{\partial \beta} u(\widehat{\beta}) = \sum_{j=1}^{n} \frac{\partial \eta_{j}}{\partial \beta} \frac{\partial \ell_{j} \{ \eta_{j}(\beta) \}}{\partial \eta_{j}} = 0, \tag{6}$$

where $u(\beta)$ is $n \times 1$ vector with jth element $\partial \ell / \partial \eta_i$.

Regression Methods

IWLS II $\widehat{\beta} = (X^{\mathrm{T}}WX)^{-1}X^{\mathrm{T}}Wz,$ where $X_{n\times p} = \partial \eta/\partial \beta^{\mathrm{T}}, \quad \text{(design matrix)}$ $W_{n\times n} = \operatorname{diag}\{\mathrm{E}(-\partial^2\ell_j/\partial \eta_j^2)\}, \quad \text{(weights)}$ $z_{n\times 1} = X\beta + W^{-1}u, \quad \text{(adjusted dependent variable)}$ Thus to obtain MLEs $\widehat{\beta}$ we use the IWLS algorithm: \square take an initial $\widehat{\beta}$. Repeat $-\operatorname{compute} X, W, u, z;$ $-\operatorname{compute} w \widehat{\beta} \text{ and replace the preceding value;}$ until changes in $\ell(\widehat{\beta})$ (or, sometimes, $\widehat{\beta}$, or both) are lower than some tolerance. \square Sometimes a line search is added, if $\ell(\widehat{\beta}_{\mathrm{new}}) < \ell(\widehat{\beta}_{\mathrm{old}})$: i.e., we half the step length and try again.

Regression Methods

Derivation of IWLS algorithm

 \square To find the maximum likelihood estimate $\widehat{\beta}$ starting from a trial value β , we make a Taylor series expansion in (6), to obtain

$$\frac{\partial \eta^{\mathrm{T}}(\beta)}{\partial \beta} u(\beta) + \left\{ \sum_{j=1}^{n} \frac{\partial \eta_{j}(\beta)}{\partial \beta} \frac{\partial^{2} \ell_{j}(\beta)}{\partial \beta^{2}} \frac{\partial \eta_{j}(\beta)}{\partial \beta^{\mathrm{T}}} + \sum_{j=1}^{n} \frac{\partial^{2} \eta_{j}(\beta)}{\partial \beta \partial \beta^{\mathrm{T}}} u_{j}(\beta) \right\} (\widehat{\beta} - \beta) \doteq 0.$$
 (7)

If we denote the $p \times p$ matrix in braces on the left by $-J(\beta)$, assumed invertible, we can rearrange (7) to obtain

$$\widehat{\beta} \doteq \beta + J(\beta)^{-1} \frac{\partial \eta^{\mathrm{T}}(\beta)}{\partial \beta} u(\beta). \tag{8}$$

This suggests that maximum likelihood estimates may be obtained by starting from a particular β , using (8) to obtain $\widehat{\beta}$, then setting β equal to $\widehat{\beta}$, and iterating (8) until convergence. This is the Newton–Raphson algorithm applied to our particular setting. In practice it can be more convenient to replace $J(\beta)$ by its expected value

$$I(\beta) = \sum_{j=1}^{n} \frac{\partial \eta_{j}(\beta)}{\partial \beta} E\left(-\frac{\partial^{2} \ell_{j}}{\partial \eta_{j}^{2}}\right) \frac{\partial \eta_{j}(\beta)}{\partial \beta^{T}};$$

the other term vanishes because $E\{u_j(\beta)\}=0$. We write

$$I(\beta) = X(\beta)^{\mathrm{T}} W(\beta) X(\beta), \tag{9}$$

where $X(\beta)$ is the $n \times p$ matrix $\partial \eta(\beta)/\partial \beta^{\mathrm{T}}$ and $W(\beta)$ is the $n \times n$ diagonal matrix whose jth diagonal element is $\mathrm{E}(-\partial^2 \ell_j/\partial \eta_i^2)$.

 \square If we replace $J(\beta)$ by $X(\beta)^{\mathrm{T}}W(\beta)X(\beta)$ and reorganize (8), we obtain

$$\widehat{\beta} = (X^{\mathrm{T}}WX)^{-1}X^{\mathrm{T}}W(X\beta + W^{-1}u) = (X^{\mathrm{T}}WX)^{-1}X^{\mathrm{T}}Wz, \tag{10}$$

say, where the dependence of the terms on the right on β has been suppressed. That is, starting from β , the updated estimate $\widehat{\beta}$ is obtained by weighted linear regression of the $n\times 1$ vector adjusted dependent variable

$$z = X(\beta)\beta + W(\beta)^{-1}u(\beta)$$

on the columns of $X(\beta)$, using weight matrix $W(\beta)$. The maximum likelihood estimates are obtained by repeating this step until the log likelihood, the estimates, or more often both, are essentially unchanged. The variable z plays the role of the response or dependent variable in the weighted least squares step.

 \square Often the structure of a model simplifies the estimation of an unknown value of ϕ . It may be estimated by a separate step between iterations of $\widehat{\beta}$, by including it in the step (8), or from the profile log likelihood $\ell_p(\phi)$.

Regression Methods

Autumn 2024 - note 1 of slide 104

Examples

Example 17 (Normal nonlinear model) Give the components of the IWLS algorithm for the normal nonlinear model.

Regression Methods

Note to Example 17

 \square Here the mean of the jth observation is $\eta_j = \eta(x_j; \beta)$. The log likelihood contribution $\ell_j(\eta_j)$ is

$$\ell_j(\eta_j, \sigma^2) \equiv -\frac{1}{2} \left\{ \log \sigma^2 + \frac{1}{\sigma^2} (y_j - \eta_j)^2 \right\},$$

SO

$$u_j(\eta_j) = \frac{\partial \ell_j}{\partial \eta_j} = \frac{1}{\sigma^2} (y_j - \eta_j), \qquad \frac{\partial^2 \ell_j}{\partial \eta_i^2} = -\frac{1}{\sigma^2};$$

the *j*th element on the diagonal of W is the constant σ^{-2} .

The jth row of the matrix $X = \partial \eta/\partial \beta^{\mathrm{T}}$ is $(\partial \eta_j/\partial \beta_0, \dots, \partial \eta_j/\partial \beta_{p-1})$, and as η_j is nonlinear as a function of β , X depends on β .

After some simplification, we see that the new value for $\widehat{\beta}$ given by (10) is

$$\widehat{\beta} \doteq (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}(X\beta + y - \eta),\tag{11}$$

where X and η are evaluated at the current β . Here $\eta \neq X\beta$ and (11) must be iterated.

The log likelihood is a function of β only through the sum of squares, $SS(\beta) = \sum_{j=1}^{n} \{y_j - \eta_j(\beta)\}^2$. The profile log likelihood for σ^2 is

$$\ell_{\rm p}(\sigma^2) = \max_{\beta} \ell(\beta, \sigma^2) \equiv -\frac{1}{2} \left\{ n \log \sigma^2 + \mathrm{SS}(\widehat{\beta}) / \sigma^2 \right\},$$

so the maximum likelihood estimator of σ^2 is $\widehat{\sigma}^2 = \mathrm{SS}(\widehat{\beta})/n$. Although $S^2 = \mathrm{SS}(\widehat{\beta})/(n-p)$ is not unbiased when the model is nonlinear, it turns out to have smaller bias than $\widehat{\sigma}^2$, and is preferable in applications.

In some cases the error variance depends on covariates, and we write the variance of the $j{\rm th}$ response as $\sigma_j^2 = \sigma^2(x_j, \gamma)$. Such models may be fitted by alternating iterative weighted least squares updates for β treating γ as fixed at a current value with those for γ with β fixed, convergence being attained when neither estimates nor log likelihood change materially.

Regression Methods

Autumn 2024 - note 1 of slide 105

Deviance

- Let $\widehat{\eta}_j = \eta_j(\widehat{\beta}, x_j)$, where $\widehat{\beta}$ is MLE of β , giving maximised log likelihood $\ell(\widehat{\beta})$ and $\widehat{\eta}^{\mathrm{T}} = (\widehat{\eta}_1, \dots, \widehat{\eta}_n)$.
- \square Let $\tilde{\eta}_j$ be the value of η_j that maximises $\log f(y_j;\eta_j)$, and let $\tilde{\eta}^{\mathrm{T}}=(\tilde{\eta}_1,\ldots,\tilde{\eta}_n)$. This corresponds to the saturated model, with

#parameters in $\eta = \#$ observations in y,

which will give the largest likelihood possible.

☐ Define the **scaled deviance**:

$$D = 2\sum_{j=1}^{n} \{ \log f(y_j; \tilde{\eta}_j) - \log f(y_j; \hat{\eta}_j) \} \ge 0.$$

- \square Small D implies $\widehat{\eta} \approx \widetilde{\eta}$, so model fits well.
- \square Large D implies poor fit like $SS(\widehat{\beta})$ in linear model.

Regression Methods

Differences of deviances

- ☐ Consider two models:
 - Model $A: \beta^{\mathrm{T}} = (\beta_1, \dots, \beta_p) \in \mathbb{R}^p$ vary freely MLEs $\widehat{\eta}^A = \eta(\widehat{\beta}^A)$;
 - Model $B: (\beta_1, \dots, \beta_q) \in \mathbb{R}^q$ vary freely, but $\beta_{q+1}, \dots, \beta_p$ are fixed hence q free parameters, MLEs $\widehat{\eta}^B = \eta(\widehat{\beta}^B)$.
- \square Model B is **nested within** model A: B can be obtained by restricting A.
- ☐ Likelihood ratio statistic for comparing the models is

$$2(\widehat{\ell}_A - \widehat{\ell}_B) = 2\sum_{j=1}^n \left\{ \log f(y_j; \widehat{\eta}_j^A) - \log f(y_j; \widehat{\eta}_j^B) \right\} = D_B - D_A,$$

and this $\stackrel{.}{\sim} \chi^2_{p-q}$ if the models are regular.

 \Box If ϕ unknown, replace it by an estimate: same distributional approximations will apply.

Example 18 (Normal linear model) Find the difference of deviances in the normal linear model.

Regression Methods

Autumn 2024 - slide 107

Note to Example 18

 \square Suppose that the y_j are normal with means η_j and known variance ϕ . Then

$$\log f(y_j; \eta_j, \phi) = -\frac{1}{2} \left\{ \log(2\pi\phi) + (y_j - \eta_j)^2 / \phi \right\}$$

is maximized with respect to η_j when $\tilde{\eta}_j = y_j$, giving $\log f(y_j; \tilde{\eta}_j, \phi) = -\frac{1}{2} \log(2\pi\phi)$. Therefore the scaled deviance for a model with fitted means $\hat{\eta}_j$ is

$$D = \phi^{-1} \sum_{j=1}^{n} (y_j - \widehat{\eta}_j)^2,$$

which is just the residual sum of squares for the model, divided by ϕ . If $\eta_j=x_j^{\rm T}\beta$ is the correct normal linear model, the distribution of the residual sum of squares is $\phi\chi^2_{n-p}$, so values of D extreme relative to the χ^2_{n-p} distribution call the model into question.

 \square The difference between deviances for nested models A and B in which β has dimensions p and q < p,

$$D_B - D_A = \phi^{-1} \sum_{j=1}^{n} \left\{ (y_j - \widehat{\eta}_j^B)^2 - (y_j - \widehat{\eta}_j^A)^2 \right\} \stackrel{\cdot}{\sim} \chi_{p-q}^2$$

when model B is correct. This distribution is exact for linear models.

If ϕ is unknown, it is replaced by an estimate. The large-sample properties of deviance differences outlined above still apply, though in small samples it may be better to replace the approximating χ^2 distribution by an F distribution with denominator degrees of freedom equal to the degrees of freedom for estimation of ϕ .

Regression Methods

Autumn 2024 - note 1 of slide 107

Model checking

- ☐ Two basic approaches:
 - overall tests either using generic statistic (e.g., chi-squared) or by model expansion (e.g., adding a term and testing for significance);
 - regression diagnostics for detecting a few possibly dodgy observations.
- \square Most widely used diagnostics in the linear model $y=X_{n\times p}\beta+\varepsilon$ are residuals $e_j=y_j-\widehat{y}_j$ and (much better) standardized residuals

$$r_j = \frac{y_j - \widehat{y}_j}{s(1 - h_{jj})^{1/2}}, \quad j = 1, \dots, n,$$

where the leverage h_{jj} is the $j{\rm th}$ diagonal element of the hat matrix $H=X(X^{\rm T}X)^{-1}X^{\rm T}$, and the Cook statistic

$$C_{j} = \frac{1}{ps^{2}} (\widehat{y} - \widehat{y}_{-j})^{\mathrm{T}} (\widehat{y} - \widehat{y}_{-j}) = \frac{r_{j}^{2} h_{jj}}{p(1 - h_{jj})},$$

which measures the effect of deleting the jth case (x_j, y_j) on the fitted model.

Regression Methods

Autumn 2024 - slide 109

Diagnostics in general case

- ☐ Linear model ideas work as approximations (2nd order Taylor series, painful expansions).
- \square Leverage h_{jj} defined as $j{
 m th}$ diagonal element of

$$H = W^{1/2} X (X^{\mathrm{T}} W X)^{-1} X^{\mathrm{T}} W^{1/2},$$

depends in general on $\widehat{\beta}$, unlike in linear model.

☐ Cook statistic is change in deviance

$$C_j = 2p^{-1} \left\{ \ell(\widehat{\beta}) - \ell(\widehat{\beta}_{-j}) \right\} \doteq \frac{h_{jj}}{p(1 - h_{jj})} r_{P_j}^2,$$

where $\widehat{\beta}_{-j}$ is MLE when $j{\rm th}$ case (x_j,y_j) is dropped, and r_{Pj} is standardized Pearson residual (see below).

 \Box There are several types of residual (see next page).

Regression Methods

Residuals in general case

□ Deviance residual:

$$d_j = \operatorname{sign}(\tilde{\eta}_j - \hat{\eta}_j) [2\{\ell_j(\tilde{\eta}_j; \phi) - \ell_j(\hat{\eta}_j; \phi)\}]^{1/2},$$

for which $\sum d_j^2 = D$ is deviance.

 \square Pearson residual: $u_j(\widehat{\beta})/\sqrt{w_j(\widehat{\beta})}$.

☐ Standardized versions

$$r_{Dj} = \frac{d_j}{(1 - h_{jj})^{1/2}}, \quad r_{Pj} = \frac{u_j(\widehat{\beta})}{\{w_j(\widehat{\beta})(1 - h_{jj})\}^{1/2}},$$

and (even better)

$$r_j^* = r_{Dj} + r_{Dj}^{-1} \log(r_{Pj}/r_{Dj}) \stackrel{\cdot}{\sim} N(0, 1)$$

for many models.

These all reduce to usual standardized residual for normal linear model.

Regression Methods

Autumn 2024 - slide 111

Example

Example 19 (Gumbel linear model) Give the components of the IWLS algorithm for fitting the linear model

$$y_j = \beta_0 + \beta_1(x_j - \overline{x}) + \tau \varepsilon_j, \quad j = 1, \dots, n,$$

with Gumbel errors having density function

$$f(y_j; \eta_j, \tau) = \tau^{-1} \exp\left\{-\frac{y_j - \eta_j}{\tau} - \exp\left(-\frac{y_j - \eta_j}{\tau}\right)\right\},$$

where $\tau > 0$ and $\eta_j = \beta_0 + \beta_1(x_j - \overline{x})$; this distribution is natural for maxima; note that τ^2 is not the variance.

Regression Methods

Note to Example 19

 \square As the data are annual maxima, it is more appropriate to suppose that y_j has the Gumbel density

$$f(y_j; \eta_j, \tau) = \tau^{-1} \exp\left\{-\frac{y_j - \eta_j}{\tau} - \exp\left(-\frac{y_j - \eta_j}{\tau}\right)\right\},\tag{12}$$

where τ is a scale parameter and $\eta_j = \beta_0 + \beta_1(x_j - \overline{x})$; here we have replaced the γ s with β s for continuity with the general discussion above.

☐ In this case

$$\ell_j(\eta_j, \tau) = -\log \tau - \frac{y_j - \eta_j}{\tau} - \exp\left(-\frac{y_j - \eta_j}{\tau}\right),\tag{13}$$

and it is straightforward to establish that

$$\frac{\partial \ell_j(\eta_j, \tau)}{\partial \eta_j} = \tau^{-1} \left\{ 1 - \exp\left(-\frac{y_j - \eta_j}{\tau}\right) \right\}, \quad E\left\{-\frac{\partial^2 \ell_j(\eta_j, \tau)}{\partial \eta_j^2}\right\} = \tau^{-2},$$

that $\partial \eta/\partial \beta^{\mathrm{T}} = X$ is the $n \times 2$ matrix whose jth row is $(1, x_j - \overline{x})$, and $W = \tau^{-2}I_n$. Hence (10) becomes $\widehat{\beta} \doteq (X^{\mathrm{T}}X)^{-1}(X\beta + \tau^2 u)$, where the jth element of u is $\tau^{-1}[1 - \exp\{-(y_j - \eta_j)/\tau\}]$.

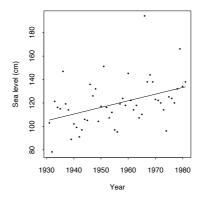
- \square Here it is simplest to fix au, to obtain $\widehat{\beta}$ by iterating (10) for each fixed value of au, and then to repeat this over a range of values of au, giving the profile log likelihood $\ell_p(au)$ and hence confidence intervals for au. Confidence intervals for eta_0 and eta_1 are obtained from the information matrix.
- With starting value chosen to be the least squares estimates of β , and with $\tau=5$, 19 iterations of (10) were required to give estimates and a maximized log likelihood whose relative change was less than 10^{-6} between successive iterations. We then took $\tau=5.5,\ldots,40$, using $\widehat{\beta}$ from the preceding iteration as starting-value for the next; in most cases just three iterations were needed. The left panel of Figure 1 shows a close-up of $\ell_p(\tau)$; its maximum is at $\widehat{\tau}=14.5$, and the 95% confidence interval for τ is (11.9,18.1). The maximum likelihood estimates of β_0 and β_1 are 111.4 and 0.563, with standard errors 2.14 and 0.137; these compare with standard errors 2.61 and 0.177 for the least squares estimates. There is some gain in precision in using the more appropriate model.

Regression Methods

Autumn 2024 - note 1 of slide 112

Venice data

Example 20 (Venice sea level data) The figure below shows annual maximum sea levels in Venice, from 1931–1981. The very large value in 1966 is not an outlier. The fit of a Gumbel model to the data using IWLS gives MLEs (SEs) $\hat{\beta}_0 = 111.4~(2.14)~(\text{cm})$ and $\hat{\beta}_1 = 0.563~(0.137)~(\text{cm/year})$. The standard errors for LSEs are 2.61,~0.177,~larger than for MLEs with Gumbel model — gain in precision through using appropriate model.

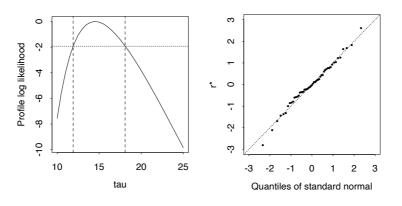


Regression Methods

Autumn 2024 - slide 113

Venice data

Figure 1: Gumbel analysis of Venice data. Left panel: profile log likelihood $\ell_p(\tau) = \max_{\beta} \ell(\beta, \tau)$, with 95% confidence interval (11.9, 18.1) (cm) for τ . Right panel: normal probability plot of residuals r_i^* .



Regression Methods

Summary

 \square For regression problems with independent responses y_j dependent on parameters β through parameter $\eta_j = \eta(x_j; \beta)$, generalise least squares estimation to maximum likelihood estimation, using iterative weighted least squares algorithm: iterate to convergence

$$\widehat{\beta} = (X^{\mathrm{T}}WX)^{-1}X^{\mathrm{T}}Wz, \quad z = X\beta + W^{-1}u,$$

where

$$X_{n \times p} \equiv X(\beta) = \frac{\partial \eta}{\partial \beta^{\mathrm{T}}}, \quad u_{n \times 1} \equiv u(\eta) = \frac{\partial \ell}{\partial \eta}, \quad W_{n \times n} \equiv W(\eta) = -\mathrm{E}\left\{\frac{\partial^2 \ell}{\partial \eta \partial \eta^{\mathrm{T}}}\right\},$$

with ℓ the log likelihood for the data.

- $\hfill \square$ Standard likelihood theory is used for confidence intervals and model comparison.
- \square Linear model diagnostics (residuals, leverage, Cook statistics, ...) generalise to this setting.
- □ Next: generalized linear models (GLMs), wide class of models with exponential family-like response distributions.

Regression Methods

Motivation

- Need to generalise linear model beyond normal responses, e.g. to data with $y \in \{0, 1, ..., m\}$, or $y \in \{0, 1, ...\}$, or y > 0.
- □ Consider **exponential family** response distributions (binomial, Poisson, ...), which have an elegant unifying theory, and encompass many possibilities (in addition to the normal)
- ☐ Basic idea is to build models such that

$$E(y) = \mu, \quad g(\mu) = \eta = x^{\mathrm{T}}\beta,$$

where g is a suitable function, and $y \sim$ exponential family (almost).

☐ Warnings:

- **Don't** confuse Generalized Linear Model (GLM) with General Linear Model (GLM, in older books, the latter is $y = X\beta + \varepsilon$, with $cov(\varepsilon) = \sigma^2 V$ not diagonal);
- **Don't** write $y = \mu + \varepsilon$, since in a GLM the distribution of ε usually depends on μ .

Regression Methods

Autumn 2024 - slide 117

Generalized linear model (GLM)

- □ Normal linear model has three key aspects:
 - structure for covariates: **linear predictor**, $\eta = x^{\mathrm{T}}\beta$;
 - response distribution: $y \sim N(\mu, \sigma^2)$;
 - linear relation $\eta = \mu$ between $\mu = E(y)$ and η .
- ☐ GLM extends last two to
 - Y has density/mass function

$$f(y;\theta,\phi) = \exp\left\{\frac{y\theta - b(\theta)}{\phi} + c(y;\phi)\right\}, \quad y \in \mathcal{Y}, \theta \in \Omega_{\theta}, \phi > 0,$$
(14)

where

- \triangleright \mathcal{Y} is the support of Y,
- \triangleright Ω_{θ} is the parameter space of valid values for $\theta \equiv \theta(\eta)$, and
- \triangleright the **dispersion parameter** ϕ is often known;
- $\eta = g(\mu)$, where g is monotone link function
 - \triangleright the **canonical link** function giving $\eta = \theta = b'^{-1}(\mu)$ has nice statistical properties;
 - \triangleright but a range of link functions are possible for each distribution of Y.

Regression Methods

Examples

Example 21 (GLM density) Show that the moment-generating function of $f(y; \theta, \phi)$ is $M_Y(t) = \exp[\{b(\theta + t\phi) - b(\theta)\}/\phi]$, and deduce that

$$E(Y) = b'(\theta) = \mu, \quad var(Y) = \phi b''(\theta) = \phi b''\{b'^{-1}(\mu)\} = \phi V(\mu);$$

the function $\mu \mapsto V(\mu)$ is known as the variance function.

Example 22 (Poisson distribution) Write the Poisson mass function as a GLM density, and find its canonical link function.

Example 23 (Normal distribution) Write the normal density function as a GLM density, and find its canonical link function.

Regression Methods

Autumn 2024 - slide 119

Note to Example 21

- \square Suppose that Y has a continuous density; if not the argument below is the same, except that integrals are replaced by summations.
- \square Let $\Omega_{\theta} = \{\theta : b(\theta) < \infty\}$. Then

$$M_Y(t) = \mathbb{E}\{\exp(tY)\}\$$

$$= \int e^{ty} \exp\left\{\frac{y\theta - b(\theta)}{\phi} + c(y;\phi)\right\} dy$$

$$= \int \exp\left\{\frac{y(\theta + t\phi) - b(\theta)}{\phi} + c(y;\phi)\right\} dy.$$

If $\theta + t\phi \in \Omega_{\theta}$, then

$$\int \exp\left\{\frac{y(\theta + t\phi) - b(\theta + t\phi)}{\phi} + c(y;\phi)\right\} dy = 1,$$

so

$$M_Y(t) = \mathbb{E}\{\exp(tY)\} = \exp\left[\left\{b(\theta + t\phi) - b(\theta)\right\}/\phi\right].$$

 \square Hence the cumulant-generating function of Y is

$$K_Y(t) = \log M_Y(t) = \{b(\theta + t\phi) - b(\theta)\}/\phi,$$

and differentiating twice with respect to t and setting t=0 yields

$$E(Y) = K'_Y(t)\big|_{t=0} = b'(\theta), \quad \text{var}(Y) = K''_Y(t)\big|_{t=0} = \phi b''(\theta).$$

One can show that $b(\theta)$ is strictly convex on Ω_{θ} . Thus $b'(\theta)$ is a monotonic increasing function of θ , so $b'^{-1}(\cdot)$ exists and is itself monotonic, so $V(\mu) = b''\{b'^{-1}(\mu)\}$ is well-defined.

Regression Methods

Autumn 2024 - note 1 of slide 119

Note to Example 22

The Poisson density may be written as

$$f(y; \mu) = \exp(y \log \mu - \mu - \log y!), \quad y = 0, 1, \dots, \quad \mu > 0,$$

which has GLM form (14) with $\theta = \log \mu$, $b(\theta) = e^{\theta}$, $\phi = 1$, and $c(y; \phi) = -\log y!$. The mean of y is $\mu = b'(\theta) = e^{\theta} = \mu$, and its variance is $b''(\theta) = e^{\theta} = \mu$, so the variance function is linear: $V(\mu) = \mu$.

Regression Methods

Autumn 2024 - note 2 of slide 119

Note to Example 23

The normal density with mean μ and variance σ^2 may be written

$$f(y; \mu, \sigma^2) = \exp\left\{-\frac{(y^2 - 2y\mu + \mu^2)}{2\sigma^2} - \frac{1}{2}\log(2\pi\sigma^2)\right\},$$

so

$$\theta = \mu$$
, $\phi = \sigma^2$, $b(\theta) = \frac{1}{2}\theta^2$, $c(y;\phi) = -\frac{1}{2\phi}y^2 - \frac{1}{2}\log(2\pi\phi)$.

As the first and second derivatives of $b(\theta)$ are θ and 1, we have $V(\mu)=1$; the variance function is constant.

Regression Methods

Autumn 2024 - note 3 of slide 119

Estimation of β

Example 24 (IWLS algorithm) Find the components of the IWLS algorithm for a GLM.

 \square If canonical link is used then $heta_j = x_j^{\mathrm{T}} eta$, so if ϕ is known, then

$$\ell(\beta) = \sum_{j=1}^{n} \left\{ \frac{y_j x_j^{\mathrm{T}} \beta - b(x_j^{\mathrm{T}} \beta)}{\phi} + c(y_j; \phi) \right\}$$
$$= \left\{ y^{\mathrm{T}} X \beta - K(\beta) \right\} / \phi + C(y; \phi),$$

say, which in terms of β is a linear exponential family with

- canonical parameter $\beta_{p\times 1}$
- canonical statistic $(X^{\mathrm{T}}y)_{p\times 1}$,

and many nice properties then hold.

- \square If X is full rank, then $\ell(\beta)$ is (almost always) strictly concave and has a unique maximum in terms of β .
- \square Problem: the maximum may be at infinity in certain (rare) cases—this can arise with binomial responses: beware of $\widehat{\theta}_r \approx \pm 36$.

Regression Methods

Note to Example 24

 \square To compute the quantities needed for the IWLS step $\widehat{\beta} = (X^{\mathrm{T}}WX)^{-1}X^{\mathrm{T}}W(X\beta + W^{-1}u)$, we need

$$X_{n \times p} = \frac{\partial \eta}{\partial \beta^{\mathrm{T}}}, \quad W_{n \times n} = \operatorname{diag}\{\mathrm{E}(-\partial^2 \ell_j/\partial \eta_j^2)\}, \quad u_{n \times 1} = \{\partial \ell_j/\partial \eta_j\},$$

where (with ϕ_i instead of ϕ for generality, see the next slide),

$$\ell_j(\beta) = \left\{ \frac{y_j \theta_j - b(\theta_j)}{\phi_j} + c(y_j; \phi_j) \right\}, \quad b'(\theta_j) = \mu_j, \quad \eta_j = g(\mu_j) = x_j^{\mathrm{T}} \beta.$$

- \square First note that $\partial \eta_j/\partial \beta_r = x_{jr}$, so $X = \partial \eta/\partial \beta^{\mathrm{T}}$ is just a matrix of constants.
- \square We need the first and second derivatives of ℓ_i with respect to η_i , so we write

$$\frac{\partial \ell_j}{\partial \eta_j} = \frac{\partial \mu_j}{\partial \eta_j} \frac{\partial \theta_j}{\partial \mu_j} \frac{\partial \ell_j}{\partial \theta_j},$$

with

$$\frac{\partial \eta_j}{\partial \mu_j} = g'(\mu_j), \quad \frac{\partial \mu_j}{\partial \theta_j} = b''(\theta_j) = V(\mu_j), \quad \frac{\partial \ell_j}{\partial \theta_j} = \frac{y_j - b'(\theta_j)}{\phi_j},$$

which yields

$$u_j = \frac{\partial \ell_j}{\partial \eta_j} = \frac{y_j - b(\theta_j)}{g'(\mu_j)\phi_j V(\mu_j)} = \frac{y_j - \mu_j}{g'(\mu_j)\phi_j V(\mu_j)} = \frac{A(\theta_j)}{B(\theta_j)},$$

say, where E(A) = 0. For the second derivative, we note that

$$\frac{\partial^2 \ell_j}{\partial \eta_j^2} = \frac{\partial}{\partial \eta_j} \frac{\partial \ell_j}{\partial \eta_j} = \left(\frac{\partial \mu_j}{\partial \eta_j} \frac{\partial \theta_j}{\partial \mu_j} \frac{\partial}{\partial \theta_j} \right) \frac{\partial \ell_j}{\partial \eta_j} = \frac{\partial \mu_j}{\partial \eta_j} \frac{\partial \theta_j}{\partial \mu_j} \left\{ \frac{A'(\theta_j)}{B(\theta_j)} - \frac{A(\theta_j) B'(\theta_j)}{B(\theta_j)^2} \right\},$$

and on noting that $B(\theta_j)$ is non-random and $A'(\theta_j) = -b''(\theta_j) = -V(\mu_j)$, we obtain

$$w_j = E\left(-\frac{\partial^2 \ell_j}{\partial \eta_j^2}\right) = \frac{1}{g'(\mu_j)} \frac{1}{V(\mu_j)} \frac{V(\mu_j)}{g'(\mu_j)\phi_j V(\mu_j)} = \frac{1}{g'(\mu_j)^2 \phi_j V(\mu_j)}.$$

Regression Methods

Autumn 2024 - note 1 of slide 120

Note to Example 24, part II

 \square From above we see that the components of the score statistic $u(\beta)$ and the weight matrix $W(\beta)$ may be expressed in terms of components μ_i of the mean vector μ as

$$u_{j} = \frac{\partial \theta_{j}}{\partial \eta_{j}} \frac{\partial \ell_{j}(\theta_{j})}{\partial \theta_{j}} = \frac{y_{j} - \mu_{j}}{g'(\mu_{j})\phi_{j}V(\mu_{j})},$$

$$w_{j} = \left(\frac{\partial \theta_{j}}{\partial \eta_{j}}\right)^{2} \frac{\partial^{2}\ell_{j}(\theta_{j})}{\partial \theta_{j}^{2}} = \frac{1}{g'(\mu_{j})^{2}\phi_{j}V(\mu_{j})},$$
(15)

where $g'(\mu_j) = dg(\mu_j)/d\mu_j$. Thus $\widehat{\beta}$ is obtained by iterative weighted least squares regression of response

$$z = X\beta + g'(\mu)(y - \mu) = \eta + g'(\mu)(y - \mu)$$

on the columns of X using weights (15).

- \square By using y as an initial value for μ and g(y) as an initial value for $\eta=X\beta$, we avoid needing an initial value for β .
- \square It may be necessary to modify y slightly for this initial step. For example if we use the log link for Poisson data, and some y_j equal zero, then we may need to replace them with some small positive value to avoid taking $\log 0$ for some components of the initial $\eta = \log y$.

Regression Methods

Autumn 2024 - note 2 of slide 120

Estimation of ϕ

- \square When ϕ unknown, it is often replaced by $\phi_j = \phi a_j$, with known a_j and a_j^{-1} treated as a weight. Then we replace the scaled deviance by the **deviance** ϕD .
- \square If the model is correct and ϕ is known, then **Pearson's statistic**

$$P = \frac{1}{\phi} \sum_{j=1}^{n} \frac{(y_j - \widehat{\mu}_j)^2}{a_j V(\widehat{\mu}_j)} \stackrel{\cdot}{\sim} \chi_{n-p}^2,$$

analogously to the sum of squares in a linear model, with $E(P) \doteq n - p$.

 \Box The MLE of ϕ can be badly behaved, so usually we prefer the method of moments estimator

$$\hat{\phi} = \frac{1}{n-p} \sum_{j=1}^{n} (y_j - \hat{\mu}_j)^2 / \{a_j V(\hat{\mu}_j)\},$$

which is obtained by solving the equation P=n-p, based on noting that $\mathrm{E}(\chi^2_{n-p})=n-p$.

If the data are sparse (e.g., many small binomial or Poisson counts), then standard asymptotic results are suspect.

Regression Methods

Example: Jacamar data

Table 3: Response (N=not sampled, S = sampled and rejected, E = eaten) of a rufous-tailed jacamar to individuals of seven species of palatable butterflies with artifically coloured wing undersides. Data from Peng Chai, University of Texas.

	Aphrissa	Phoebis	Dryas	Pierella	Consul	Siproeta
	boisduvalli	argante	iulia	luna	fabius	stelenes†
	N/S/E	N/S/E	N/S/E	N/S/E	N/S/E	N/S/E
Unpainted	0/0/14	6/1/0	1/0/2	4/1/5	0/0/0	0/0/1
Brown	7/1/2	2/1/0	1/0/1	2/2/4	0/0/3	0/0/1
Yellow	7/2/1	4/0/2	5/0/1	2/0/5	0/0/1	0/0/3
Blue	6/0/0	0/0/0	0/0/1	4/0/3	0/0/1	0/1/1
Green	3/0/1	1/1/0	5/0/0	6/0/2	0/0/1	0/0/3
Red	4/0/0	0/0/0	6/0/0	4/0/2	0/0/1	3/0/1
Orange	4/2/0	6/0/0	4/1/1	7/0/1	0/0/2	1/1/1
Black	4/0/0	0/0/0	1/0/1	4/2/2	7/1/0	0/1/0

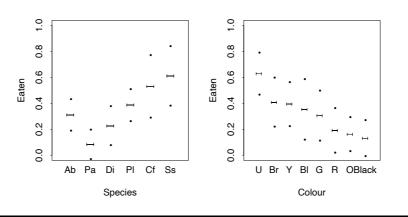
[†] includes *Philaethria dido* also.

Regression Methods

Autumn 2024 - slide 122

Jacamar data

Figure 2: Proportion of butterflies eaten $(\pm 2SE)$ for different species and wing colour.



Regression Methods

Jacamar data

- \square How does a bird respond to the species s and wing colour c of its prey?
- Response has 3 (ordered) categories: not attacked (N), attacked but then rejected (S), attacked and eaten (E)
- \square The data form an 8×6 layout, with a 3-category response in each cell, total m_{cs}
- ☐ Assume that the number in category E (response) is binomial:

$$R_{cs} \sim B(m_{cs}, \pi_{cs}), \quad c = 1, \dots, 8, s = 1, \dots, 6,$$

where c is colour and s is species, with probability that bird attacks and eats butterfly is

$$\pi_{cs} = \frac{\exp(\alpha_c + \gamma_s)}{1 + \exp(\alpha_c + \gamma_s)}, \quad c = 1, \dots, 8, s = 1, \dots, 6,$$

so

- large α_c corresponds to colours that the jacamar likes to eat,
- large γ_s corresponds to species that it likes.
- \square This is a GLM with response $y_{cs}=r_{cs}/m_{cs}$, $\mathrm{E}(y_{cs})=\pi_{cs}$, and canonical (logit) link function

$$\eta = \log{\{\pi/(1-\pi)\}}, \quad \eta_{cs} = \alpha_c + \gamma_s.$$

Regression Methods

Autumn 2024 - slide 124

Jacamar data: Analysis of deviance

Table 4: Deviances and analysis of deviance for models fitted to jacamar data. The lower part shows results for the reduced data, without two outliers.

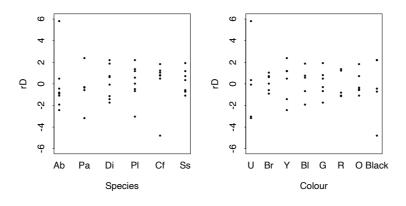
	F	ull data	With	out outliers
Terms	df	df Deviance		Deviance
1	43	134.24	35	73.68
1+Species	38	114.59	31	46.04
1+Colour	36	108.46	28	63.20
1+Species+Colour	31	67.28	24	28.02

Terms	df	Deviance	Terms	df	Deviance
		reduction			reduction
Species (unadj. for Colour)	5	19.64	Species (adj. for Colour)	5	41.18
Colour (adj. for Species)	7	47.31	Colour (unadj. for Species)	7	25.78
Species (unadj. for Colour)	4	27.63	Species (adj. for Colour)	4	35.18
Colour (adj. for Species)	7	18.03	Colour (unadj. for Species)	7	10.48

Regression Methods

Jacamar data: Residuals

Figure 3: Standardized deviance residuals $\emph{r}_{\emph{D}}$ for binomial two-way layout fitted to jacamar data.



Regression Methods

Autumn 2024 - slide 126

Jacamar data: Parameter estimates

Table 5: Estimated parameters and standard errors for the jacamar data, without 2 outliers.

_	Aphrissa	Phoebis	Dryas	Pierella	Consul	Siproe	ta
	boisduvalli	argante	iulia	luna	fabius	stelene	es
_	-1.99 (0.79)	-2.22 (0.85)	-0.56 (0.67)	0.16 (0.54)	_	1.50 (0.	78)
Brown	Yellow	Blue	Green	Red	Or	ange	Black
0.16 (0.73	3) 0.33 (0.68)	-0.53 (0.81)	-0.83 (0.75)	-1.93 (0.88)) -1.94	1 (0.85)	-1.26 (0.86)

- ☐ Interpretation
- $\hfill \square$ Residual deviance: 28.02, with 24 df
- ☐ Pearson statistic: 25.58, with 24 df
- \square Standardized residuals in range -2.03 to 1.96: OK.

Regression Methods

Example: Chimpanzee data

Table 6: Times in minutes taken by four chimpanzees to learn ten words.

Chimpanzee		Word								
	1	2	3	4	5	6	7	8	9	10
1	178	60	177	36	225	345	40	2	287	14
2	78	14	80	15	10	115	10	12	129	80
3	99	18	20	25	15	54	25	10	476	55
4	297	20	195	18	24	420	40	15	372	190

- ☐ A two-way layout.
- ☐ Times vary from 2 to 476 minutes need transformation (e.g., logarithm) if use linear model.

Regression Methods

Autumn 2024 - slide 128

Chimpanzee data

- \square How does learning time depend on word w and chimp c?
- \square Response is continuous and positive, so we try fitting the gamma distribution with mean μ and shape parameter ν , i.e.,

$$f(y; \mu, \nu) = \frac{1}{\Gamma(\nu)} y^{\nu-1} \left(\frac{\nu}{\mu}\right)^{\nu} \exp(-\nu y/\mu), \quad y > 0, \quad \nu, \mu > 0,$$

so dispersion parameter is $\phi=1/\nu$ ($\phi=\nu=1$ for exponential).

□ Possible link functions:

 $\eta = \log \mu$, (log, most common), $\eta = 1/\mu$, (reciprocal, canonical)

☐ Linear model structure:

$$\eta_{cw} = \alpha_c + \gamma_w, \quad c = 1, \dots, 4, w = 1, \dots, 10,$$

but the interpretation of the α_c and γ_w will depend on the link function.

 \square With the log link, the deviances for models 1, 1+Chimp, 1+Word, and 1+Chimp+Word are 60.38, 53.43, 21.19, and 14.97. How many df are there for each model?

Regression Methods

Chimpanzee data: Analysis of deviance

Table 7: Analysis of deviance for models fitted to chimpanzee data.

Term	df	Deviance reduction	Term	df	Deviance reduction
Chimp (unadj. for Word)	3	6.95	Chimp (adj. for Word)	3	6.22
Word (adj. for Chimp)	9	38.46	Word (unadj. for Chimp)	9	39.19

- \square Method of moments estimate is $\widehat{\phi} = 0.432$, so $\widehat{\nu} = 1/\widehat{\phi} = 2.31$.
- \square Use F tests to assess effects of Word and Chimp, for example obtaining

$$\frac{6.22/3}{0.423} = 4.78 \stackrel{.}{\sim} F_{3,27}$$

if there is no difference between the chimps. What is the corresponding statistic for testing differences between words?

Residuals suggest that this model, or one with the inverse link, are both adequate, and both are better than fitting a normal linear model to the log times.

Regression Methods

Autumn 2024 - slide 130

Summary

- ☐ Generalized linear models extend the classical linear model in two ways:
 - the response distribution is (almost) exponential family, so includes binomial, Poisson, gamma and other distributions in addition to the normal;
 - the relation between the linear predictor $\eta=x^{\mathrm{T}}\beta$ and the mean μ is determined by a wide range of possible link functions.
- ☐ Canonical link functions give particularly simple models and are widely used.
- \square Estimates of β are obtained by IWLS, which has a simple form, with no need for initial values.
- \square A simple estimate of the dispersion parameter ϕ is available using the method of moments.
- ☐ Models are compared using the analysis of deviance, which generalises the analysis of variance in the classical linear model.
- □ Standard likelihood theory results are used for inference (standard errors, confidence intervals, etc.)
- ☐ Standard diagnostics (residuals, ...) extend in a natural way to this setting.

Regression Methods

Binary response

 $\hfill \square$ Response Y has Bernoulli distribution with

$$P(Y = 1) = \pi$$
, $P(Y = 0) = 1 - \pi$, $0 < \pi < 1$.

and $E(Y) = \mu = \pi$, $var(Y) = \pi(1 - \pi)$.

- \square Linear link function $\pi = \eta = x^T\beta$ can give $\pi \notin [0,1]$, so not usually a good idea.
- \square Y can be interpreted in terms of a hidden variable/tolerance distribution: let $Z=x^{\mathrm{T}}\gamma+\sigma\varepsilon$, where $\varepsilon\sim F$. Set Y=I(Z>0), and note that

$$\pi = P(Y = 1) = P(x^{\mathrm{T}}\gamma + \sigma\varepsilon > 0) = P(\varepsilon > -x^{\mathrm{T}}\gamma/\sigma) = 1 - F(-x^{\mathrm{T}}\beta),$$

say. Note that $\beta=\gamma/\sigma$ is estimable, but γ and σ are not.

☐ The corresponding link function is given by

$$\eta = x^{\mathrm{T}} \beta = -F^{-1} (1 - \pi) = g(\pi),$$

so different choices of F yield different possible link functions.

Regression Methods

Autumn 2024 - slide 133

Link functions

Tolerance distributions and corresponding link functions for binary data.

Dist	tribution F	Link function				
Logistic	$e^u/(1+e^u)$	Logit	$\eta = \log\{\pi/(1-\pi)\}$			
Normal	$\Phi(u)$	Probit	$\eta = \Phi^{-1}(\pi)$			
Log Weibull	$1 - \exp(-\exp(u))$	Log-log	$\eta = -\log\{-\log(\pi)\}$			
Gumbel	$\exp\{-\exp(-u)\}$	Complementary log-log	$\eta = \log\{-\log(1-\pi)\}$			

- ☐ The logit and probit links are symmetric.
- □ Logit (canonical link) is usual choice, good for medical studies (later), with nice interpretation, but the probit is very similar to it and may be preferred in some cases, for its relation to the normal distribution.
- ☐ The log-log and complementary log-log links are asymmetric.

Regression Methods

Logistic regression

☐ Commonest choice of link function for proportion data is the logit, which gives

$$P(Y = 1) = \pi = \frac{\exp(x^{T}\beta)}{1 + \exp(x^{T}\beta)}, \quad P(Y = 0) = 1 - \pi = \frac{1}{1 + \exp(x^{T}\beta)},$$

leading to a linear model for the log odds of success,

$$\log \left\{ \frac{P(Y=1)}{P(Y=0)} \right\} = \log \left(\frac{\pi}{1-\pi} \right) = x^{\mathrm{T}} \beta, \quad \beta \in \mathbb{R}^{p}.$$

 \square The likelihood for β based on independent responses y_1, \ldots, y_n with covariate vectors x_1, \ldots, x_n and corresponding probabilities π_1, \ldots, π_n is

$$L(\beta) = \prod_{j=1}^{n} \pi_{j}^{y_{j}} (1 - \pi_{j})^{1 - y_{j}} = \dots = \frac{\exp\left(\sum_{j=1}^{n} y_{j} x_{j}^{\mathrm{T}} \beta\right)}{\prod_{j=1}^{n} \left\{1 + \exp\left(x_{j}^{\mathrm{T}} \beta\right)\right\}},$$

which is a regular exponential family with $s(y) = X^{\mathrm{T}}y$ and \log likelihood

$$\ell(\beta) = (X^{\mathrm{T}}y)^{\mathrm{T}}\beta - \sum_{j=1}^{n} \log \left\{ 1 + \exp\left(x_{j}^{\mathrm{T}}\beta\right) \right\}, \quad \beta \in \mathbb{R}^{p},$$

known as the logistic regression model.

Regression Methods

Autumn 2024 - slide 135

Nodal involvement data

Data on nodal involvement: 53 patients with prostate cancer have nodal involvement (r), with five binary covariates age, stage, etc.

$\begin{array}{c ccccccccccccccccccccccccccccccccccc$							
6 1 0 0 0 0 1 4 0 1 1 1 0 0 4 2 1 1 0 0 1 4 0 0 0 0 0 0 3 2 0 1 1 0 1 3 1 1 0 0 0 0 3 0 1 0 0 0 1 3 0 1 0 0 0 0 2 0 1 0 0 0 1 2 1 0 1 0 0 1 2 1 0 0 1 0 0 1 1 1 1 1 1 1 1 1 0 0 1 0 1 1 1 1 0 0 1 0 1 1 1 1 1 1	m	r	age	stage	grade	xray	acid
4 0 1 1 1 0 0 4 2 1 1 0 0 1 4 0 0 0 0 0 0 3 2 0 1 1 0 1 3 1 1 1 0 0 0 3 0 1 0 0 0 0 1 3 0 1 0 1 0 0 1 <	6	5	0	1	1	1	1
4 2 1 1 0 0 1 4 0 0 0 0 0 0 3 2 0 1 1 0 1 3 1 1 1 0 0 0 3 0 1 0 0 0 0 1 3 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 1 0 0 1 <td< td=""><td>6</td><td>1</td><td>0</td><td>0</td><td>0</td><td>0</td><td>1</td></td<>	6	1	0	0	0	0	1
4 0 0 0 0 0 0 3 2 0 1 1 0 1 3 1 1 1 0 0 0 0 3 0 1 0 0 0 0 0 0 2 0 1 0 0 1 0 0 1 0 0 1	4	0	1	1	1	0	0
3 2 0 1 1 0 1 3 1 1 1 0 0 0 3 0 1 0 0 0 0 0 2 0 1 0 0 1 0 0 1 2 1 0 1 0 0 1 0 0 1 2 1 0 0 1 0 0 1 0 0 1 <td< td=""><td>4</td><td>2</td><td>1</td><td>1</td><td>0</td><td>0</td><td>1</td></td<>	4	2	1	1	0	0	1
3 1 1 1 0 0 0 3 0 1 0 0 0 1 3 0 1 0 0 0 0 2 0 1 0 0 1 0 2 1 0 0 1 0 0 1 2 1 0 0 1 0 0 1 2 1 0 0 1 0 0 0 1 1 1 1 1 1 1 1 1 1 0 0 1 0 1 1 1 1 0 0 0 1 1 1	4	0	0	0	0	0	0
3 1 1 1 0 0 0 3 0 1 0 0 0 1 3 0 1 0 0 0 0 2 0 1 0 0 1 0 2 1 0 0 1 0 0 1 2 1 0 0 1 0 0 1 2 1 0 0 1 0 0 0 1 1 1 1 1 1 1 1 1 1 0 0 1 0 1 1 1 1 0 0 0 1 1 1	3	2	0	1	1	0	1
3 0 1 0 0 0 1 3 0 1 0 0 0 0 2 0 1 0 0 1 0 2 1 0 0 1 0 0 1 2 1 0 0 1 0 0 1 0 0 1 <td< td=""><td></td><td>1</td><td>1</td><td>1</td><td>0</td><td>0</td><td>0</td></td<>		1	1	1	0	0	0
2 0 1 0 0 1 0 2 1 0 1 0 0 1 2 1 0 0 1 0 0 1 1 1 1 1 1 1 1 1 0 0 1 0 1 1 0 0 0 0 1 1	3	0	1	0	0	0	1
2 1 0 1 0 0 1 2 1 0 0 1 0 0 1 1 1 1 1 1 1 1 1 0 0 1 0 1 1 0 0 0 0 1 1	3	0	1	0	0	0	0
2 1 0 0 1 0 0 1 1 1 1 1 1 1 : <td>2</td> <td>0</td> <td>1</td> <td>0</td> <td>0</td> <td>1</td> <td>0</td>	2	0	1	0	0	1	0
2 1 0 0 1 0 0 1 1 1 1 1 1 1 : <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td></td>							
2 1 0 0 1 0 0 1 1 1 1 1 1 1 : <td>2</td> <td>1</td> <td>0</td> <td>1</td> <td>0</td> <td>0</td> <td>1</td>	2	1	0	1	0	0	1
: : : : : 1 1 0 0 1 0 1 1 0 0 0 0 1 1		1	0	0	1	0	0
1 0 0 0 0 1 1	1	1	1	1	1	1	1
1 0 0 0 0 1 1	:			:	:	:	
1 0 0 0 0 1 1		•	•	•	•		
	1	1	0	0	1	0	1
1 0 0 0 0 1 0	1	0	0	0	0	1	1
	1	0	0	0	0	1	0

Regression Methods

Deviances for nodal involvement models

Scaled deviances ${\cal D}$ for 32 logistic regression models for nodal involvement data. + denotes a term included in the model.

age	st	gr	xr	ac	df	D	age	st	gr	xr	ac	df	D
					52	40.71	+	+	+			49	29.76
+					51	39.32	+	+		+		49	23.67
	+				51	33.01	+	+			+	49	25.54
		+			51	35.13	+		+	+		49	27.50
			+		51	31.39	+		+		+	49	26.70
				+	51	33.17	+			+	+	49	24.92
+	+				50	30.90		+	+	+		49	23.98
+		+			50	34.54		+	+		+	49	23.62
+			+		50	30.48		+		+	+	49	19.64
+				+	50	32.67			+	+	+	49	21.28
	+	+			50	31.00	+	+	+	+		48	23.12
	+		+		50	24.92	+	+	+		+	48	23.38
	+			+	50	26.37	+	+		+	+	48	19.22
		+	+		50	27.91	+		+	+	+	48	21.27
		+		+	50	26.72		+	+	+	+	48	18.22
			+	+	50	25.25	+	+	+	+	+	47	18.07

Regression Methods

Autumn 2024 - slide 137

Model selection

- ☐ We have 32 competing models, and would like to select the 'best', or a few 'near-best'.
- \square In general we have 2^p models, so automatic selection of some sort is helpful.
- ☐ Could use likelihood ratio tests (differences of deviances) to compare competing models, but this involves many correlated tests, so may lead to spurious results.
- Usually minimise an information criterion, which accounts for the number of parameters in each model, such as

$$AIC \equiv D + 2p$$
, $BIC \equiv D + p \log n$,

where D is the deviance.

- \square Recall their properties, with p fixed and as $n \to \infty$:
 - AIC tends to overfit, i.e., it has a positive probability of choosing a model that is too complex,;
 - BIC applies a stronger penalty, so if the true model is among those fitted, it will choose it with probability one;
 - BIC usually yields less complex models than AIC, but they may predict less well.
- ☐ There are many other information criteria, but these are most used in practice.

Regression Methods

Example: Nodal involvement

☐ Model with lowest AIC has stage, xray, acid:

$$x^{ \mathrm{\scriptscriptstyle T} } \widehat{\beta} = -3.05 + 1.65 I_{\mathtt{stage}} + 1.91 I_{\mathtt{xray}} + 1.64 I_{\mathtt{acid}},$$

where $I_{\mbox{stage}}=1$ indicates that stage takes its higher level, etc.

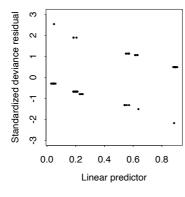
- ☐ Interpretation of this model:
 - for an individual with stage, xray and acid at their lowest levels, the fitted probability of nodal involvement is $e^{-3.05}/(1+e^{-3.05}) \doteq 0.045$ (though there are no such people in the data, so this involves extrapolation);
 - for someone with only $I_{\rm stage}=1$, the odds of nodal involvement are $e^{-3.05+1.65}=e^{-1.4}\doteq0.25$, a probability of 0.2;
 - for someone with $I_{\rm stage}=I_{\rm xray}=I_{\rm acid}=1$, the odds of nodal involvement are $e^{-3.05+1.65+1.91+1.64}\doteq 8.6$, a probability of 0.9;
- \square Problems with interpretation of residual deviance of 19.64: how many df? can amalgamate independent binary responses with same covariates.
- \square Likewise problems with residuals ...

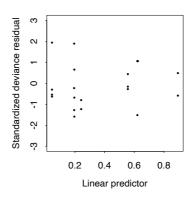
Regression Methods

Autumn 2024 - slide 139

Nodal involvement residuals

Figure 4: Standardized deviance residuals for nodal involvement data, for ungrouped responses (left) and grouped responses (right).





Regression Methods

Su	mmary
	Proportion data are often modelled using the Bernoulli/binomial response distributions.
	Link functions (logit, probit,) have interpretations in terms of underlying continuous variables that have been dichotomized.
	The canonical and most commonly-used link is the logit, and fitting using this yields logistic regression, in which
	 the canonical parameter is the log odds;
	– classical data structures (e.g., the $2 imes 2$ table) have nice interpretations.
	The deviance can be used to compare models (so can AIC, BIC, \dots), but using its absolute value to assess fit can be dangerous (exercise).
	Residuals for binary data are not very informative.

Regression Methods

2.5 Count Data slide 142

Types of count data

- \Box $y \in \{0, 1, 2, \ldots\}$, perhaps with upper bound m, depending on sampling scheme:
 - counts, with no fixed total;
 - m individuals, subdivided into various categories:
 - ▶ nominal response—unordered categories (gender, nationality, . . .)
 - ▶ ordinal response—ordered categories (pain level, spiciness of curry, ...)
- ☐ Simplest models:
 - single unbounded response, or Poisson approximation to binomial, takes $Y \sim \text{Pois}(\mu)$;
 - group of responses (Y_1, \ldots, Y_d) with fixed total $\sum Y_j = m$ has multinomial distribution, probabilities (π_1, \ldots, π_d) and denominator m.
- ☐ Previous examples:
 - Doll and Hill data on smoking had response y Poisson with $\mu = T\lambda(x; \beta)$;
 - Jacamar data had ordinal (?) response N/S/E with total N+S+E fixed—multinomial with d=3

Regression Methods

Autumn 2024 - slide 143

Poisson and multinomial distributions

 \square $Y \sim \operatorname{Pois}(\mu)$ implies that

$$f(y;\mu) = \frac{\mu^y}{y!}e^{-\mu}, \quad y = 0, 1, 2, \dots, \quad \mu > 0.$$

- Exponential family with natural parameter $\theta = \log \mu$, GLM with canonical logarithmic link, $x^{\mathrm{T}}\beta = \eta = \log \mu$.
- \Box If Y is number of events in Poisson process of rate λ observed for period of length T, then $\mu=\lambda T$ and we set $\eta=x^{ \mathrm{\scriptscriptstyle T} }\beta+\log T$
 - offset $\log T$ is fixed part of linear predictor η
- □ If $Y_r \stackrel{\text{ind}}{\sim} \operatorname{Pois}(\mu_r)$, $r = 1, \dots, d$, then the joint distribution of Y_1, \dots, Y_d given $Y_1 + \dots + Y_d = m$ is **multinomial**, with denominator m, and probabilities

$$\pi_1 = \frac{\mu_1}{\sum_{r=1}^d \mu_r}, \quad \dots, \quad \pi_d = \frac{\mu_d}{\sum_{r=1}^d \mu_r}.$$

 \square If $(Y_1,\ldots,Y_d)\sim \mathrm{Mult}(m;\pi_1,\ldots,\pi_d)$, then marginal and conditional distributions, e.g., of

$$(Y_1 + Y_2, Y_3 + Y_4 + Y_5, Y_6, \dots, Y_d), \quad (Y_1, Y_2, Y_4) \mid (Y_3, Y_5, \dots, Y_d),$$

are also multinomial.

Regression Methods

Log-linear and logistic regressions

 $\hfill\Box$ Special case: if d=2, then

$$Y_2 \mid Y_1 + Y_2 = m \quad \sim \quad B\left(m, \pi = \frac{\mu_2}{\mu_1 + \mu_2}\right)$$

 $\Box \quad \text{If } \mu_1 = \exp(\gamma + x_1^{ \mathrm{\scriptscriptstyle T} } \beta) \text{, } \mu_2 = \exp(\gamma + x_2^{ \mathrm{\scriptscriptstyle T} } \beta) \text{, then }$

$$\pi = \frac{\exp(\gamma + x_2^{\mathrm{T}}\beta)}{\exp(\gamma + x_1^{\mathrm{T}}\beta) + \exp(\gamma + x_2^{\mathrm{T}}\beta)} = \frac{\exp\{(x_2 - x_1)^{\mathrm{T}}\beta\}}{1 + \exp\{(x_2 - x_1)^{\mathrm{T}}\beta\}},$$

which corresponds to a logistic regression model for Y_2 with denominator m and probability π .

 \square Can estimate β using log linear model or logistic model—but can't estimate γ from logistic model.

Regression Methods

Autumn 2024 - slide 145

2.6 Poisson Regression

slide 146

>							
>							
	soccer						
	month	day	year	team1	team2	score1	score2
1	Aug	19	2000	Charlton	${\tt ManchesterC}$	4	0
2	Aug	19	2000	Chelsea	WestHam	4	2
3	Aug	19	2000	Coventry	Middlesbr	1	3
4	Aug	19	2000	Derby	${\tt Southampton}$	2	2
5	Aug	19	2000	Leeds	Everton	2	0
6	Aug	19	2000	Leicester	AstonVilla	0	0
7	Aug	19	2000	Liverpool	Bradford	1	0
8	Aug	19	2000	Sunderland	Arsenal	1	0
9	Aug	19	2000	Tottenham	Ipswich	3	1
10	Aug	20	2000	ManchesterU	Newcastle	2	0
11	Aug	21	2000	Arsenal	Liverpool	2	0
12	Aug	22	2000	Bradford	Chelsea	2	0
13	Aug	22	2000	Ipswich	ManchesterU	1	1
14	Aug	22	2000	Middlesbr	Tottenham	1	1
15	Aug	23	2000	Everton	Charlton	3	0
16	Aug	23	2000	ManchesterC	Sunderland	4	2
17	Aug	23	2000	Newcastle	Derby	3	2
18	Aug	23	2000	Southampton	Coventry	1	2
19	Aug	23	2000	WestHam	Leicester	0	1
20	Aug	26	2000	Arsenal	Charlton	5	3

Regression Methods

Premier	League	data
I I CITIICI	LCaguc	uata

- □ 380 soccer matches in English Premier League in 2000–2001 season.
- Data: home score y_{ij}^h and away score y_{ij}^a when team i is at home to team j, for $i,j,=1,\ldots,20$, $i\neq j$.
- ☐ Treat these as Poisson counts with means

$$\mu_{ij}^h = \exp(\Delta + \alpha_i - \beta_j), \quad \mu_{ij}^a = \exp(\alpha_j - \beta_i)$$

where

- Δ represents the home advantage;
- α_i and β_i represent the offensive and defensive strengths of team i.
- ☐ Two possibilities for fitting:
 - Poisson GLM, with 39 parameters;
 - binomial GLM, with 20 parameters.

Regression Methods

Autumn 2024 - slide 148

Premier League data: Analysis of deviance

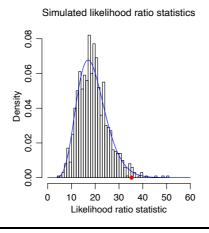
Poi	sson m	nodel	Bind	Binomial model		
Terms	df	Deviance	Terms	df	Deviance	
		reduction			reduction	
Home	1	33.58	Home	1	33.58	
Defence	19	39.21	Team	19	79.63	
Offence	19	58.85				
Residual	720	801.08	Residual	332	410.65	

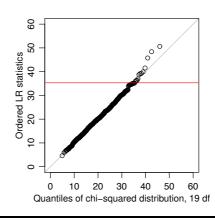
- ☐ There's a strong effect of playing at home, and lots of evidence of differences among the teams—more in offence than defence.
- \square Both residual deviances are a little large, but since the counts are small, we don't expect the large-sample χ^2 distribution to apply well to the residual deviance.
- ☐ Simulations from the fitted model suggest that the residual deviances are not unusually large, so there's no evidence of a lack of fit.

Regression Methods

Premier League data: Null deviance for defence effect

Defence effect deviance (in red) for the Poisson model is large(ish) relative to χ^2_{19} distribution, but the asymptotics seem OK, based on simulations from a model without this effect (i.e., Home + Offence). It seems we can trust asymptotic distributions for differences of deviances, even though the counts are small.



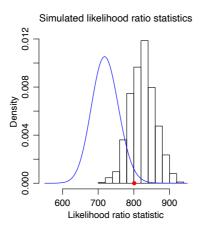


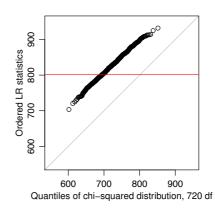
Regression Methods

Autumn 2024 - slide 150

Premier League data: Residual deviance

Residual deviance of 801 (in red) for the Poisson model seems large(ish) relative to χ^2_{720} distribution, but the asymptotics are suspect because most of the counts are small. Comparison of observed deviance with χ^2_{720} distribution shows that 801 is in fact somewhat smaller than average for datasets simulated from the fitted model.





Regression Methods

Premier League data: Estimates

	Overall (δ)	Offensive (α)	Defensive (β)
Manchester United	0.39	0.22	0.15
Liverpool	0.13	0.12	-0.08
Arsenal		0.04	
Chelsea	-0.09	0.08	-0.22
Leeds	-0.10	0.02	-0.17
lpswich	-0.16	-0.10	-0.13
Sunderland	-0.33	-0.31	-0.10
Aston Villa	-0.48	-0.31	-0.15
West Ham	-0.53	-0.33	-0.30
Middlesborough	-0.53	-0.35	-0.17
Charlton	-0.55	-0.21	-0.43
Tottenham	-0.58	-0.28	-0.38
Newcastle	-0.59	-0.35	-0.30
Southampton	-0.60	-0.45	-0.25
Everton	-0.75	-0.32	-0.46
Leicester	-0.77	-0.47	-0.31
Manchester City	-0.90	-0.40	-0.56
Coventry	-0.93	-0.53	-0.52
Derby	-0.93	-0.51	-0.45
Bradford	-1.29	-0.71	-0.62
SEs	0.29	0.20	0.20

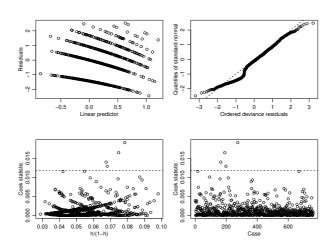
Home advantage: $\widehat{\Delta} = 0.37~(0.07),~\exp(\widehat{\Delta}) = 1.45.$

Regression Methods

Autumn 2024 - slide 152

Premier League data: Assessment of fit

Diagnostic plots for fitted model: residuals against $\widehat{\eta}$ (top left); normal QQ-plot of residuals (top right); Cook statistic C_j against leverage ratio $h_j/(1-h_j)$ (lower left); Cook statistic C_j against case number (lower right).



Regression Methods

Sampling schemes

- A **contingency table** contains individuals (sampling units) cross-classified by various categorical variables.
 - Example: the jacamar data cross-classify butterflies by

6 species
$$\times$$
 8 colours \times 3 fates

for a total of 144 categories, each with its number of butterflies $0, 1, \dots, 14$.

- \Box The sampling scheme underlying a table may fix certain totals. Suppose a pollster wants to find out how people will vote. She might
 - wait in the street for a morning, and get opinions from those people willing to talk to her;
 - wait until she has the views of a fixed number, say m, of people;
 - wait until she has the views of fixed numbers of men and women.

Example 25 Find the likelihoods for each of these sampling schemes, under (unrealistic!) assumptions of independence of voters.

Regression Methods

Autumn 2024 - slide 155

Note to Example 25

- \square An $R \times C$ table arises by randomly sampling a population over a fixed period and then classifying the resulting individuals.
- \square In the first scheme there are no constraints on the row and column totals, and a simple model is that the count in the (r,c) cell, y_{rc} , has a Poisson distribution with mean μ_{rc} . The resulting likelihood is

$$\prod_{r} \left\{ \frac{\mu_{rc}^{y_{rc}}}{y_{rc}!} e^{-\mu_{rc}} \right\};$$

this is simply the Poisson likelihood for the counts in the RC groups.

The pollster may set out with the intention of interviewing a fixed number m of individuals, stopping only when $\sum_{rc} y_{rc} = m$. In this case the data are multinomially distributed, with likelihood

$$\frac{m!}{\prod_{r,c} y_{rc}!} \prod_{r,c} \pi_{rc}^{y_{rc}}, \quad \sum_{r,c} \pi_{rc} = 1,$$

with $\pi_{rc} = \mu_{rc}/\sum_{s,t} \mu_{st}$ the probability of falling into the (r,c) cell.

 \square A third scheme is to interview fixed numbers of men and of women, thus fixing the row totals $m_r = \sum_c y_{rc}$ in advance. In effect this treats the row categories as subpopulations, and the column categories as the response. This yields independent multinomial distributions for each row, and product multinomial likelihood

$$\prod_{r} \left\{ \frac{m_r!}{\prod_c y_{rc}!} \prod_c \pi_{rc}^{y_{rc}} \right\}, \quad \sum_{c} \pi_{1c} = \dots = \sum_{c} \pi_{Rc} = 1,$$

in which $\pi_{rc} = \mu_{rc}/\sum_t \mu_{rt}$.

Regression Methods

Autumn 2024 - note 1 of slide 155

Contingency tables and Poisson response models

- ☐ Multinomial models can be fitted using Poisson errors, provided the appropriate baseline terms are always included in the linear predictor.
- Write the data as two-way layout, with C columns and R rows with fixed totals (e.g., $6 \times 8 = 48$ rows each with 3 columns for the jacamar data).
- \square Consider Poisson model with means $\mu_{rc} = \exp(\gamma_r + x_{rc}^{\mathrm{T}} eta)$:
 - the row parameters $\gamma_1, \ldots \gamma_R$ are **nuisance parameters**, not of interest;
 - we want inference for the parameter of interest, β .
- \square Corresponding multinomial model has fixed row totals m_r and probabilities

$$\pi_{rc} = \frac{\mu_{rc}}{\sum_{d=1}^{C} \mu_{rd}} = \frac{\exp(\gamma_r + x_{rc}^{\mathrm{T}}\beta)}{\sum_{d=1}^{C} \exp(\gamma_r + x_{rd}^{\mathrm{T}}\beta)} = \frac{\exp(x_{rc}^{\mathrm{T}}\beta)}{\sum_{d=1}^{C} \exp(x_{rd}^{\mathrm{T}}\beta)},$$

for r = 1, ..., R, c = 1, ..., C; i.e., one multinomial variable for each row.

 \square The resulting multinomial log likelihood is

$$\ell_{\text{Mult}}(\beta; y \mid m) \equiv \sum_{r=1}^{R} \sum_{c=1}^{C} y_{rc} \log \pi_{rc}$$

$$= \sum_{r=1}^{R} \left\{ \sum_{c=1}^{C} y_{rc} x_{rc}^{\mathsf{T}} \beta - m_r \log \left(\sum_{d=1}^{C} e^{x_{rd}^{\mathsf{T}} \beta} \right) \right\}.$$

Regression Methods

Autumn 2024 - slide 156

Contingency tables and Poisson response models, II

Lemma 26 If parameters τ_r for the row margins are included in the above setup, then we can write

$$\ell_{\text{Poiss}}(\beta, \tau) = \ell_{\text{Poiss}}(\tau; m) + \ell_{\text{Mult}}(\beta; y \mid m).$$

- ☐ Implications:
 - the MLEs of β and τ based on the LHS are the same as those from separate maximisations of the terms on the right:
 - \triangleright $\widehat{\beta}$ equals the MLE for the multinomial model,
 - $\triangleright \quad \widehat{\tau}_r = m_r$
 - the observed and expected information matrices for β, τ are block diagonal.
 - SEs based on the multinomial and Poisson models are equal (exercise).
- \Box General conclusion: inferences on β are the same for multinomial and Poisson models,

provided the parameters associated to the margins fixed under the multinomial model, i.e., the γ_r , are included in the Poisson fit.

Regression Methods

Note to Lemma 26

 \Box The Poisson model has no conditioning, so with $\log \mu_{rc} = \gamma_r + x_{rc}^{\mathrm{T}} \beta$ the log likelihood is

$$\ell_{\text{Poiss}}(\beta, \gamma) \equiv \sum_{r,c} \left(y_{rc} \log \mu_{rc} - \mu_{rc} \right) = \sum_{r=1}^{R} \left(m_r \gamma_r + \sum_{c=1}^{C} y_{rc} x_{rc}^{\mathsf{T}} \beta - e^{\gamma_r} \sum_{c=1}^{C} e^{x_{rc}^{\mathsf{T}} \beta} \right).$$

 \square Now we reparametrise in terms of the row totals $\tau_r = \sum_c \mu_{rc}$, noting that

$$\tau_r = e^{\gamma_r} \sum_{c=1}^C e^{x_{rc}^{\mathrm{T}} \beta}, \quad \gamma_r = \log \tau_r - \log \left\{ \sum_{c=1}^C \exp(x_{rc}^{\mathrm{T}} \beta) \right\},$$

SO

$$\ell_{\text{Poiss}}(\beta, \tau) \equiv \sum_{r=1}^{R} (m_r \log \tau_r - \tau_r) + \sum_{r=1}^{R} \left\{ \sum_{c=1}^{C} y_{rc} x_{rc}^{\text{T}} \beta - m_r \log \left(\sum_{c=1}^{C} e^{x_{rc}^{\text{T}} \beta} \right) \right\},$$

$$= \ell_{\text{Poiss}}(\tau; m) + \ell_{\text{Mult}}(\beta; y \mid m),$$

which is the log likelihood corresponding to

- independent Poisson row totals m_r with means τ_r , and, independent of this,
- the multinomial log likelihood for the contingency table.

Regression Methods

Autumn 2024 - note 1 of slide 157

Jacamar data

Response (N=not sampled, S = sampled and rejected, E = eaten) of a rufous-tailed jacamar to individuals of seven species of palatable butterflies with artifically coloured wing undersides. Data from Peng Chai, University of Texas.

	Aphrissa	Phoebis	Dryas	Pierella	Consul	Siproeta
	boisduvalli	argante	iulia	luna	fabius	stelenes†
	N/S/E	N/S/E	N/S/E	N/S/E	N/S/E	N/S/E
Unpainted	0/0/14	6/1/0	1/0/2	4/1/5	0/0/0	0/0/1
Brown	7/1/2	2/1/0	1/0/1	2/2/4	0/0/3	0/0/1
Yellow	7/2/1	4/0/2	5/0/1	2/0/5	0/0/1	0/0/3
Blue	6/0/0	0/0/0	0/0/1	4/0/3	0/0/1	0/1/1
Green	3/0/1	1/1/0	5/0/0	6/0/2	0/0/1	0/0/3
Red	4/0/0	0/0/0	6/0/0	4/0/2	0/0/1	3/0/1
Orange	4/2/0	6/0/0	4/1/1	7/0/1	0/0/2	1/1/1
Black	4/0/0	0/0/0	1/0/1	4/2/2	7/1/0	0/1/0

† includes *Philaethria dido* also.

Regression Methods

Jacamar data: Models

- \square Let factors F, S, C represent the 3 fates, the 6 species, and the 8 colours.
- \square The models C * S, C * S + F, and C * S + C * F mean we set

$$\log \mu_{csf} = \alpha_{cs}$$
, $\log \mu_{csf} = \alpha_{cs} + \gamma_f$, $\log \mu_{csf} = \alpha_{cs} + \gamma_{cf}$.

 \square The vector of probabilities corresponding to the model with terms C*S is

$$(\pi_{cs1}, \pi_{cs2}, \pi_{cs3}) = \left(\frac{\mu_{cs1}}{\sum_{f=1}^{3} \mu_{csf}}, \frac{\mu_{cs2}}{\sum_{f=1}^{3} \mu_{csf}}, \frac{\mu_{cs3}}{\sum_{f=1}^{3} \mu_{csf}}\right) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}),$$

and that corresponding to the model with terms $C \ast S + F$ is

$$(\pi_{cs1}, \pi_{cs2}, \pi_{cs3}) = \left(\frac{\mu_{cs1}}{\sum_{f=1}^{3} \mu_{csf}}, \frac{\mu_{cs2}}{\sum_{f=1}^{3} \mu_{csf}}, \frac{\mu_{cs3}}{\sum_{f=1}^{3} \mu_{csf}}\right)$$
$$= \frac{1}{e^{\gamma_1} + e^{\gamma_2} + e^{\gamma_3}} \left(e^{\gamma_1}, e^{\gamma_2}, e^{\gamma_3}\right).$$

 \square Exercise: similar computations for C*S+C*F and C*S+C*F+S*F.

Regression Methods

Autumn 2024 - slide 159

Jacamar data: Analysis of deviance

Deviances for log-linear models fitted to jacamar data.

Terms	df	Deviance
C * S	88	259.42
C * S + F	86	173.86
C * S + C * F	72	139.62
C * S + S * F	76	148.23
C * S + C * F + S * F	62	90.66
C * S * F	0	0

- \square The null model C * S is not of interest.
- \square The first model it is sensible to fit is C * S + F.
- The best model seems to be C * S + C * F + S * F, corresponding to independent effects of species and colour, though its deviance is high (but remember the two outlying cells!)

Regression Methods

Pneumoconiosis data

Period of exposure x and prevalence of pneumoconiosis amongst coalminers.

-									
	Period of exposure (years)								
	5.8	15	21.5	27.5	33.5	39.5	46	51.5	
Normal	98	51	34	35	32	23	12	4	
Present	0	2	6	5	10	7	6	2	
Severe	0	1	3	8	9	8	10	5	

☐ Here

Normal < Present < Severe,

so these are ordinal responses with d=3 categories and the total in each group (corresponding to each period of exposure) fixed.

☐ We imagine that the assigned category stems from an underlying continuous variable, even if this cannot be quantified very well.

Regression Methods

Autumn 2024 - slide 162

Models

 \square Assume we have n independent individuals whose responses I_1,\ldots,I_n fall into the set $\{1,\ldots,L\}$, corresponding to L ordered categories, and that

$$\gamma_l = P(I_i \le l) = \pi_1 + \dots + \pi_l, \quad l = 1, \dots, L, \quad \gamma_L = 1,$$

 \square The corresponding likelihood is $\prod_{j=1}^n \pi_{I_j}$, where usually the contribution $\pi_{I_j} \equiv \pi_{I_j}(\eta_j)$ for individual j will depend on covariates x_j through a linear predictor $\eta_j = x_j^{\mathrm{T}} \beta$.

☐ We often want the interpretation of the parameters not to change if we merge adjacent categories, and we can do this using an underlying tolerance distribution, with

$$I_j = l \quad \Leftrightarrow \quad x_j^{\mathrm{T}} \beta + \varepsilon_j \in (\zeta_{l-1}, \zeta_l], \quad \zeta_0 = -\infty < \zeta_1 < \dots < \zeta_{L-1} < \zeta_L = \infty,$$

where the tolerance distribution F of ε_j is often taken to be logistic, giving the **proportional** odds model, in which

$$\pi_l(x_j^{\mathrm{T}}\beta) = \mathrm{P}(\zeta_{l-1} < x_j^{\mathrm{T}}\beta + \varepsilon \le \zeta_l) = F(\zeta_l - x_j^{\mathrm{T}}\beta) - F(\zeta_{l-1} - x_j^{\mathrm{T}}\beta), \quad l = 1, \dots, L;$$

here $\zeta_1,\ldots,\zeta_{L-1}$ are aliased with an intercept β_0 and are not usually of interest.

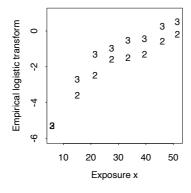
 \square Another standard tolerance distribution is $F(u) = 1 - \exp\{-\exp(u)\}$.

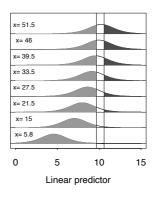
 \square To fit, we just apply IWLS to the multinomial likelihood $\prod_{j=1}^n \pi_{I_j}$.

Regression Methods

Pneumoconiosis data

Pneumoconiosis data analysis, showing how the implied fitted logistic distributions depend on x. Left: plots of empirical logistic transforms for comparing categories 1 with 2+3 and 1+2 with 3; the nonlinearity suggests using $\log x$ as covariate. Right: fitted model, showing probabilities for the three groups with an underlying logistic distribution.





Regression Methods

Autumn 2024 - slide 164

Comments on count data

Log-linear models are mathematically elegant and useful defaults for count data, with close links to logistic regression, based on the relation between the Poisson and multinomial distributions.

☐ Interpretation of log-linear models can be difficult, especially for contingency tables, because marginal and conditional parameters cannot be disentangled.

☐ Other models exist that are less elegant mathematically, but are more interpretable statistically.

Also possible to fit models for ordinal data, using multinomial models and tolerance distribution interpretation used for binomial data.

Regression Methods

Overdispersion

- Often find that discrete response data are more variable than might be expected from a simple Poisson or binomial model, so we see
 - residual deviances larger than expected
 - residuals more variable than expected under the model

but otherwise no evidence of systematic lack of fit

☐ This is **overdispersion**, perhaps due to effect of unmeasured explanatory variables on the responses.

Regression Methods

Autumn 2024 - slide 167

1988	≥14 : 6 3 3 2 2 1		: 8 3 6 8	: 2 11 4	: 3 8 18	: 9 9	: 16 27	: 80	:	:	Year
1988 1 31 80 16 9 3 2 8 6 2 26 99 27 9 8 11 3 3 3 31 95 35 13 18 4 6 3 4 36 77 20 26 11 3 8 2 1989 1 32 92 32 10 12 19 12 2 2 15 92 14 27 22 21 12 1 3 34 104 29 31 18 8 6 4 38 101 34 18 9 15 6 1990 1 31 124 47 24 11 15 8 2 32 132 36 10 9 7 6 3 49 107 51 </th <th>6 3 3 2 2</th> <th></th> <th>8 3 6 8</th> <th>2 11 4</th> <th>3 8 18</th> <th>9</th> <th>16 27</th> <th>80</th> <th>•</th> <th></th> <th></th>	6 3 3 2 2		8 3 6 8	2 11 4	3 8 18	9	16 27	80	•		
2 26 99 27 9 8 11 3 3 3 31 95 35 13 18 4 6 3 4 36 77 20 26 11 3 8 2 1989 1 32 92 32 10 12 19 12 2 2 15 92 14 27 22 21 12 1 3 34 104 29 31 18 8 6 4 38 101 34 18 9 15 6 1990 1 31 124 47 24 11 15 8 2 32 132 36 10 9 7 6 3 49 107 51 17 15 8 9 4 44 153 41 16 11	3 3 2 2		3 6 8	11 4	8 18	9	27		31		
3 31 95 35 13 18 4 6 33 4 36 77 20 26 11 3 8 2 1989 1 32 92 32 10 12 19 12 2 2 15 92 14 27 22 21 12 1 3 34 104 29 31 18 8 6 4 38 101 34 18 9 15 6 1990 1 31 124 47 24 11 15 8 2 32 132 36 10 9 7 6 3 49 107 51 17 15 8 9 4 44 153 41 16 11 6 5 1991 1 41 137 29 33	3 2 2		6 8	4	18			00	<u> </u>	1	1988
4 36 77 20 26 11 3 8 2 1989 1 32 92 32 10 12 19 12 2 2 15 92 14 27 22 21 12 1 3 34 104 29 31 18 8 6 4 38 101 34 18 9 15 6 1990 1 31 124 47 24 11 15 8 2 32 132 36 10 9 7 6 3 49 107 51 17 15 8 9 4 44 153 41 16 11 6 5 1991 1 41 137 29 33 7 11 6	2 2		8			13		99	26	2	
1989 1 32 92 32 10 12 19 12 2 2 15 92 14 27 22 21 12 1 3 34 104 29 31 18 8 6 6 4 38 101 34 18 9 15 6 1990 1 31 124 47 24 11 15 8 2 32 132 36 10 9 7 6 3 49 107 51 17 15 8 9 4 44 153 41 16 11 6 5 1991 1 41 137 29 33 7 11 6	2			3			35	95	31	3	
2 15 92 14 27 22 21 12 ··· 13 34 104 29 31 18 8 6 ··· 4 38 101 34 18 9 15 6 ··· 14 1990 1 31 124 47 24 11 15 8 ··· 2 32 132 36 10 9 7 6 ··· 3 49 107 51 17 15 8 9 ··· 4 44 153 41 16 11 6 5 ··· 1991 1 41 137 29 33 7 11 6 ···		• • •	12		11	26	20	77	36	4	
3 34 104 29 31 18 8 6 4 38 101 34 18 9 15 6 1990 1 31 124 47 24 11 15 8 2 32 132 36 10 9 7 6 3 49 107 51 17 15 8 9 4 44 153 41 16 11 6 5 1991 1 41 137 29 33 7 11 6	1			19	12	10	32	92	32	1	1989
4 38 101 34 18 9 15 6 1990 1 31 124 47 24 11 15 8 2 32 132 36 10 9 7 6 3 49 107 51 17 15 8 9 4 44 153 41 16 11 6 5 1991 1 41 137 29 33 7 11 6	_	• • •	12	21	22	27	14	92	15	2	
1990 1 31 124 47 24 11 15 8 2 32 132 36 10 9 7 6 3 49 107 51 17 15 8 9 4 44 153 41 16 11 6 5 1991 1 41 137 29 33 7 11 6			6	8	18	31	29	104	34	3	
2 32 132 36 10 9 7 6 ··· 3 49 107 51 17 15 8 9 ··· 4 44 153 41 16 11 6 5 ··· 1991 1 41 137 29 33 7 11 6 ···			6	15	9	18	34	101	38	4	
3 49 107 51 17 15 8 9 ··· 4 44 153 41 16 11 6 5 ··· 1991 1 41 137 29 33 7 11 6 ···			8	15	11	24	47	124	31	1	1990
4 44 153 41 16 11 6 5 ··· 1991 1 41 137 29 33 7 11 6 ···			6	7	9	10	36	132	32	2	
1991 1 41 137 29 33 7 11 6			9	8	15	17	51	107	49	3	
			5	6	11	16	41	153	44	4	
2 56 124 39 14 12 7 10 ···			6	11	7	33	29	137	41	1	1991
			10	7	12	14	39	124	56	2	
3 53 175 35 17 13 11 2			2	11	13	17	35	175	53	3	
4 63 135 24 23 12 1				1	12	23	24	135	63	4	
1992 1 71 161 48 25 5					5	25	48	161	71	1	1992
2 95 178 39 6						6	39	178	95	2	
3 76 181 16							16	181	76	3	

Regression Methods

AIDS data

- $\hfill\square$ UK monthly reports of AIDS diagnoses 1983–1992, with reporting delay up to several years!
- ☐ Example of incomplete contingency table (very common in insurance)
- \square Chain-ladder model: number of reports in row j and column k is Poisson, with mean

$$\mu_{jk} = \exp(\alpha_j + \beta_k).$$

☐ Analysis of deviance:

Model	df	Deviance reduction	df	Deviance
			464	14184.3
Time (rows)	37	6114.8	427	8069.5
Delay (cols)	14	7353.0	413	716.5

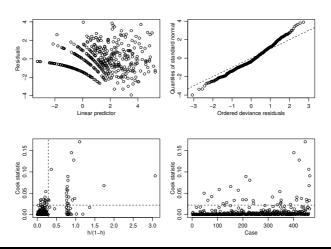
- \square Residual deviance is obviously far too large for a Poisson model to be OK, but the model is also too complex, since we expect smooth variation in the α_j .
- Residuals on next page show no obvious problems, just generic overdispersion.

Regression Methods

Autumn 2024 - slide 169

AIDS data: Assessment of fit

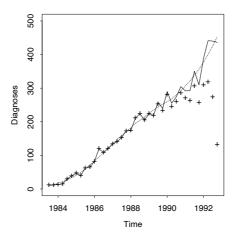
Diagnostic plots for fitted model: residuals against $\widehat{\eta}$ (top left); normal QQ-plot of residuals (top right); Cook statistic C_j against leverage ratio $h_j/(1-h_j)$ (lower left); Cook statistic C_j against case number (lower right).



Regression Methods

AIDS data

- \Box Data (+) and predicted true numbers based on simple Poisson model (solid) and GAM (dots).
- ☐ The Poisson model and data agree up to where data start to be missing.



Regression Methods

Autumn 2024 - slide 171

Dealing with overdispersion

- ☐ Two basic approaches:
 - parametric modelling
 - quasi-likelihood estimation, based only on the variance function

Example 27 (Linear and quadratic variance functions) Suppose that, conditional on $\varepsilon > 0$, $Y \sim \operatorname{Pois}(\mu \varepsilon)$, where $\operatorname{E}(\varepsilon) = 1$ and $\operatorname{var}(\varepsilon) = \xi$. Show that this can lead to either linear or quadratic variance functions, but a lot of data may be needed to distinguish them.

Comparison of variance functions for overdispersed count data. The linear and quadratic variance functions are $V_L(\mu)=(1+\xi_L)\mu$ and $V_Q(\mu)=\mu(1+\xi_Q\mu)$, with $\xi_L=0.5$ and ξ_Q chosen so that $V_L(15)=V_Q(15)$.

μ	1	2	5	10	15	20	30	40	60
Linear	1.5	3.0	7.5	15.0	22.5	30	45	60	90
Quadratic	1.0	2.1	5.8	13.3	22.5	33	60	93	180

Regression Methods

Note to Example 27

Let ε have unit mean and variance $\xi > 0$, and to be concrete suppose that conditional on ε , Y has the Poisson distribution with mean $\mu \varepsilon$. Then

$$E(Y) = E_{\varepsilon} \{ E(Y \mid \varepsilon) \}, \quad var(Y) = var_{\varepsilon} \{ E(Y \mid \varepsilon) \} + E_{\varepsilon} \{ var(Y \mid \varepsilon) \},$$

so the response has mean and variance

$$E(Y) = E_{\varepsilon}(\mu \varepsilon) = \mu, \quad var(Y) = var_{\varepsilon}(\mu \varepsilon) + E_{\varepsilon}(\mu \varepsilon) = \mu(1 + \xi \mu).$$

If on the other hand the variance of ε is ξ/μ , then $\mathrm{var}(Y)=(1+\xi)\mu$. In both cases the variance of Y is greater than its value under the standard Poisson model, for which $\xi=0$. In the first case the variance function is quadratic, and in the second it is linear.

Regression Methods

Autumn 2024 - note 1 of slide 172

Negative binomial model

Example 28 (Negative binomial) In Example 27, if ε is gamma with shape parameter $1/\nu$, show that

$$f(y; \mu, \nu) = \frac{\Gamma(y + \nu)}{\Gamma(\nu)y!} \frac{\nu^{\nu} \mu^{y}}{(\nu + \mu)^{\nu + y}}, \quad y = 0, 1, \dots, \quad \mu, \nu > 0,$$

and that quadratic and linear variance functions are obtained on setting $\nu=1/\xi$ and $\nu=\mu/\xi$ respectively.

The log link function $\log \mu = x^{\mathrm{T}} \beta$ is most natural.

 ξ is estimated by maximum likelihood or through Pearson's statistic.

Example 29 (AIDS data)

- \square MLE $\hat{\xi}_Q = 22.7 (5.5)$
- \square Analysis of Deviance (with $\hat{\xi}_Q$ fixed):

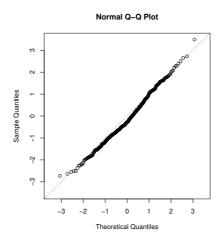
Model	df	Deviance reduction	df	Deviance
			464	7998.3
Time (rows)	37	3582.5	427	4415.8
Delay (cols)	14	3892.2	413	523.6

☐ Still somewhat overdispersed?

Regression Methods

AIDS data: Deviance residuals for NB model

Clear improvement over previous plots, even if not perfect.



Regression Methods

Autumn 2024 - slide 174

Quasi-likelihood

- ☐ Recall two basic assumptions for the linear model:
 - the responses are uncorrelated with means $\mu_j = x_j^{\mathrm{T}} \beta$ and equal variances σ^2 ;
 - in addition to this, the responses are normally distributed.
- ☐ To avoid parametric modelling, we generalise the second-order assumptions, to

$$E(Y_j) = \mu_j, \quad var(Y_j) = \phi_j V(\mu_j), \quad g(\mu_j) = \eta_j = x_j^{\mathrm{T}} \beta,$$

where the variance function $V(\cdot)$ and the link function are taken as known.

 \square We obtain estimates $\tilde{\beta}$ by solving the estimating equation

$$h(\beta; Y) = X^{\mathrm{T}} u(\beta) = \sum_{j=1}^{n} x_j u_j(\beta) = \sum_{j=1}^{n} x_j \frac{Y_j - \mu_j}{g'(\mu_j)\phi_j V(\mu_j)} = 0.$$

□ If the mean structure is correct, then $E(Y_j) = \mu_j$, so $E\{h(\beta; Y)\} = 0$, and under mild conditions $\tilde{\beta}$ is consistent (but maybe not efficient) as $n \to \infty$.

Regression Methods

Quasi-likelihood II

Recall that the general variance of an estimator $\tilde{\beta}$ defined by an estimating equation $h(\beta;Y)_{p\times 1}=0_p$ has sandwich form

$$\mathrm{E}\left\{-\frac{\partial h(\beta;Y)}{\partial \beta^{\mathrm{\scriptscriptstyle T}}}\right\}^{-1}\mathrm{var}\left\{h(\beta;Y)\right\}\mathrm{E}\left\{-\frac{\partial h(\beta;Y)^{\mathrm{\scriptscriptstyle T}}}{\partial \beta}\right\}^{-1}.$$

Lemma 30 If $V(\mu)$ is correctly specified, then $\mathrm{var}(\tilde{\beta}) \doteq (X^{\mathrm{T}}WX)^{-1}$, where W is diagonal with (j,j) element $\{g'(\mu_j)^2\phi_jV(\mu_j)\}^{-1}$.

- \Box If $\phi_j=\phi a_j$, with known $a_j>0$ and unknown $\phi>0$, then we obtain
 - $\tilde{\beta}$ by fitting the GLM with variance function $V(\mu)$ and link $g(\mu)$;
 - standard errors by multiplying the standard errors for this fit by $\widehat{\phi}^{1/2}$, where

$$\widehat{\phi} = \frac{1}{n-p} \sum_{j=1}^{n} \frac{(y_j - \widehat{\mu}_j)^2}{a_j g'(\mu_j)^2 V(\widehat{\mu}_j)}.$$

Regression Methods

Note to Lemma 30

 \square Note first that we can write

$$u_j(\beta) \equiv u_j(\mu_j) = \frac{A_j(\mu_j)}{B_j(\mu_j)},$$

where $A_j(\mu_j)=Y_j-\mu_j$ and $B_j(\mu_j)=g'(\mu_j)\phi_jV(\mu_j)$. Only A_j is random and $\mathrm{E}\{A_j(\mu_j)\}=0$. Hence if we let prime denote derivative with respect to μ_j ,

$$\frac{\partial u_j(\mu_j)}{\partial \mu_j} = \frac{A'_j(\mu_j)}{B_j(\mu_j)} - \frac{A_j(\mu_j)B'_j(\mu_j)}{B_j^2(\mu_j)}$$

has expectation $E\{A'_j(\mu_j)\}/B_j(\mu_j) = -1/B_j(\mu_j)$.

 \square We require $\mathrm{E}\{-\partial h(\beta;Y)/\partial \beta^{\mathrm{T}}\}$ and $\mathrm{var}\{h(\beta;Y)\}$. Now

$$\frac{\partial u_j(\beta)}{\partial \beta^{\mathrm{T}}} = \frac{\partial \eta_j}{\partial \beta^{\mathrm{T}}} \frac{\partial \mu_j}{\partial \eta_i} \frac{\partial u_j(\beta)}{\partial \mu_j} = x_j^{\mathrm{T}} \frac{1}{g'(\mu_j)} u'_j(\mu_j),$$

which gives

$$E\left\{-\frac{\partial h(\beta;Y)}{\partial \beta^{\mathrm{T}}}\right\} = -\sum_{j=1}^{n} x_{j} E\left\{\frac{\partial u_{j}(\beta)}{\partial \beta^{\mathrm{T}}}\right\} = \sum_{j=1}^{n} x_{j} x_{j}^{\mathrm{T}} \frac{1}{g'(\mu_{j})^{2} \phi_{j} V(\mu_{j})} = X^{\mathrm{T}} W X,$$

where W is the $n \times n$ diagonal matrix with jth element $\{g'(\mu_j)^2\phi_jV(\mu_j)\}^{-1}$. Moreover if in addition the variance function has been correctly specified, then $\text{var}(Y_i) = \phi_iV(\mu_i)$, and hence

$$var\{h(\beta; Y)\} = X^{\mathrm{T}} var\{u(\beta)\} X = \sum_{j=1}^{n} x_{j} x_{j}^{\mathrm{T}} \frac{var(Y_{j})}{g'(\mu_{j})^{2} \phi_{j}^{2} V(\mu_{j})^{2}} = X^{\mathrm{T}} W X.$$

Thus the sandwich equals $(X^{\mathrm{T}}WX)^{-1}$.

Had the variance function been wrongly specified, the variance matrix of $\tilde{\beta}$ would have been $(X^{\mathrm{T}}WX)^{-1}(X^{\mathrm{T}}W'X)(X^{\mathrm{T}}WX)^{-1}$, where W' is a diagonal matrix involving the true and assumed variance functions. Only if the variance function has been chosen very badly will this sandwich matrix differ greatly from $(X^{\mathrm{T}}WX)^{-1}$, which therefore provides useful standard errors unless a plot of absolute residuals against fitted means is markedly non-random. In that case the choice of variance function should be reconsidered.

Regression Methods

Autumn 2024 - note 1 of slide 176

Quasi-likelihood III

- Under an exponential family model, $h(\beta; Y)$ is the score statistic, so $\tilde{\beta}$ is the MLE and is efficient (i.e., it has the smallest possible variance in large samples).
- If not, inference is valid provided g and V are correctly chosen, and $\tilde{\beta}$ is optimal among estimators based on linear combinations of the $Y_j \mu_j$, by extending the Gauss–Markov theorem.
- \square In fact we can define a quasi-likelihood Q and its score through

$$Q(\beta;Y) = \sum_{j=1}^{n} \int_{Y_j}^{\mu_j} \frac{Y_j - u}{\phi a_j V(u)} du, \quad h(\beta;Y) = \frac{\partial}{\partial \beta} Q(\beta;Y),$$

and a (quasi-)deviance as $D=-2\phi Q(\beta;Y)$.

 \Box To compare models A, B with numbers of parameters $p_B < p_A$ and deviances $D_B > D_A$, we use the fact that

$$\frac{(D_B - D_A)/(p_A - p_B)}{\widehat{\phi}_A} \quad \dot{\sim} \quad F_{p_A - p_B, n - p_A},$$

if the simpler model B is adequate. This is easy in R.

Regression Methods

Autumn 2024 - slide 177

AIDS example

```
> aids.ql <- glm(y~factor(time)+factor(delay),family=quasipoisson,data=aids.in)</pre>
```

> anova(aids.ql,test="F")
Analysis of Deviance Table

Model: quasipoisson, link: log

Response: y

Terms added sequentially (first to last)

```
Df Deviance Resid. Df Resid. Dev F Pr(>F)

NULL 464 14184.3

factor(time) 37 6114.8 427 8069.5 92.638 < 2.2e-16 ***
factor(delay) 14 7353.0 413 716.5 294.402 < 2.2e-16 ***
---

Signif. codes: 0 '***' 0.001 '**' 0.05 '.' 0.1 ' ' 1
```

Regression Methods

Su	mmary
	Overdispersion is widespread in count and proportion data.
	We deal with it either by
	 parametric modelling, or
	 quasi-likelihood (QL) estimation, which involves assumptions only on the mean-variance relationship.
	QL estimators equal the ML ones, but SEs are inflated by $\widehat{\phi}^{1/2}.$
	(Quasi-)deviance can also be defined, and used for model comparison, with F tests replacing χ^2 tests.

Regression Methods