Regression Methods

Anthony Davison

©2024

http://stat.epfl.ch

1 The Linear Model	2
1.1 Introduction	3
1.2 Inference	17
1.3 Analysis of Variance	33
1.4 Diagnostics	42
1.5 Model Building	59
1.6 Variable Selection	63
1.7 Robustness and Estimating Functions	78
2 General Models	89
2.1 Inference	95
2.2 Model Checking	108
2.3 Generalized Linear Models	116
2.4 Proportion Data	132
2.5 Count Data	142
2.6 Poisson Regression	146
2.7 Contingency Tables	154
2.8 Ordinal Responses	161
2.9 Overdispersion	166
3 Regularisation	180
3.1 Basic Notions	181

3.2 Simple Applications	197
3.3 Lasso	214
3.4 Splines	225
3.5 General Framework	239
3.6 Components of Variance	262
3.7 Linear Mixed Model	274
3.8 Generalized Additive Models	286

slide 2

1.1 Introduction slide 3

Dictionary

- Regression: (statistics) a measure of the relation between the mean value of
 - one variable (e.g., output), denoted y (the response variable) and
 - corresponding values of other variables (e.g., time and cost), denoted x (explanatory variables).
- ☐ The explanatory variables are also called **covariates** or **features** (ML).
- \square We avoid the terms dependent variable (Y) and independent variable (x) used in older books.
- ☐ Questions we try and answer:
 - (description/explanation) how does y depend on x? How much of the variation of y is due to x? Do I need all of x to explain the variation in y?
 - (prediction) what will y be if $x = x_+$?
 - (causation) if I change x, what will happen to y?
- \square The causation question presupposes that we can change (some of) x, which is not always true.

Regression Methods

Autumn 2024 - slide 4

Linear model

 \square Simplest explanation of y in terms of x is linear model:

$$y = g(x) = x_1 \beta_1 + \dots + x_p \beta_p = x^{\mathrm{T}} \beta,$$

where

$$y \in \mathbb{R}, \quad x^{\mathrm{T}} = (x_1, \dots, x_p) \in \mathbb{R}^p, \quad \beta^{\mathrm{T}} = (\beta_1, \dots, \beta_p) \in \mathbb{R}^p.$$

 \square The data consist of n instances/examples/cases (x_i, y_i) for i = 1, ..., n, so

$$y_{n\times 1} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad X_{n\times p} = \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & \\ x_{n1} & \cdots & x_{np} \end{pmatrix}, \quad \beta_{p\times 1} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$

and we write

$$y = X\beta$$
.

- \square Key point: linearity refers to linearity in β , not in terms of elements of X, which might be polynomials, or basis functions, or . . .
- \square Sometimes we can transform to a linear model. For example, the multiplicative expression $y=\gamma x_1^{\beta_1}x_2^{\beta_2}$ becomes

$$\log y = \log \gamma + \beta_1 \log x_1 + \beta_2 \log x_2.$$

Regression Methods

No	tation
	Vectors are column vectors
	We write $X_{n imes p}$ to give the dimensions of a matrix or vector
	a^{T} (row vector) is the transpose of a (column vector)
	$j \in \{1, \dots, n\}$ (or sometimes i) indexes the rows of y (cases/examples)
	$x_j^{\mathrm{\scriptscriptstyle T}}$ is the j th row of X
	$r,s,t,\ldots \in \{1,\ldots,p\}$ indexes the columns of X (covariates/features)
	Roman letters (y, X, z, \dots) denote observed quantities, and may be the realisations of random variables
	Greek letters $(eta, \gamma, heta, \sigma, \ldots)$ denote unknown (often vector) parameters of models
	\widehat{eta} denotes an estimate of eta
	lpha denotes the level of significance tests and confidence intervals
	If Q is scalar (or a row vector) and β is a vector, then $\partial Q/\partial \beta$ denotes the vector (or matrix) the same shape as β with elements $\partial Q/\partial \beta_r$.
	If Q is scalar and β, γ are vectors, then $\partial^2 Q/\partial \beta \partial \gamma^{\mathrm{T}}$ denotes the matrix with (r,s) element $\partial^2 Q/\partial \beta_r \partial \gamma_s$.
	$u\perp v$ means that the vectors u and v are orthogonal (i.e., $u^{ \mathrm{\scriptscriptstyle T} }v=0$); ditto for matrices.
	$Y \perp \!\!\! \perp Z$ means that the random variables Y and Z are independent.

Regression Methods

Autumn 2024 - slide 6

Useful matrix decompositions

 \square Singular value decomposition (SVD): any real matrix X can be written in the form

$$X_{n \times p} = U_{n \times n} D_{n \times p} V_{p \times p}^{\mathrm{T}}$$

where

- $U=(u_1,\ldots,u_n)$ and $V=(v_1,\ldots,v_p)$ are orthogonal (i.e., $U^{\mathrm{T}}U=UU^{\mathrm{T}}=I_n$, $V^{\mathrm{T}}V=VV^{\mathrm{T}}=I_p$) and D is $n\times p$ rectangular diagonal with real diagonal entries (singular values) $d_1\geq\cdots\geq d_m\geq 0$, where $m=\min(n,p)$,
- if one or more $d_j = 0$, then X is singular, and
- the u_i and v_r respectively span the column and row spaces of X.
- $\hfill \square$ The SVD implies that the ranks of X, $X^{ \mathrm{\scriptscriptstyle T} } X$ and $XX^{ \mathrm{\scriptscriptstyle T} }$ are equal and at most m.
- \square Spectral theorem: any real symmetric matrix H can be written as

$$H_{n \times n} = U_{n \times n} D_{n \times n} U_{n \times n}^{\mathrm{T}},$$

where

- $D = \operatorname{diag}(d_1, \dots, d_n)$ contains the eigenvalues of H;
- $-\ U$ is an orthogonal matrix whose columns are the corresponding eigenvectors; and
- if H is positive semi-definite then $d_1 \ge \cdots \ge d_n \ge 0$.

Regression Methods

Least squares fit

☐ Assume that

$$y = X\beta$$

and find the 'best fit' by choosing β to minimise the (squared) Euclidean distance between y and $X\beta$, i.e., the sum of squares

$$||y - X\beta||^2 = (y - X\beta)^{\mathrm{T}}(y - X\beta) = \sum_{j=1}^{n} (y_j - x_j^{\mathrm{T}}\beta)^2.$$

- \square In vector space terms, $y \in \mathbb{R}^n$ and $X\beta \in \operatorname{span}(X) \subset \mathbb{R}^n$.
- \square The 'best fit' vector \widehat{y} is the vector in $\mathrm{span}(X)$ closest to y; Pythagoras' theorem (sketch) gives $\widehat{y} \perp (y \widehat{y})$ (but see below).
- \square We call $\widehat{y} \in \mathbb{R}^n$ the fitted value(s) and $e = y \widehat{y} \in \mathbb{R}^n$ the residual (vector).

Lemma 1 When X has rank p and $n \ge p$ then $\widehat{y} = X\widehat{\beta} = Hy$, where

$$\widehat{\beta} = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y, \quad H = X(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}.$$

The 'hat matrix' H has rank p, is symmetric and idempotent, and satisfies HX = X: it gives the orthogonal projection of \mathbb{R}^n onto $\mathrm{span}(X)$.

Regression Methods

Autumn 2024 - slide 8

Note to Lemma 1

- \square If X has rank p, so too does the $p \times p$ matrix X^TX , which is therefore invertible.
- \Box The sum of squares

$$Q = (y - X\beta)^{\mathrm{T}}(y - X\beta) = y^{\mathrm{T}}y - \beta^{\mathrm{T}}X^{\mathrm{T}}y - y^{\mathrm{T}}X\beta + \beta^{\mathrm{T}}X^{\mathrm{T}}X\beta = y^{\mathrm{T}}y - 2y^{\mathrm{T}}X\beta + \beta^{\mathrm{T}}X^{\mathrm{T}}X\beta$$

has first and second derivatives (respectively a $p \times 1$ vector and $p \times p$ matrix)

$$\frac{\partial Q}{\partial \beta} = -2X^{ \mathrm{\scriptscriptstyle T} } y + 2X^{ \mathrm{\scriptscriptstyle T} } X \beta, \quad \frac{\partial^2 Q}{\partial \beta \partial \beta^{ \mathrm{\scriptscriptstyle T} }} = 2X^{ \mathrm{\scriptscriptstyle T} } X$$

with respect to β . Setting $\partial Q/\partial\beta=0$ implies that $(X^{\mathrm{T}}X)\beta=X^{\mathrm{T}}y$, and as $X^{\mathrm{T}}X$ is invertible we can write

$$\widehat{\beta} = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y, \quad \widehat{y} = X\widehat{\beta} = X(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y = Hy,$$

say. The matrix $X^{\mathrm{T}}X$ is positive definite, so $(y-X\beta)^{\mathrm{T}}(y-X\beta)$ is minimised at $\widehat{\beta}$.

- The $n \times n$ 'hat matrix' H (which 'puts a hat' on y) satisfies $H^{\mathrm{T}} = H$, $H^2 = H$, so it is symmetric and idempotent, i.e., its eigenvalues equal 0 or 1, and their multiplicities must be n-p and p, as its rank is p. H is the matrix that projects \mathbb{R}^n orthogonally onto the span of the columns of X, $\mathrm{span}(X)$.
- The inner product between \widehat{y} and $y-\widehat{y}$ equals zero, because $\widehat{y}=Hy$, $y-\widehat{y}=(I-H)y$, and $\widehat{y}^{\mathrm{T}}(y-\widehat{y})=y^{\mathrm{T}}H^{\mathrm{T}}(I-H)y=y^{\mathrm{T}}(H-H)y=0$. Hence \widehat{y} and $y-\widehat{y}$ are orthogonal.
- \square Clearly $HX=X(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}X=X$, so $H(X\beta)=X\beta$ for any $\beta\in\mathbb{R}^p$, i.e., a vector in $\mathrm{span}(X)$ is left unchanged by multiplication by H.

Regression Methods

Autumn 2024 - note 1 of slide 8

Analysis of variance I

Lemma 2 Let $X_{n\times p}=(X_0,X_1,\ldots,X_R)$ have rank p, where $p\le n$, and let H_r denote the projection matrices formed using (X_0,\ldots,X_r) , for $r=0,\ldots,R$; hence $H_R=H$. Define $P_r=H_r-H_{r-1}$ for $r=1,\ldots,R$ and $P_{R+1}=I-H$. Then (i) $H_rH_s=H_r$ whenever $r\le s$, (ii) $H_0P_r=0$ for any r, and (iii) the matrices P_r are symmetric and idempotent, with $P_rP_s=0$ when $r\ne s$.

 \square In the setup of Lemma 2 suppose we fit the models with projection matrices $H_0, \dots, H_R = H$ and corresponding fitted values $\widehat{y}_r = H_r y$. Then

$$y = \widehat{y}_0 + (\widehat{y}_1 - \widehat{y}_0) + \dots + (\widehat{y}_R - \widehat{y}_{R-1}) + (y - \widehat{y}_R)$$

= $H_0 y + (H_1 - H_0) y + \dots + (H_R - H_{R-1}) y + (I - H) y$
= $H_0 y + P_1 y + \dots + P_R y + P_{R+1} y$,

and Lemma 2 implies that the terms on the RHS are orthogonal, i.e.,

$$(H_0 y)^{\mathrm{T}}(P_r y) = 0, \quad (P_s y)^{\mathrm{T}}(P_r y) = 0, \quad r \neq s.$$

☐ Hence Pythagoras' theorem gives the analysis of variance (ANOVA) decomposition

$$||y||^2 = ||\widehat{y}_0||^2 + \sum_{r=1}^R ||\widehat{y}_r - \widehat{y}_{r-1}||^2 + ||y - \widehat{y}||^2.$$

Regression Methods

Autumn 2024 - slide 9

Note to Lemma 2

 \square (i) Let $\mathcal{V}_0 \subset \cdots \subset \mathcal{V}_R$ denote the linear spaces onto which \mathbb{R}^n is projected by $H_0, \ldots, H_R = H$, and suppose that $r \leq s$. Now $H_r y \in \mathcal{V}_r$ for any $y \in \mathbb{R}^n$, so as $\mathcal{V}_r \subset \mathcal{V}_s$, $H_r y \in \mathcal{V}_s$. Hence $H_s H_r y = H_r y$ for any $y \in \mathbb{R}^n$, so $H_s H_r = H_r$. This implies that

$$H_s H_r = H_r = H_r^{\mathrm{T}} = (H_s H_r)^{\mathrm{T}} = H_r^{\mathrm{T}} H_s^{\mathrm{T}} = H_r H_s, \quad s \ge r.$$

- \Box (ii) For $r=1,\dots,R$, (i) yields $H_0P_r=H_0H_r-H_0H_{r-1}=H_0-H_0=0,$ and $H_0P_{R+1}=H_0(I-H_R)=0.$
- \square (iii) The matrices P_1, \ldots, P_R are symmetric because

$$P_r^{\mathrm{T}} = (H_r - H_{r-1})^{\mathrm{T}} = H_r^{\mathrm{T}} - H_{r-1}^{\mathrm{T}} = H_r - H_{r-1} = P_r,$$

and idempotent because (i) gives

$$P_r^2 = (H_r - H_{r-1})(H_r - H_{r-1})$$

$$= H_r^2 - H_r H_{r-1} - H_{r-1} H_r + H_{r-1}^2$$

$$= H_r - H_{r-1} - H_{r-1} + H_{r-1}$$

$$= H_r - H_{r-1} = P_r.$$

Moreover if $r < s \le R$, then

$$P_r P_s = (H_r - H_{r-1})(H_s - H_{s-1})$$

$$= H_r H_s - H_r H_{s-1} - H_s H_{r-1} + H_{r-1} H_{s-1}$$

$$= H_r - H_r - H_{r-1} + H_{r-1}$$

$$= 0.$$

The corresponding results for P_{R+1} are equally easy to check.

Regression Methods

Autumn 2024 - note 1 of slide 9

Analysis of variance II

 \square Usually $X_0=1_n$; then $\widehat{y}_0=1_n(1_n^{\mathrm{T}}1_n)^{-1}1_n^{\mathrm{T}}y=\overline{y}1_n$ and

$$||y||^2 - ||\widehat{y}_0||^2 = \sum_{j=1}^n y_j^2 - \sum_{j=1}^n \overline{y}^2 = \sum_{j=1}^n (y_j - \overline{y})^2,$$

equals n times the empirical variance of y_1, \ldots, y_n . Hence

$$\sum_{j=1}^{n} (y_j - \overline{y})^2 = ||y||^2 - ||\widehat{y}_0||^2 = \sum_{r=1}^{R} ||\widehat{y}_r - \widehat{y}_{r-1}||^2 + ||y - \widehat{y}||^2$$

decomposes ('analyses') the variability of y around its average \overline{y} into

- the contributions $\|\widehat{y}_r \widehat{y}_{r-1}\|^2$ due to adding the columns of X_r to X_0, \ldots, X_{r-1} ,
- the **residual sum of squares** $||y \hat{y}||^2$ left after fitting $X = (X_0, \dots, X_R)$.
- \square Large $\|\widehat{y}_r \widehat{y}_{r-1}\|^2$ implies that X_r explains a lot of the variation of y even after allowing for that explained by X_0, \ldots, X_{r-1} .
- ☐ The
 - **degrees of freedom** of a fit is the rank ν_r of the corresponding H_r , and the
 - residual degrees of freedom is $n \nu_R = n p$.

Regression Methods

Autumn 2024 - slide 10

Terms

 $\ \square$ A constant column $X_0=1_n$ is almost always present in the design matrix, so

$$X\beta = \begin{pmatrix} 1_n & X_1 & \cdots & X_R \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_R \end{pmatrix} = 1_n \beta_0 + X_1 \beta_1 + \cdots + X_R \beta_R,$$

where the matrices X_1, \ldots, X_R , the **terms**, are successively included.

 \square The baseline model with only 1_n has fitted value and residual vector

$$\widehat{y}_0 = \overline{y}1_n, \quad y - \widehat{y}_0 = y - \overline{y}1_n.$$

- \square Starting from the baseline we ask which terms lead to large reductions in the residual sum of squares, i.e., best explain the variation of y.
- \square The successive residual degrees of freedom, i.e., the ranks of the matrices $I-H_r$, are

$$n-1 = n - \nu_0 > n - \nu_1 > \cdots > n - \nu_R$$
.

 \square When the columns of X_{r+1} depend linearly on those of $1_n, X_1, \ldots, X_r$, we have $\nu_{r+1} = \nu_r$, so inclusion of X_{r+1} does not change the fitted value or improve the fit.

Regression Methods

Model formulae

 \square A mean vector such as $1_n\beta_0+X_1\beta_1+X_2\beta_2$ is often written as the right-hand side of

$$y \sim X1 + X2$$

where

- the columns of 1s is (silently) included first by default,
- X1 and X2 represent the vector subspaces of \mathbb{R}^n generated by the corresponding terms, and
- + represents addition of vector subspaces.
- \square Software generally drops any column of a design matrix that is linearly dependent on previous columns, and this affects which elements of β can be estimated and the meaning of estimates corresponding to later columns.
- \square Carefully choosing the order of terms in a model can give easily interpreted estimates of the parameters of interest for example, if X_2 is full-rank and a column of 1s lies in $\mathrm{span}(X_1) + \mathrm{span}(X_2)$ then

$$y \sim X1 + X2$$
, $y \sim X2 + X1 - 1$,

span the same linear space but the second estimates the parameters of β_2 (unadjusted for the mean) and the parameters of β_1 , adjusted for the presence of X_2 .

Regression Methods

Autumn 2024 - slide 12

ANOVA

Terms	Residual df	Residual SS	Term added	Reduction in	Reduction in SS	Mean square
				residual df		
1_n	$n - \nu_0 = n - 1$	SS_0				
$1_n, X_1$	$n-\nu_1$	SS_1	X_1	$\nu_1 - \nu_0$	$SS_0 - SS_1$	$\frac{\mathrm{SS}_0 - \mathrm{SS}_1}{\nu_1 - \nu_0}$
$1_n, X_1, X_2$	$n-\nu_2$	SS_2	X_2	$\nu_2 - \nu_1$	$SS_1 - SS_2$	$\frac{\substack{\nu_1-\nu_0\\\mathrm{SS}_1-\mathrm{SS}_2}}{\substack{\nu_2-\nu_1}}$
:	:	:	i i	:	i:	÷
$1_n, X_1, \ldots, X_R$	$n - \nu_R = n - p$	SS_R	X_R	$\nu_R - \nu_{R-1}$	$SS_{R-1} - SS_R$	$\frac{SS_{R-1}-SS_R}{\nu_R-\nu_{R-1}}$

 \square The sum of squares when including terms $1_n, X_1, \ldots, X_r$ is

$$SS_r = \|y - \widehat{y}_r\|^2.$$

 \square The 'mean square' for term X_r ,

$$MS_r = \frac{SS_{r-1} - SS_r}{\nu_r - \nu_{r-1}}$$

is the average reduction in SS_r per degree of freedom when X_r is added to the model.

☐ Usually show only the RHS of the table and the bottom line of its LHS (next slide).

Regression Methods

ANOVA table

Term added	df	Reduction in SS	Mean square
X_1	$\nu_1 - \nu_0$	$SS_0 - SS_1$	$MS_1 = (SS_0 - SS_1)/(\nu_1 - \nu_0)$
X_2	$\nu_2 - \nu_1$	$SS_1 - SS_2$	$MS_2 = (SS_1 - SS_2)/(\nu_2 - \nu_1)$
:	:	:	<u>:</u>
X_R	$\nu_R - \nu_{R-1}$	$SS_{R-1} - SS_R$	$MS_R = (SS_{R-1} - SS_R)/(\nu_R - \nu_{R-1})$
Residual	$n-\nu_R$	SS_R	$MS_{Res} = SS_R/(n - \nu_R)$

- Used to screen which terms give the largest reductions, comparing MS_r with the residual mean square MS_{Res} .
- \square Judge 'significance' of reductions relative to residual using F-tests (later).
- ☐ Problem: the order of adding terms matters, so there is no unique reduction in general.

Regression Methods

Autumn 2024 - slide 14

Coefficient of determination

 \square Coefficient of determination \mathbb{R}^2 measures reduction in variance of y as

$$R^2 = \frac{\|\widehat{y} - \overline{y}1_n\|^2}{\|y - \overline{y}1_n\|^2} = \frac{\{(H - H_0)y\}^{\mathrm{T}}(H - H_0)y}{\{(I - H_0)y\}^{\mathrm{T}}(I - H_0)y} = \frac{y^{\mathrm{T}}(H - H_0)y}{y^{\mathrm{T}}(I - H_0)y},$$

where H_0 and H are the hat matrices for regression on 1_n and X, and $1_n \in \text{span}(X)$.

- \square $R^2 \in [0,1]$ is the squared empirical correlation between y and \hat{y} , so $R^2 \approx 1$ implies that most of the variation in y is explained by \hat{y} .
- ☐ There is a geometric interpretation, as the terms on the right of

$$(I_n - H_0)y = (I_n - H)y + (H - H_0)y$$

are orthogonal (check this).

 \square Adding columns to X must increase R^2 , unlike the adjusted R^2 ,

$$R_a^2 = R^2 + (1 - R^2) \frac{n-1}{n-p}.$$

 \square If $1_n \notin \operatorname{span}(X)$, use

$$R_0^2 = \frac{\widehat{y}^{\mathrm{T}} \widehat{y}}{y^{\mathrm{T}} y}, \quad R_{0,a}^2 = R_0^2 + (1 - R_0^2) \frac{n}{n - p}.$$

9

Regression Methods

L.	റ	m	m	e	n	ts

- \square We have supposed that $X_{n \times p}$ has rank p:
 - if X is rank-deficient, then a least squares algorithm usually drops columns that lie in the span of preceding ones, but care is needed to construct X so that the resulting $\widehat{\beta}$ is easy to interpret;
 - if X is nearly rank-deficient, then regularisation may be needed. More later ...
- ☐ Everything so far as purely numerical:
 - least squares estimation is a numerical technique for using X to approximate y;
 - $\widehat{y} = X\widehat{\beta}$ is the resulting approximation, which lies in span(X);
 - $-\widehat{\beta}$ gives the coefficients of the columns of X for the best approximation;
 - the coefficient of determination \mathbb{R}^2 measures how much of the overall variation of y was explained by X; and
 - the ANOVA decomposition summarises how much of the variation in y is explained by different subsets of columns of X (terms).
- \square For statistics we need to add some distributional assumptions ... shortly ...
- ☐ First some reminders . . .

Regression Methods

Autumn 2024 - slide 16

1.2 Inference slide 17

Reminder: Moment-generating function

Definition 3 The moment-generating function (MGF) of a random vector $Y_{n\times 1}$ is

$$M_Y(t) = E(e^{t^T Y}) = E(e^{\sum_{j=1}^n t_j Y_j}), \quad t \in \mathcal{T} = \{t \in \mathbb{R}^n : M_Y(t) < \infty\},$$

and the cumulant-generating function of Y is $K_Y(t) = \log M_Y(t)$, $t \in \mathcal{T}$.

Then

- \square $0 \in \mathcal{T}$, so $M_Y(0) = 1$ and $K_Y(0) = 0$;
- $\hfill\Box$ if ${\mathcal T}$ contains an open set, then

$$\mu = \mathrm{E}(Y) = K_Y'(0) = \left. \frac{\partial K_Y(t)}{\partial t} \right|_{t=0}, \quad \Omega = \mathrm{var}(Y) = \left. \frac{\partial^2 K_Y(t)}{\partial t \partial t^\mathrm{T}} \right|_{t=0};$$

 \square if \mathcal{A},\mathcal{B} are disjoint subsets of $\{1,\ldots,n\}$ and $Y_{\mathcal{A}}$ denotes the sub-vector of Y containing $\{Y_j:j\in\mathcal{A}\}$, etc., then $Y_{\mathcal{A}}\perp\!\!\!\perp Y_{\mathcal{B}}$ if and only if

$$M_Y(t) = \mathrm{E}(e^{t_{\mathcal{A}}^{\mathrm{T}} Y_{\mathcal{A}} + t_{\mathcal{B}}^{\mathrm{T}} Y_{\mathcal{B}}}) = M_{Y_{\mathcal{A}}}(t_{\mathcal{A}}) M_{Y_{\mathcal{B}}}(t_{\mathcal{B}}), \quad t \in \mathcal{T};$$

- \square the MGF of Y_A equals $M_Y(t)$ evaluated with $t_B=0$;
- if we recognise an MGF, then we know the probability distribution that gave it.

Regression Methods

Reminder: Multivariate normal distribution

A random variable $Y_{n\times 1}$ with real components has the multivariate normal distribution, $Y \sim \mathcal{N}_n(\mu,\Omega)$, if $a^{\mathrm{T}}Y \sim \mathcal{N}(a^{\mathrm{T}}\mu,a^{\mathrm{T}}\Omega a)$ for every constant vector $a_{n\times 1}$, and then

(a) Ω is symmetric semi-positive definite with real components and

$$E(Y) = \mu_{n \times 1}, \quad var(Y) = \Omega_{n \times n}, \quad M_Y(t) = \exp(t^T \mu + \frac{1}{2} t^T \Omega t), \quad t \in \mathbb{R}^n,$$

where we call μ the **mean vector** and Ω the **(co)variance matrix** of X;

- (b) for any constants $a_{m\times 1}$ and $B_{m\times n}$, $a+BY\sim \mathcal{N}_m\left(a+B\mu,B\Omega B^{\mathrm{T}}\right)$;
- (c) if $Y^{\mathrm{T}}=(Y_1^{\mathrm{T}},Y_2^{\mathrm{T}})$, where Y_1 is $m\times 1$, and μ and Ω are partitioned correspondingly, then the marginal and conditional distributions of Y_1 are also multivariate normal:

$$Y_1 \sim \mathcal{N}_m(\mu_1, \Omega_{11}), \quad Y_1 \mid Y_2 = y_2 \sim \mathcal{N}_m \left\{ \mu_1 + \Omega_{12} \Omega_{22}^{-1} (y_2 - \mu_2), \Omega_{11} - \Omega_{12} \Omega_{22}^{-1} \Omega_{21} \right\};$$

- (d) $Y_1 \perp \!\!\!\perp Y_2$ iff $\Omega_{12} = 0$, and $a + BY \perp \!\!\!\perp c + DY$ iff $B\Omega D^{\mathrm{T}} = 0$;
- (e) if $Y_1, \ldots, Y_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, then $Y_{n \times 1} \sim \mathcal{N}_n(\mu 1_n, \sigma^2 I_n)$; and finally,
- (f) Y has a density on \mathbb{R}^n iff Ω is positive definite (i.e., has rank n), and then

$$f(y;\mu,\Omega) = \frac{1}{(2\pi)^{n/2}|\Omega|^{1/2}} \exp\left\{-\frac{1}{2}(y-\mu)^{\mathrm{T}}\Omega^{-1}(y-\mu)\right\}, \quad y \in \mathbb{R}^n.$$
 (1)

Regression Methods

Note: Multivariate normal distribution

- (a) Let e_i denote the *n*-vector with 1 in the *j*th place and zeros everywhere else.
- \square Then $Y_j = e_i^{\mathrm{T}} Y \sim N(\mu_j, \omega_{jj})$, giving the mean and variance of Y_j .
- \square Now $\operatorname{var}(Y_j + Y_k) = \operatorname{var}(Y_j) + \operatorname{var}(Y_k) + 2\operatorname{cov}(Y_j, Y_k)$, and

$$Y_j + Y_k = (e_j + e_k)^{\mathrm{T}} Y \sim \mathcal{N}(\mu_j + \mu_k, \omega_{jj} + \omega_{kk} + 2\omega_{jk}),$$

which implies that $cov(Y_i, Y_k) = \omega_{ik} = \omega_{ki}$. This gives the mean and covariance matrix of Y.

- - (b) The MGF of a + BY equals

$$\begin{split} \mathbf{E} \left[\exp\{t^{\mathsf{T}}(a + BY)\} \right] &= \mathbf{E} \left[\exp\{t^{\mathsf{T}}a + (B^{\mathsf{T}}t)^{\mathsf{T}}Y)\} \right] \\ &= e^{t^{\mathsf{T}}a} M_Y(B^{\mathsf{T}}t) \\ &= \exp\{t^{\mathsf{T}}a + (B^{\mathsf{T}}t)^{\mathsf{T}}\mu + \frac{1}{2}(B^{\mathsf{T}}t)^{\mathsf{T}}\Omega(B^{\mathsf{T}}t)\} \\ &= \exp\left\{t^{\mathsf{T}}(a + B\mu) + \frac{1}{2}t^{\mathsf{T}}(B\Omega B^{\mathsf{T}})t\right\}, \end{split}$$

which is the MGF of the $\mathcal{N}_m(a+B\mu,B\Omega B^{\mathrm{T}})$ distribution. Hence linear combinations of normal variables are themselves normal.

(c) Write $Y^{\mathrm{T}}=(Y_1^{\mathrm{T}},Y_2^{\mathrm{T}})$ and partition μ and Ω conformally. Then

$$M_Y(t) = \exp\left\{t_1^{\mathrm{T}} \mu_1 + t_2^{\mathrm{T}} \mu_2 + \frac{1}{2} \left(t_1^{\mathrm{T}} \Omega_{11} t_1 + 2 t_1^{\mathrm{T}} \Omega_{12} t_2 + t_2^{\mathrm{T}} \Omega_{22} t_2\right)\right\}$$

and by setting $t_2=0$ and $t_1=0$ we see that $M_{Y_1}(t_1)=\exp\left(t_1^{\mathrm{T}}\mu_1+\frac{1}{2}t_1^{\mathrm{T}}\Omega_{11}t_1\right)$ and $M_{Y_2}(t_2)=\exp\left(t_2^{\mathrm{T}}\mu_2+\frac{1}{2}t_2^{\mathrm{T}}\Omega_{22}t_2\right)$. Hence the marginal distribution of Y_1 is $\mathcal{N}_m(\mu_1,\Omega_{11})$. For the conditional distribution, note that $W=Y_1-\Omega_{12}\Omega_{22}^{-1}Y_2$ is a linear combination of Y and

$$E(W) = \mu_1 - \Omega_{12}\Omega_{22}^{-1}\mu_2, \quad \text{var}(W) = \Omega_{11} - \Omega_{12}\Omega_{22}^{-1}\Omega_{21}, \quad \text{cov}(W, Y_2) = \Omega_{12} - \Omega_{12}\Omega_{22}^{-1}\Omega_{22} = 0.$$

Hence $W \perp \!\!\! \perp Y_2$. As $Y_1 = W + \Omega_{12}\Omega_{22}^{-1}Y_2$ and conditioning on Y_2 does not change the distribution of W,

$$E(Y_1 \mid Y_2 = y_2) = E(W) + \Omega_{12}\Omega_{22}^{-1}y_2, \quad var(Y_1 \mid Y_2 = y_2) = var(W + \Omega_{12}\Omega_{22}^{-1}y_2) = var(W).$$

Putting the pieces together gives the stated conditional distribution.

(d) The joint MGF given in (c) factorises iff the variables are independent, and on inspecting it we see that

$$M_Y(t) = M_{Y_1}(t_1)M_{Y_2}(t_2) \iff \Omega_{12} = 0.$$

The variance matrix of

$$\begin{pmatrix} a \\ c \end{pmatrix} + \begin{pmatrix} B \\ D \end{pmatrix} Y$$

is

$$\begin{pmatrix} B\Omega B^{\mathrm{T}} & B\Omega D^{\mathrm{T}} \\ D\Omega B^{\mathrm{T}} & D\Omega D^{\mathrm{T}} \end{pmatrix}$$

so $a + BY \perp \!\!\! \perp c + DY$ iff $B\Omega D^{\mathrm{T}} = 0$.

Regression Methods

Autumn 2024 - note 1 of slide 19

Note: Multivariate normal distribution II

- (e) Each Y_j has mean μ and variance σ^2 , and since they are independent, $\operatorname{cov}(Y_j,Y_k)=0$ for $j\neq k$. If $u\in\mathbb{R}^n$, then $u^{\mathrm{T}}Y$ is a linear combination of normal variables, with mean $\sum_{j=1}^n u_j\mu=u^{\mathrm{T}}\mu 1_n$ and variance $\sum_{j=1}^n u_j^2\sigma^2=u^{\mathrm{T}}\sigma^2 I_n u$, so $Y\sim\mathcal{N}_n(\mu 1_n,\sigma^2 I_n)$, as required.
- (f) Since Ω is symmetric and positive semi-definite, the spectral theorem tells us that we may write $\Omega = ADA^{\mathrm{T}}$, where $D = \mathrm{diag}(d_1,\ldots,d_n)$ contains the (real) eigenvalues of Ω , with $d_1 \geq \cdots \geq d_n \geq 0$, and A is a $n \times n$ orthogonal matrix, i.e., $A^{\mathrm{T}}A = AA^{\mathrm{T}} = I_n$ and |A| = 1. The columns A_1,\ldots,A_n of A are the eigenvectors corresponding to the respective eigenvalues,

$$\Omega = ADA^{\mathrm{T}} = \sum_{j=1}^{n} d_j a_j a_j^{\mathrm{T}},$$

with $|\Omega|=|ADA^{\rm T}|=|A|\times |D|\times |A^{\rm T}|=|D|$ and $\Omega^{-1}=AD^{-1}A^{\rm T}$ if the inverse exists.

 \square Now let $Z_1,\ldots,Z_n\stackrel{\mathrm{iid}}{\sim}\mathcal{N}(0,1)$ $Z=(Z_1,\ldots,Z_n)^{\mathrm{\scriptscriptstyle T}}$, and $u\in\mathbb{R}^n$, set and consider

$$u^{\mathrm{T}}(\mu + AD^{1/2}Z) = u^{\mathrm{T}}\mu + \sum_{j=1}^{n} Z_{j}u^{\mathrm{T}}a_{j}d_{j}^{1/2}.$$

This is a linear combination of normal variables, so it has a normal distribution, with mean $u^{\mathrm{T}}\mu$ and variance

$$\operatorname{var}\left(u^{\mathrm{T}}\mu + \sum_{j=1}^{n} Z_{j}u^{\mathrm{T}}a_{j}d_{j}^{1/2}\right) = \sum_{j=1}^{n} d_{j}(u^{\mathrm{T}}a_{j})^{2}\operatorname{var}(Z_{j}) = u^{\mathrm{T}}\left(\sum_{j=1}^{n} d_{j}a_{j}a_{j}^{\mathrm{T}}\right)u = u^{\mathrm{T}}\Omega u,$$

so we can write $X = \mu + AD^{1/2}Z \sim N_n(\mu, \Omega)$.

 \square If Ω has rank n, then $d_n > 0$. The change of variables $z \mapsto x = \mu + AD^{1/2}z$ has Jacobian

$$\left| \frac{\partial x}{\partial z} \right| = |AD^{1/2}| = |A||D|^{1/2} = 1 \times |D|^{1/2} = |\Omega|^{1/2} > 0.$$

Moreover $z=D^{-1/2}A^{\mathrm{\scriptscriptstyle T}}(x-\mu)$, and therefore $z^{\mathrm{\scriptscriptstyle T}}z=(x-\mu)^{\mathrm{\scriptscriptstyle T}}\Omega^{-1}(x-\mu)$. Hence using the joint density of Z, $f_Z(z)=(2\pi)^{-n/2}\exp(-\sum_{j=1}^n z_j^2/2)$,

$$f_X(x) = f_Z(z)|_{z=D^{-1/2}A^{\mathrm{T}}(x-\mu)} \left| \frac{\partial z}{\partial x} \right| = (2\pi)^{-n/2} \exp\left(-\frac{z^{\mathrm{T}}z}{2}\right) \Big|_{z=D^{-1/2}A^{\mathrm{T}}(x-\mu)} |\Omega|^{-1/2},$$

which reduces to (1). If $d_n = 0$, then the Jacobian is zero, so the transformation $z \mapsto x$ is singular and X does not have a density on \mathbb{R}^n .

 \square Now suppose that $d_m > d_{m+1} = 0$, so just m eigenvalues of Ω are positive. Then

$$X = \mu + \sum_{j=1}^{m} Z_j a_j d_j^{1/2} \in \mathcal{S} = \mu + \text{span}(a_1, \dots, a_m),$$

where S is a hyperplane of dimension m passing through μ and generated by the vectors a_1, \ldots, a_m . In this case the previous argument shows that X has an m-dimensional Gaussian density on S, but places no probability elsewhere.

Regression Methods

Autumn 2024 - note 2 of slide 19

Reminder: χ^2 distribution

Definition 4 If $Y_j \stackrel{\mathrm{ind}}{\sim} \mathcal{N}(\mu_j, \sigma^2)$, then $W = Y_1^2 + \dots + Y_\nu^2$ has the non-central chi-square distribution with ν degrees of freedom (df) and non-centrality parameter $\delta^2 = (\mu_1^2 + \dots + \mu_\nu^2)/\sigma^2$; we write $W \sim \sigma^2 \chi_\nu^2(\delta^2)$. Then

$$M_W(t) = \exp\left(\frac{t\sigma^2\delta^2}{1 - 2t\sigma^2}\right)(1 - 2\sigma^2 t)^{-\nu/2}, \quad t < 1/(2\sigma^2).$$

If $\delta^2=0$ and $\sigma^2=1$ then W has the (central) chi-square distribution with ν df, we write $W\sim\chi^2_{\nu}$, its MGF is $M_W(t)=(1-2t)^{-\nu/2}$, and its p-quantile is $c_{\nu}(p)$.

Chi-square variables satisfy

- \square $E(W) = \sigma^2(\nu + \delta^2)$, $var(W) = 2\sigma^4(\nu + 2\delta^2)$;
- $\qquad \qquad \square \quad \text{if } W_1 \sim \chi^2_{\nu_1} \perp \!\!\! \perp W_2 \sim \chi^2_{\nu_2} \text{, then } W_1 + W_2 \sim \chi^2_{\nu_1 + \nu_2}; \\$
- $\square \quad W \sim \chi^2_{\nu}$ implies that W has the gamma density

$$f(w) = \frac{\beta^{\alpha} w^{\alpha - 1}}{\Gamma(\alpha)} e^{-\beta w}, \quad w > 0, \quad \alpha, \beta > 0,$$

with $\alpha = \nu/2$ and $\beta = 1/2$.

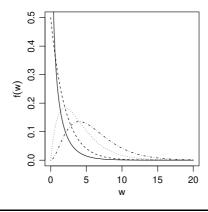
Regression Methods

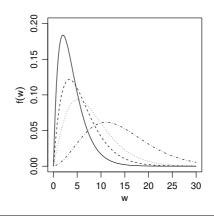
Autumn 2024 - slide 20

Reminder: χ^2_{ν} densities

Left: central densities with $\nu=1,2,4,6$ (solid, large dashes, small dashes, dot-dash).

Right: non-central densities with $\nu=4$ and $\delta=0,2,4,10$ (solid, large dashes, small dashes, dot-dash).





Regression Methods

Reminder: Student t distribution

Definition 5 If $Z \sim \mathcal{N}(0,1) \perp \!\!\! \perp W \sim \chi^2_{\nu}$, then $T = Z/(W/\nu)^{1/2}$ has the **Student** t distribution with ν df, $T \sim t_{\nu}$, and we write $t_{\nu}(p)$ for the corresponding p-quantile. The density function of T is

$$f_T(t) = \frac{\Gamma\{(\nu+1)/2\}}{\sqrt{\nu\pi}\Gamma(\nu/2)} \frac{1}{(1+t^2/\nu)^{(\nu+1)/2}}, \quad -\infty < t < \infty, \ \nu = 1, 2, \dots$$

Properties:

 \square the mean and variance exist only for $\nu \geq 2$ and $\nu \geq 3$ respectively, and then

$$E(T) = 0, \quad var(T) = \frac{\nu}{\nu - 2};$$

 \square with $\nu = 1$ we have the **Cauchy density**,

$$\frac{1}{\pi(1+t^2)}, \quad -\infty < t < \infty,$$

and then T has no moments;

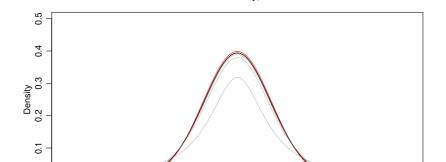
as $\nu \to \infty$, the limiting distribution of T is $\mathcal{N}(0,1)$; usually the approximation is 'good enough' for $\nu > 25$ (say).

Regression Methods

Autumn 2024 - slide 22

Reminder: Student t densities

Student t density functions with $\nu=1,5,10,20$ (black, $\nu=20$), and the standard normal density (red):



0

Student t density, nu=20

Regression Methods

Reminder: F distribution

Definition 6 If $W_1, W_2 \stackrel{\mathrm{ind}}{\sim} \chi^2_{\nu_1}, \chi^2_{\nu_2}$, then

$$F = \frac{W_1/\nu_1}{W_2/\nu_2}$$

has the F distribution with ν_1 and ν_2 df: we write $F \sim F_{\nu_1,\nu_2}$. The density function is

$$f_F(u) = \frac{\Gamma\left(\frac{1}{2}\nu_1 + \frac{1}{2}\nu_2\right)\nu_1^{\nu_1/2}\nu_2^{\nu_2/2}}{\Gamma\left(\frac{1}{2}\nu_1\right)\Gamma\left(\frac{1}{2}\nu_2\right)} \frac{u^{\frac{1}{2}\nu_1 - 1}}{(\nu_2 + \nu_1 u)^{(\nu_1 + \nu_2)/2}}, \quad u > 0, \ \nu_1, \nu_2 = 1, 2, \dots,$$

and the p-quantile is written $F_{\nu_1,\nu_2}(p)$.

Regression Methods

Autumn 2024 - slide 24

Reminder: Computation

- Quantiles of the $\mathcal{N}(\mu, \sigma^2)$, χ^2_{ν} , t_{ν} , F_{ν_1,ν_2} distributions can be found in tables, or in environments such as R (see http://www.r-project.org/), where they can also be simulated.
- ☐ Examples:

R: Copyright 2005, The R Foundation for Statistical Computing Version 2.2.1 (2005-12-20 r36812)

. . .

> qnorm(0.025) # this is a comment; access normal quantiles

[1] -1.959964 # the [1] means this is the first element of a vector

> ?qnorm # help on use of function qnorm()

> qchisq(0.025,df=3) # chi-squared quantiles, nu=3

[1] 0.2157953

> qt(0.025,df=3) # t quantiles, nu=3

[1] -3.182446

> qf(0.025,df1=3,df2=4) # F quantiles, nu1=3, nu2=4

[1] 0.06622087

Regression Methods

Autumn 2024 - slide 25

Statistical models

- \square Least squares fitting gives a deterministic description of the variation in some numbers y in terms of other numbers X.
- \square A **statistical model** is a description of data y in terms of a collection of probability distributions on the sample space for y.
- ☐ We distinguish
 - primary aspects of a model, which specify what questions we aim to answer, from
 - secondary aspects, which complete the model, indicate what analysis might be suitable, and determine the precision of conclusions.
- Often the primary aspects are embodied in one or more parameters of the model.
- ☐ (Almost) all models are **tentative**, and we must check that they are reasonable.

Second-order and normal assumptions

- ☐ Two distributional assumptions are in general use for the linear model:
 - second-order assumptions,

$$y \sim (X\beta, \sigma^2 V)$$
, i.e., $E(y) = X\beta$, $var(y) = \sigma^2 V_{n \times n}$;

normal assumptions,

$$y \sim \mathcal{N}_n(X\beta, \sigma^2 V),$$

i.e., y has a multivariate normal distribution with mean vector $X\beta$ and positive definite (co)variance matrix σ^2V .

- \square X is called the **design matrix**: more later.
- \square V is assumed known. Unless stated otherwise we set $V=I_n$, so the y_j are uncorrelated; if normal they are therefore independent.
- \square If $V \neq I_n$, then we can perform weighted least squares (WLS) estimation, minimising

$$||y - X\beta||_V^2 = (y - X\beta)^T W(y - X\beta),$$

where $W = V^{-1}$ is the **weight matrix**.

Above the **linearity** is (usually) primary, whereas the **distributional assumption**, use of weights, . . . , are (usually) secondary.

Regression Methods

Autumn 2024 - slide 27

Consequences of second-order assumptions

Lemma 7 Under the second-order assumptions, $\widehat{\beta}$ is an unbiased estimator of β ,

$$E(\widehat{\beta}) = \beta$$
, $var(\widehat{\beta}) = \sigma^2 (X^T X)^{-1}$.

and $S^2 = (n-p)^{-1} \|y - \widehat{y}\|^2$ is an unbiased estimator of σ^2 .

Theorem 8 (Gauss–Markov) The least squares estimator $\widehat{\beta}$ has the smallest variance among all estimators $\widetilde{\beta} = A_{p \times n} y$; it is the best linear unbiased estimator (BLUE) of β .

- ☐ Obviously these results hold under the (stronger) normal assumptions.
- \square The Gauss–Markov theorem only concerns <u>linear</u> estimators. Nonlinear estimators of β might have smaller variance than $\sigma^2(X^{\mathrm{T}}X)^{-1}$ (and in fact the optimal maximum likelihood estimators of β for non-normal models will be nonlinear in y).

Regression Methods

Note to Lemma 7

- \square Recall that expectation is linear, and that $var(A_{p\times n}y) = Avar(y)A^{\mathrm{T}}$.
- \square Set $A_{p \times n} = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}$ and note that

$$\begin{split} & \mathrm{E}(\widehat{\beta}) &= \mathrm{E}(Ay) = A\mathrm{E}(y) = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}X\beta = \beta, \\ & \mathrm{var}(\widehat{\beta}) &= A\mathrm{var}(y)A^{\mathrm{T}} = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}I_{n}\sigma^{2}\{(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}\}^{\mathrm{T}} = \sigma^{2}(X^{\mathrm{T}}X)^{-1}. \end{split}$$

 \square Recall that $\mathrm{E}(yy^{\mathrm{T}}) = \mathrm{var}(y) + \mathrm{E}(y)\mathrm{E}(y)^{\mathrm{T}} = \sigma^2 I_n + X\beta\beta^{\mathrm{T}}X^{\mathrm{T}}$, and note that

$$||y - \hat{y}||^2 = (y - \hat{y})^{\mathrm{T}}(y - \hat{y}) = y^{\mathrm{T}}(I_n - H)^{\mathrm{T}}(I_n - H)y = y^{\mathrm{T}}(I_n - H)y = \operatorname{tr}\{(I_n - H)yy^{\mathrm{T}}\}.$$

Hence $\mathrm{E}(\|y-\widehat{y}\|^2)$ equals

$$E[tr\{(I_n - H)yy^{T}\}] = tr\{(I_n - H)E(yy^{T})\} = tr\{(I_n - H)(\sigma^2 I_n + X\beta\beta^{T}X^{T})\} = \sigma^2 tr(I_n - H),$$

because $(I_n - H)X = 0$. Moreover $tr(I_n) = n$ and

$$tr(H) = tr\{X(X^{T}X)^{-1}X^{T}\} = tr\{(X^{T}X)^{-1}X^{T}X\} = tr(I_{p}) = p,$$

so $E(S^2) = \sigma^2$, because

$$E(||y - \widehat{y}||^2) = \sigma^2 \operatorname{tr}(I_n - H) = \sigma^2(n - p).$$

Regression Methods

Autumn 2024 - note 1 of slide 28

Note to Theorem 8

 \square Let $\tilde{\beta}$ denote any unbiased estimator of β that is linear in y. Then a $p \times n$ matrix A exists such that $\tilde{\beta} = Ay$, and unbiasedness implies that $\mathrm{E}(\tilde{\beta}) = AX\beta = \beta$ for any parameter vector β ; this entails $AX = I_p$. Now

$$\operatorname{var}(\widetilde{\beta}) - \operatorname{var}(\widehat{\beta}) = A\sigma^{2}I_{n}A^{\mathsf{T}} - \sigma^{2}(X^{\mathsf{T}}X)^{-1}$$

$$= \sigma^{2} \left\{ AA^{\mathsf{T}} - AX(X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}A^{\mathsf{T}} \right\}$$

$$= \sigma^{2}A(I_{n} - H)A^{\mathsf{T}}$$

$$= \sigma^{2}A(I_{n} - H)(I_{n} - H)^{\mathsf{T}}A^{\mathsf{T}}$$

and this $p \times p$ matrix is positive semidefinite. Thus $\widehat{\beta}$ has smallest variance in finite samples among all linear unbiased estimators of β .

 \square This is a finite-sample result that holds for all n and X (of rank p, with $n \ge p$).

Regression Methods

Autumn 2024 - note 2 of slide 28

Second-order assumptions and large samples

 $\Box \quad \text{We can write } y_j = x_j^{ \mathrm{\scriptscriptstyle T} } \beta + \sigma \varepsilon_j \text{, where } \varepsilon_j \overset{\text{ind}}{\sim} (0,1) \text{, so}$

$$\widehat{\beta} = (X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}y = \sum_{j=1}^{n} (X^{\mathsf{T}}X)^{-1}x_{j}y_{j} = \beta + \sigma n^{-1}\sum_{j=1}^{n} a_{j}\varepsilon_{j},$$

say, where a_1, \ldots, a_n are $p \times 1$ vectors. We have $E(\widehat{\beta}) = \beta$ and $var(\widehat{\beta}) = (X^T X)^{-1}$, but is $\widehat{\beta}$ approximately normal for large n?

 \square The a_j , or equivalently X, must be such that no single y_j can dominate in $n^{-1} \sum a_j \varepsilon_j$.

Theorem 9 (no proof) Let $\{X_n\}$ be a sequence of $n \times p$ design matrices each of rank p, let $h_{11}^n, \ldots, h_{nn}^n$ be the diagonal elements of the hat matrices $X_n(X_n^{\mathrm{T}}X_n)^{-1}X_n^{\mathrm{T}}$ and let $y_n \sim (X_n\beta, \sigma^2 I_n)$ for each n. If

$$\lim_{n \to \infty} \max_{j=1,\dots,n} h_{jj} = 0,$$

then the corresponding sequence of least squares estimators \widehat{eta}_n satisfies

$$(X_n^{\mathrm{T}} X_n)^{1/2} (\widehat{\beta}_n - \beta) \xrightarrow{D} \mathcal{N}_n(0, \sigma^2 I_n), \quad n \to \infty,$$

i.e., if H has a 'well-behaved' diagonal, then $\widehat{\beta} \sim \mathcal{N}_p\{\beta, \sigma^2(X^{\mathrm{T}}X)^{-1}\}$ in large samples.

Regression Methods

Autumn 2024 - slide 29

Normal-theory linear model

The following results allow exact inferences for β and σ^2 , and in analysis of variance.

Theorem 10 Under the normal-theory linear model,

$$\widehat{\beta} \sim \mathcal{N}_p\{\beta, \sigma^2(X^{\mathrm{T}}X)^{-1}\}$$
 $\perp \perp \frac{(n-p)S^2}{\sigma^2} \sim \chi_{n-p}^2$.

Lemma 11 If $y \sim \mathcal{N}_n(\mu, \sigma^2 I_n)$ and H is symmetric and idempotent with rank p, then $y^{\mathrm{T}}Hy \sim \sigma^2 \chi_p^2(\delta^2)$, where $\sigma^2 \delta^2 = \mu^{\mathrm{T}}H\mu$.

Theorem 12 If $y \sim \mathcal{N}_n(\mu, \sigma^2 I_n)$ and a linear model is fitted whose design matrix X is structured as in Lemma 2, then the sums of squares in the ANOVA decomposition

$$\sum_{j=1}^{n} (y_j - \overline{y})^2 = \sum_{r=1}^{R} \|\widehat{y}_r - \widehat{y}_{r-1}\|^2 + \|y - \widehat{y}\|^2 = \sum_{r=1}^{R+1} \|P_r y\|^2$$

are independent and $\|P_ry\|^2 \sim \sigma^2 \chi^2_{\nu_{r-1}-\nu_r}(\delta_r^2/\sigma^2)$, where $\sigma^2 \delta_r^2 = \mu^{\mathrm{T}} P_r \mu$. If X_r does not explain any variation in μ after allowing for X_0, \ldots, X_{r-1} , then $P_r\mu = 0$, so $\delta_r^2 = 0$.

Theorem 12 implies that the sums of squares for terms that explain variation in y will tend to be larger than sums of squares for other terms, which can be used to estimate σ^2 .

Regression Methods

Note to Theorem 10

- \square The first part is easy, because $\widehat{\beta}$ is a linear combination of normal variables so it is normal, and its mean and variance matrix were given by Lemma 7.
- Likewise the residual $e = y \widehat{y} = (I H)y$ is a linear combination of y with mean 0_n and variance $(I H)\sigma^2$, so $e \sim \mathcal{N}_n\{0_p, (I H)\sigma^2\}$.
- \square As $cov(\widehat{\beta}, e)$ equals

$$cov\{(X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}y, (I-H)y\} = (X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}cov(y)(I-H)^{\mathsf{T}} = \sigma^{2}(X^{\mathsf{T}}X)^{-1}\{(I-H)X\}^{\mathsf{T}} = 0,$$

we see that $\widehat{\beta}$ is independent of (any function of) e, and therefore in particular of

$$(n-p)S^2/\sigma^2 = ||y-\widehat{y}||^2/\sigma^2 = e^{\mathrm{T}}e/\sigma^2.$$

The eigenvalues of H are p 1's and n-p 0's, so those of I-H are n-p 1's and p 0's. The spectral decomposition implies that there exists an $n\times n$ orthogonal matrix U such that $I-H=UDU^{\rm T}$, where $D={\rm diag}(1,\dots,1,0,\dots,0)$ and $UU^{\rm T}=U^{\rm T}U=I_n$. Thus $Z=U^{\rm T}e/\sigma$ has mean vector 0_n and variance matrix

$$\operatorname{var}(Z) = U^{\mathrm{T}} \operatorname{var}(e) U / \sigma^{2} = U^{\mathrm{T}} (I - H) \sigma^{2} U / \sigma^{2} = U^{\mathrm{T}} U D U^{\mathrm{T}} U = D,$$

i.e. the Z_1, \ldots, Z_n are independent normal variables, n-p of them have variance 1 and p of them have variance 0 and therefore equal 0 with probability one. Hence, as required,

$$(n-p)S^2/\sigma^2 = e^{\mathrm{T}}e/\sigma^2 = (UZ)^{\mathrm{T}}(UZ) = Z^{\mathrm{T}}U^{\mathrm{T}}UZ = \sum_{j=1}^{n-p} Z_j^2 \sim \chi_{n-p}^2.$$

Regression Methods

Autumn 2024 - note 1 of slide 30

Note to Lemma 11

The spectral decomposition of H is UDU^{T} , where D is diagonal with p 1's and n-p 0's, and $Z=U^{\mathrm{T}}y\sim\mathcal{N}_n(U^{\mathrm{T}}\mu,\sigma^2I_n)$; note that the Z_j are independent. Now

$$y^{\mathrm{T}}Hy = (U^{\mathrm{T}}y)^{\mathrm{T}}D(U^{\mathrm{T}}y) = \sum_{j=1}^{n} d_{j}Z_{j}^{2} = \sum_{j:d_{j}=1}Z_{j}^{2},$$

which has a (possibly non-central) χ^2 distribution with $p={\rm tr}(H)$ degrees of freedom, scale parameter σ^2 and

$$\sigma^2 \delta^2 = \sum_{j:d_j=1} E(Z_j)^2 = \sum_{j=1}^n d_j E(Z_j)^2 = (U^{\mathrm{T}} \mu)^{\mathrm{T}} D(U^{\mathrm{T}} \mu) = \mu^{\mathrm{T}} H \mu.$$

Regression Methods

Autumn 2024 - note 2 of slide 30

Note to Theorem 12

- As $P_r P_s = 0$ for $r \neq s$, we have $cov(P_r y, P_s y) = P_r var(y) P_s^{\mathrm{T}} = \sigma^2 P_r P_s = 0$, i.e., $P_r y$ and $P_s y$ are independent. Hence the terms in the ANOVA decomposition are independent.
- \square P_r is a symmetric idempotent matrix, so Lemma 11 gives

$$||P_r y||^2 \sim \sigma^2 \chi_{\nu}^2 (\delta_r^2 / \sigma^2), \quad \delta_r^2 = \mu^{\mathrm{T}} P_r \mu,$$

where $\nu=\mathrm{rank}(P_r)$. These ranks are $\nu_{r-1}-\nu_r$ for $r=1,\ldots,R$, and $\nu_{R+1}=n-p$ for $P_{R+1}=I_n-H$.

If X_r does not explain any variation in μ after allowing for X_0, \ldots, X_{r-1} , then $H_r\mu = H_{r-1}\mu \in \mathcal{V}_{r-1}$, i.e., $P_r\mu = 0$, and thus $\delta_r^2 = 0$.

Regression Methods

Autumn 2024 - note 3 of slide 30

Inference on β

 \square Theorem 10 implies that for any constant $c_{p\times 1}$, $c^{\mathrm{T}}\widehat{\beta}\sim \mathcal{N}\{c^{\mathrm{T}}\beta,\sigma^{2}c^{\mathrm{T}}(X^{\mathrm{T}}X)^{-1}c\}$, so

$$Z = \frac{c^{\mathrm{T}} \widehat{\beta} - c^{\mathrm{T}} \beta}{\sigma \sqrt{c^{\mathrm{T}} (X^{\mathrm{T}} X)^{-1} c}} \sim \mathcal{N}(0, 1) \quad \perp \perp \quad (n - p) S^2 / \sigma^2 = W \sim \chi_{n - p}^2,$$

and thus

$$\frac{c^{\mathrm{T}}\widehat{\beta}_r - c^{\mathrm{T}}\beta_r}{S\sqrt{c^{\mathrm{T}}(X^{\mathrm{T}}X)^{-1}c}} = \frac{Z}{\sqrt{W/(n-p)}} \sim t_{n-p}.$$

- \square Let v_{rs} denote the (r,s) element of $(X^{\mathrm{T}}X)^{-1}$, so v_{rr} denotes its rth diagonal element.
- \square Different choices of c allow inferences on the elements of β .
- \square For example, if $c^{\mathrm{T}}=(c_1,\ldots,c_p)$, $c_r=1$ and $c_s=0$ for $s\neq r$, then $c^{\mathrm{T}}\beta=\beta_r$, and we
 - test the hypothesis that $\beta_r=\beta_r^0$ by comparing $(\widehat{\beta}_r-\beta_r^0)/(Sv_{rr}^{1/2})$ to the t_{n-p} distribution, and
 - a $(1-\alpha)$ confidence interval for β_r has limits

$$\hat{\beta}_r \pm S v_{rr}^{1/2} t_{n-p} (1 - \alpha/2), \quad 0 < \alpha < 1.$$

 \square Likewise we can compare β_r and β_s by setting $c_r=1$, $c_s=-1$ and all other $c_t=0$.

Regression Methods

Autumn 2024 - slide 31

Prediction

Inference for the value of a further random variable Y_+ with known $p \times 1$ covariate vector x_+ and satisfying the linear model, so $Y_+ \sim \mathcal{N}(x_+^{\mathrm{\scriptscriptstyle T}}\beta, \sigma^2)$ independent of the other variables, is performed by noting that $Y_+ \perp\!\!\!\perp \widehat{\beta}, S^2$ and

$$Y_{+} - x_{+}^{\mathrm{T}} \widehat{\beta} \sim \mathcal{N} \left[0, \sigma^{2} \{ 1 + x_{+}^{\mathrm{T}} (X^{\mathrm{T}} X)^{-1} x_{+} \} \right],$$

SO

$$\frac{Y_{+} - x_{+}^{\mathrm{T}} \widehat{\beta}}{S\{1 + x_{+}^{\mathrm{T}} (X^{\mathrm{T}} X)^{-1} x_{+}\}^{1/2}} \sim t_{n-p},$$

which leads to prediction intervals for Y_+ once $\widehat{\beta}$ and S have been observed.

Although we expect inferences for β and σ^2 to hold as approximations under second-order assumptions, this is not the case for inference on Y_+ . (Why not?)

1.3 Analysis of Variance

slide 33

Analysis of variance

 \square We previously saw that

$$\sum_{j=1}^{n} (y_j - \overline{y})^2 = ||y||^2 - ||\widehat{y}_0||^2 = \sum_{r=1}^{R} ||\widehat{y}_r - \widehat{y}_{r-1}||^2 + ||y - \widehat{y}||^2$$

decomposes ('analyses') the variability of y around its average \overline{y} into

- the contributions $\|\widehat{y}_r \widehat{y}_{r-1}\|^2$ due to adding the columns of X_r to X_0, \ldots, X_{r-1} ,
- the residual sum of squares $||y \hat{y}||^2$ left after fitting $X = (X_0, \dots, X_R)$.
- \square Large $\|\widehat{y}_r \widehat{y}_{r-1}\|^2$ implies that X_r explains a lot of the variation of y even after allowing for that explained by X_0, \ldots, X_{r-1} .
- Theorem 12 implies that under the normal assumptions, and if $E(y) = \mu$ lies in the column space of X, the sums of squares on the RHS above are independent and satisfy

$$\|\widehat{y}_r - \widehat{y}_{r-1}\|^2 = \|P_r y\|^2 \sim \sigma^2 \chi^2_{\nu_{r-1} - \nu_r} (\delta_r^2 / \sigma^2)$$
 $\perp \!\!\! \perp \|y - \widehat{y}\|^2 \sim \sigma^2 \chi^2_{n-p}.$

Hence if $\delta_r^2 = 0$, i.e., $\mu \in \operatorname{span}(X_0, \dots, X_{r-1})$, then

$$\frac{\|\widehat{y}_r - \widehat{y}_{r-1}\|^2/(\nu_{r-1} - \nu_r)}{\|y - \widehat{y}\|^2/(n-p)} \sim F_{\nu_{r-1} - \nu_r, n-p}.$$

Regression Methods

Autumn 2024 - slide 34

ANOVA table

Term added	df	Reduction in SS	Mean square
X_1	$n - 1 - \nu_1$	$SS_0 - SS_1$	$MS_1 = (SS_0 - SS_1)/(n - 1 - \nu_1)$
X_2	$\nu_1 - \nu_2$	$SS_1 - SS_2$	$MS_2 = (SS_1 - SS_2)/(\nu_1 - \nu_2)$
:	:	:	:
X_R	$\nu_{R-1} - \nu_R$	$SS_{R-1} - SS_R$	$MS_R = (SS_{R-1} - SS_R)/(\nu_{R-1} - \nu_R)$
Residual	$\nu_R = n - p$	SS_R	$MS_{Res} = SS_R/\nu_R$

- \square If $\mu \in \operatorname{span}(X)$ then the residual mean square $\operatorname{MS}_{\operatorname{Res}}$ gives an estimate of σ^2 .
- \square We test for an effect of term X_r by noting that
 - if X_r explains no more than (X_0,\ldots,X_{r-1}) , then

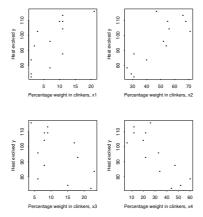
$$F_r = \frac{\mathrm{MS}_r}{\mathrm{MS}_{\mathrm{Res}}} \sim F_{\nu_{r-1} - \nu_r, \nu_R},$$

- if X_r does have additional explanatory power, then the distribution of MS_r is shifted to the right, and we expect F_r to be large relative to its null distribution.

Regression Methods

Example: Cement data

Percentage weights in clinkers of 4 four constitutents of cement (x_1, \ldots, x_4) and heat evolved y in calories, in n = 13 samples.



Regression Methods

Autumn 2024 - slide 36

Example: Cement data

> cement

x1 x2 x3 x4 y L 7 26 6 60 78.5

2 1 29 15 52 74.3

3 11 56 8 20 104.3

4 11 31 8 47 87.6

5 7 52 6 33 95.9

6 11 55 9 22 109.2

7 3 71 17 6 102.7

8 1 31 22 44 72.5

9 2 54 18 22 93.1

10 21 47 4 26 115.9

11 1 40 23 34 83.8

12 11 66 9 12 113.3 13 10 68 8 12 109.4

Regression Methods

Example: Cement data

- $\hfill\square$ Reductions in overall sum of squares when terms entered in the order given.
- $\ \square$ Clearly x_1 and x_2 should be included, maybe not the others.

Term	df	Reduction in	Mean square	F
		sum of squares		
x_1	1	1450.1	1450.1	242.5
x_2	1	1207.8	1207.8	202.0
x_3	1	9.79	9.79	1.64
x_4	1	0.25	0.25	0.04
Residual	8	47.86	5.98	

Regression Methods

Autumn 2024 - slide 38

Example: Cement data

 \square What if we change the order of the terms?

Term	df Reduction in		Mean square	\overline{F}
		sum of squares		
x_4	1	1831.9	1831.9	306.2
x_3	1	708.1	708.1	118.4
x_2	1	101.9	101.9	17.04
x_1	1	26.0	26.0	4.34
Residual	8	47.86	5.98	

 $\ \square$ Should x_1 and x_2 be included or not?

Regression Methods

Orthogonality

- ☐ In general, the ANOVA and ANOVA table depend on the order of inclusion of terms.
- Its interpretation is unclear if X_r is significant when included early, but not when it is included late. Is the term important or not?
- \square In a model with orthogonal terms,

$$X\beta = 1_n\beta_0 + X_1\beta_1 + X_2\beta_2, \quad X_r^{\mathrm{T}}X_s = X_r^{\mathrm{T}}1_n = 0, \quad r \neq s.$$

we obtain

$$\begin{pmatrix} \widehat{\beta}_0 \\ \widehat{\beta}_1 \\ \widehat{\beta}_2 \end{pmatrix} = \begin{pmatrix} 1^{\mathrm{T}} 1 & 0 & 0 \\ 0 & X_1^{\mathrm{T}} X_1 & 0 \\ 0 & 0 & X_2^{\mathrm{T}} X_2 \end{pmatrix}^{-1} \begin{pmatrix} 1 & X_1 & X_2 \end{pmatrix}^{\mathrm{T}} y$$

so since $\widehat{y}=X\widehat{\beta}$, we have

$$y^{\mathrm{T}}y - \widehat{y}^{\mathrm{T}}\widehat{y} = y^{\mathrm{T}}y - n\overline{y}^{2} - \widehat{\beta}_{1}^{\mathrm{T}}X_{1}^{\mathrm{T}}X_{1}\widehat{\beta}_{1} - \widehat{\beta}_{2}^{\mathrm{T}}X_{2}^{\mathrm{T}}X_{2}\widehat{\beta}_{2},$$

and the residual sums of squares for the sub-models $1_n\beta_0$, $1_n\beta_0+X_1\beta_1$, $1_n\beta_0+X_2\beta_2$ are

$$y^{\mathrm{T}}y - n\overline{y}^{2}, \quad y^{\mathrm{T}}y - n\overline{y}^{2} - \widehat{\beta}_{1}^{\mathrm{T}}X_{1}^{\mathrm{T}}X_{1}\widehat{\beta}_{1}, \quad y^{\mathrm{T}}y - n\overline{y}^{2} - \widehat{\beta}_{2}^{\mathrm{T}}X_{2}^{\mathrm{T}}X_{2}\widehat{\beta}_{2},$$

so the reductions do not depend on the order of inclusion. Hooray!

Regression Methods

Autumn 2024 - slide 40

Balance

- \square Balanced design matrices induce orthogonality after fitting 1_n (or a more complex design X_0).
- ☐ Gram—Schmidt orthogonalisation with respect to early terms makes later terms mutually orthogonal, leading to a clear interpretation of the ANOVA for the later terms.
- \square If we write $H_0 = X_0 (X_0^{\mathrm{T}} X_0)^{-1} X_0^{\mathrm{T}}$ and let

$$Z_r = P_0 X_r = (I_n - H_0) X_r, \quad r = 1, 2,$$

denote the versions of X_1 and X_2 after adjusting for X_0 , then

$$X_0\beta_0 + X_1\beta_1 + X_2\beta_2 = (X_0\beta_0 + H_0X_1\beta_1 + H_0X_2\beta_2) + P_0X_1\beta_1 + P_0X_2\beta_2$$

= $Z_0\gamma_0 + Z_1\beta_1 + Z_2\beta_2$,

say, and $Z_1^{\mathrm{T}} Z_0 = Z_2^{\mathrm{T}} Z_0 = 0$, because $P_0 X_0 = P_0 H_0 = 0$.

 \square If the design satisfies $Z_1^{\mathrm{T}}Z_2=0$, then the order of inclusion of X_1 , X_2 is irrelevant, provided X_0 is already present in the fit.

Example 13 (3×2 layout) Observations and their means written as

$$y_{11}$$
 y_{12} $\mu + \alpha_1 + \delta_1$ $\mu + \alpha_2 + \delta_1$

$$y_{21}$$
 y_{22} , $\mu + \alpha_1 + \delta_2$ $\mu + \alpha_2 + \delta_2$.

$$y_{31}$$
 y_{32} $\mu + \alpha_1 + \delta_3$ $\mu + \alpha_2 + \delta_3$

Regression Methods

Note to Example 13

 \square In terms of the parameter vector $(\mu, \alpha_1, \alpha_2, \delta_1, \delta_2, \delta_3)^T$, the design matrix is

$$X_{6\times 6}^* = \begin{pmatrix} 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 0 & 1 \end{pmatrix}, \quad \text{with responses} \quad y = \begin{pmatrix} y_{11} \\ y_{12} \\ y_{21} \\ y_{22} \\ y_{31} \\ y_{32} \end{pmatrix},$$

with $X_0 \equiv 1_6$ the first column of X^* , columns 2–3 the term X_1^* for columns, and columns 4–6 the term X_2^* for rows.

This model has six parameters, but they cannot all be estimated, because X_0 lies in the column spaces of X_1^* and X_2^* , and it is easy to check that X^* has rank 4. The usual way to deal with this is to set $\alpha_1 = \delta_1 = 0$, corresponding to dropping columns 2 and 4 of X^* , giving the so-called corner-point parametrization in which the means are

$$y_{11} \quad y_{12} \qquad \qquad \mu \qquad \mu + \alpha_2 y_{21} \quad y_{22}, \qquad \qquad \mu + \delta_2 \quad \mu + \alpha_2 + \delta_2, y_{31} \quad y_{32} \qquad \qquad \mu + \delta_3 \quad \mu + \alpha_2 + \delta_3$$

i.e.,

- the 'grand mean' μ corresponds to the mean of observations with the first level of every factor,
- α_2 corresponds to the mean difference between column 2 and column 1,
- δ_2 corresponds to the mean difference between row 2 and row 1, and
- δ_3 corresponds to the mean difference between row 3 and row 1.

This is the default in R. More rarely we might set $\sum_{c} \alpha_{c} = \sum_{r} \delta_{r} = 0$.

☐ Even after these columns are dropped to give

$$X = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 1 \end{pmatrix},$$

the terms X_1 for columns and X_2 for rows are not orthogonal, and they are not orthogonal to 1_n . On the other hand Z_1 and Z_2 in the corresponding centred matrix,

$$\begin{pmatrix} 1 & -\frac{1}{2} & -\frac{1}{3} & -\frac{1}{3} \\ 1 & \frac{1}{2} & -\frac{1}{3} & -\frac{1}{3} \\ 1 & -\frac{1}{2} & \frac{2}{3} & -\frac{1}{3} \\ 1 & \frac{1}{2} & \frac{2}{3} & -\frac{1}{3} \\ 1 & -\frac{1}{2} & -\frac{1}{3} & \frac{22}{3} \\ 1 & \frac{1}{2} & -\frac{1}{3} & \frac{22}{3} \end{pmatrix}$$

are orthogonal to the constant by construction and to each other because the design is balanced: δ_2 and δ_3 each occur equally often with α_2 and without α_2 . This balance implies that if μ is fitted first, the reductions in sums of squares due to X_1 and X_2 , or equivalently Z_1 and Z_2 , are unique.

1.4 Diagnostics slide 42

Assumptions and model checking

- ☐ How heavily do our conclusions depend on our assumptions?
- ☐ In any given context,
 - **primary** aspects relate to the questions our analysis will address,
 - secondary aspects relate to how we go about finding answers to them.
- ☐ Concerns about primary aspects suggest that we should start again.
- ☐ Concerns about secondary aspects suggest that we modify the analysis.
- Regression diagnostics check that a fitted model is adequate:
 - Does y depend linearly on the columns of X?
 - Does y depend systematically on variables omitted from X?
 - Are the variances constant?
 - Are the responses uncorrelated/independent?
 - Are there outliers or otherwise unusual data?
 - Are the responses normally distributed?
- ☐ Usually these involve plots, sometimes tests beware over-interpretation!
- ☐ Key question: 'how would the failure I see/suspect change my conclusions?'

Regression Methods

Autumn 2024 - slide 43

Residuals

☐ The raw residuals

$$e = y - \widehat{y} = y - X\widehat{\beta} = (I_n - H)y$$

have $\mathrm{E}(e)=0$, $\mathrm{var}(e)=\sigma^2(I_n-H)$ if model correct, so

$$var(e_j) = \sigma^2(1 - h_{jj}) \quad cov(e_j, e_k) = -\sigma^2 h_{jk}, \ j \neq k.$$

☐ To (roughly) equalise the variances we define **standardized residuals**

$$r_j = \frac{e_j}{s(1 - h_{ij})^{1/2}} = \frac{y_j - x_j^{\mathrm{T}} \widehat{\beta}}{s(1 - h_{ij})^{1/2}}, \quad j = 1, \dots, n,$$

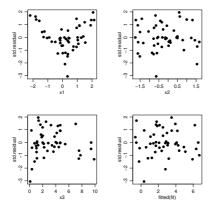
with s replacing σ . Then $E(r_j) = 0$ and $var(r_j) \doteq 1$.

- \square Although $e^{\mathrm{T}} \widehat{y} = \mathrm{cov}(e, \widehat{y}) = 0$ (check!), this only implies no <u>linear</u> relation between e and \widehat{y} .
- ☐ We check
 - linearity by plotting r_i against the covariates (those in X and those not in X);
 - constant variance by plotting r_j (or $|r_j|$) against fitted values \widehat{y}_j ;
 - independence by ACF of residuals (if data time-ordered);
 - for outliers, which are visible as unusual residuals; and
 - normality using a normal QQ-plot of r_i .

Regression Methods

Checking linearity

Plot r against each covariate, included or not in the model, and against \hat{y} , which is uncorrelated with e (as $\hat{y}^{\mathrm{T}}e=0$):



Regression Methods

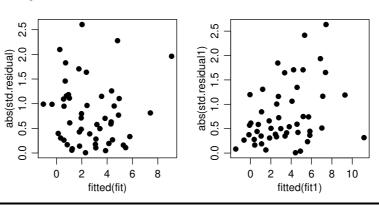
Autumn 2024 - slide 45

Checking the variance

 \square Does var(y) depend on E(y)?

Variance function shows how var(y) depends on $\mu = E(y)$. For normal linear model should have $var(y) = \sigma^2$, so variance is constant function of μ

 \square Plot r or |r| against \widehat{y} :



Regression Methods

Checking independence

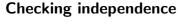
- ☐ Dependence can greatly increase uncertainty of final conclusions.
- ☐ Substantive knowledge is helpful in suggesting whether it might be present:
 - were the data gathered in temporal/spatial/...order?
 - were the data sampled/gathered in groups (e.g., spatial, several observations on different individuals, . . .)?
 - was randomisation used? If so, how?
- ☐ If observations are time-ordered, try using correlogram (ACF) and partial correlogram (PACF) to estimate serial correlations and partial correlations

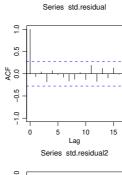
$$corr(r_j, r_{j+t}), \quad corr(r_j, r_{j+t} \mid r_{j+1}, \dots, r_{j+t-1}), \quad t = 1, \dots$$

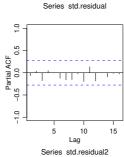
On next page, top panels show uncorrelated residuals, lower ones show evidence of correlation, suggesting use of a time series model.

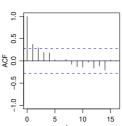
Regression Methods

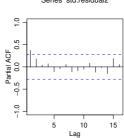
Autumn 2024 - slide 47











Regression Methods

Checking for outliers and normality

 \square Normal Q-Q plot for $Y_1,\ldots,Y_n\stackrel{\mathrm{iid}}{\sim}\mathcal{N}(\mu,\sigma^2)$ graphs ordered values

$$Y_{(1)} \le Y_{(2)} \le \dots \le Y_{(n)}$$

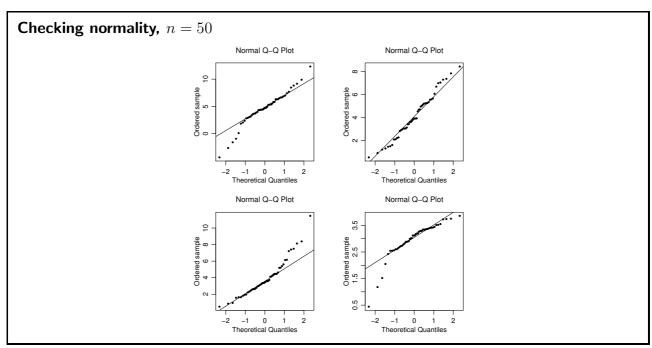
against (approximate) expected normal order statistics

$$\Phi^{-1}\{1/(n+1)\}, \Phi^{-1}\{2/(n+1)\}, \dots, \Phi^{-1}\{n/(n+1)\}.$$

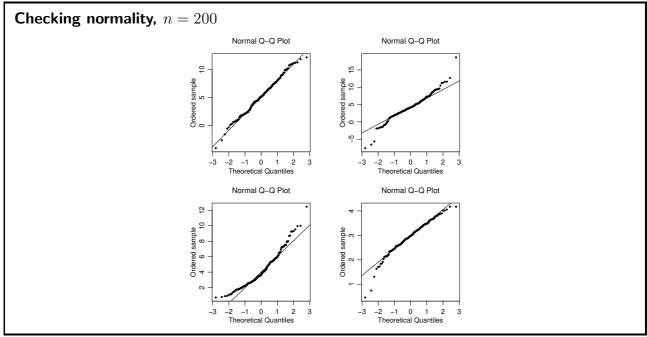
- \square Normality roughly straight line, slope σ , intercept μ .
- ☐ Outliers, skewness, heavy tails (easily) seen.
- \square Beware over-interpretation of such plots when n is small often useful to add simulation envelope.
- $\ \square$ Apply to standardized residuals r_j from regression model.

Regression Methods

Autumn 2024 - slide 49



Regression Methods



Regression Methods

Autumn 2024 - slide 51

Leverage and influence

- \square Does case (x_j, y_j) strongly influence the fitted model (picture)?
- ☐ As

$$\operatorname{var}(y_j - \widehat{y}_j) = \operatorname{var}(y_j - x_j^{\mathrm{T}}\beta) = \sigma^2(1 - h_{jj}),$$

having leverage $h_{jj} \doteq 1$ implies that $\widehat{y}_j \approx y_j$ — need one parameter to fit this case.

- \square As $\operatorname{tr}(H) = \sum_{j=1}^n h_{jj} = p$, the average h_{jj} is p/n. If $h_{jj} > 2p/n$, then jth case should be checked (rule of thumb), e.g. by refitting without (x_j, y_j) .
- \square Let \widehat{y}_{-j} be fitted values for (all) data when (x_j,y_j) is dropped and use Cook's distance

$$C_j = \frac{1}{ps^2} (\widehat{y} - \widehat{y}_{-j})^{\mathrm{T}} (\widehat{y} - \widehat{y}_{-j}) = \dots = \frac{r_j^2 h_{jj}}{p(1 - h_{jj})}$$

to measure the difference between \widehat{y} and \widehat{y}_{-j} .

- \square Large C_j implies large r_j and/or large h_{jj} .
- \square Cases with $C_j > 8/(n-2p)$ worth a closer look (rule of thumb).
- \Box High leverage and/or influence need not be bad, just need to be aware of it.
- ☐ These ideas are not very useful in large samples, since the plots become uninformative.

Regression Methods

Response transformation

- \square Linear model for y may be better applied for some transformation g(y), especially if some y are much larger than others, or the variance is non-constant.
- Survival times y_{ptj} in 10-hour units of animals in a 3×4 factorial experiment with four replicates, with (below) average (standard deviation) for the poison \times treatment combinations:
 - generally see higher SD and mean together,
 - times must be positive, so linear model inappropriate?

Treatment	Poison 1	Poison 2	Poison 3
Α	0.31, 0.45, 0.46, 0.43	0.36, 0.29, 0.40, 0.23	0.22, 0.21, 0.18, 0.23
В	0.82, 1.10, 0.88, 0.72	0.92, 0.61, 0.49, 1.24	0.30, 0.37, 0.38, 0.29
C	0.43, 0.45, 0.63, 0.76	0.44, 0.35, 0.31, 0.40	0.23, 0.25, 0.24, 0.22
D	0.45, 0.71, 0.66, 0.62	0.56, 1.02, 0.71, 0.38	0.30, 0.36, 0.31, 0.33

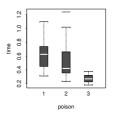
Treatment	Poison 1	Poison 2	Poison 3	Average
А	0.41 (0.07)	0.32 (0.08)	0.21 (0.02)	0.31
В	0.88 (0.16)	0.82 (0.34)	0.34 (0.05)	0.68
C	0.57 (0.16)	0.38 (0.06)	0.24 (0.01)	0.39
D	0.61 (0.11)	0.67 (0.27)	0.33 (0.03)	0.53
Average	0.62	0.55	0.28	0.48

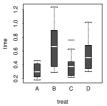
Regression Methods

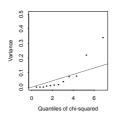
Autumn 2024 - slide 53

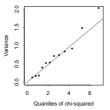
Example: Poison data

Upper panels: dependence of responses on the factor levels. Lower left: χ^2_3 probability plots of the $3s^2_{pt}$, where s^2_{pt} is the sample variance of y_{ptj} . Lower right: same for y^{-1}_{ptj} .









Regression Methods

Box-Cox transformation

 \square For y > 0, the **Box–Cox transformation**

$$y^{(\lambda)} = \begin{cases} \frac{y^{\lambda} - 1}{\lambda}, & \lambda \neq 0, \\ \log y, & \lambda = 0, \end{cases}$$

includes the inverse $(\lambda=-1)$, $\log\ (\lambda=0)$, cube and square roots $(\lambda=\frac{1}{3},\frac{1}{2})$, original scale $(\lambda=1)$ and square $(\lambda=2)$; sometimes map $y\mapsto y+c>0$.

☐ Suppose normal linear model

$$y^{(\lambda)} \sim \mathcal{N}_n(X\beta, \sigma^2 I_n)$$

applies for some β , σ and λ to be determined. Here X contains 1_n , so use of $y^{(\lambda)}$ just changes intercept and rescales β and σ .

 \square Use profile log likelihood for λ to choose 'best' transformation (usually from list above to aid interpretation).

Interpretation of β depends on λ , so usually we ignore the fact that λ was estimated, unless we are not interested in β (e.g., when performing 'automatic' prediction).

Regression Methods

Autumn 2024 - slide 55

Example: Poison data

☐ Fits of two-way layout model, with interaction:

$$y_{tpj}^{(\lambda)} \sim \mathcal{N}(\mu + \alpha_t + \beta_p + \gamma_{tp}, \sigma^2), \quad t = 1, 2, 3, 4, \ p = 1, 2, 3, \ j = 1, 2, 3, 4.$$

 $\ \square$ Analyses of variance with responses y and y^{-1} . For MS and F read 'Mean square' and 'F statistic'.

 \square The terms explain appreciably more of the variation of y^{-1} , suggesting that this is a preferable choice of response.

Term	df	Response y			Re	sponse y	I^{-1}	
		SS	MS	F		SS	MS	F
Poisons	2	1.033	0.517	23.22	,	34.88	17.44	72.63
Treatments	3	0.921	0.307	13.81		20.41	6.80	28.34
$Treatments \times Poisons$	6	0.250	0.042	1.87		1.57	0.26	1.09
Residual	36	0.801	0.022			8.64	0.24	

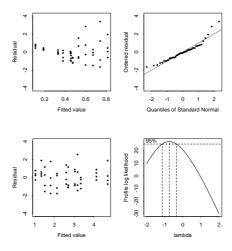
Regression Methods

Example: Poison data

Top: residuals for model without interactions γ_{tp} ; clearly problematic.

Lower right: profile log likelihood for Box–Cox λ , showing 95% confidence interval.

Lower left: residuals for the two-way layout model (no interactions) for 1/y.



Regression Methods

Autumn 2024 - slide 57

Summary on model-checking

□ Recall the distinction between primary and secondary assumptions. Use of the standard linear model when the secondary assumptions fail leads to inefficient estimation and over-confident uncertainty assessment, but is not usually disastrous per se.

☐ When they fail . . .

- **Linearity** (primary): add terms (e.g., x^2) to the model, transform the covariate (e.g., to $\log x$), or question the basic setup;
- Constant variance (secondary): use a response transformation (below), weighted least squares, or question primary aspects. Non-constant variance affects uncertainty assessment, but not estimation;
- Lack of correlation (independence) (secondary): use a correlated error model (e.g., time series or random effects). Dependence affects uncertainty assessment, but not estimation;
- normality (secondary): often does not matter, because the CLT applies to the estimators. It does matter for prediction, which is affected by the distribution of individual responses;
- ☐ Checking leverage and influence may be useful in small and moderate samples, but rarely in large samples. In any case, automatic dropping of outlying and/or influential cases is dangerous!

Regression Methods

Goals	
	What to do faced with a set of data?
	Two main aims:
	 understand (science) — maybe have prior idea/hypotheses on how response depends on explanatory variables. Interpretation is key.
	- predict/control (technology) — don't really care how y depends on X . Interpretation not critical (though this describes only prediction in the narrowest of senses).
	There is no reason that a single model will do both, or even that there must be a single 'best' model:
	 maybe two models with different interpretations both fit about equally well, and then future work might aim to choose between them;
	 prediction with a mixture of models might be better than using a single model.
Regression Methods Autumn 2024 – slide 60	
Meta-algorithm	
	Collect data intended to answer question of interest;
	examine data (graphs, look for outliers, problems with sampling scheme);
	<pre>choose/construct response variable (transformations? independence?);</pre>
	consider what models are coherent with context of problem (limiting properties, units, similar problems/datasets, covariates that must be included,);
	iterate:
	 fit models, compare quality of fits;
	– check interpretations of $\widehat{eta}, \widehat{\sigma}^2$ and
	check fit (diagnostics, outliers,)
	until satisfied; finally
	give conclusions —careful interpretation of best model(s) in terms of original problem, consider deficiencies, and explain what extra data might overcome them.
Regression Methods Autumn 2024 – slide 61	
Initial examination of data	
	Plot y against covariates, look for outliers, non-constant variance, nonlinearity, etc.
	Plot covariates against each other, look for dependence.
	Try to understand covariates (e.g., dimensions), are transformations needed?
	May need to reduce dimension of X by regularisation — many ways to do this (later).

Regression Methods

Albert Einstein (1879–1955)



'Everything should be made as simple as possible, but no simpler.'

Regression Methods

Autumn 2024 - slide 64

William of Occam (?1285-1347/9)



Occam's razor: *Pluralitas non est ponenda sine necessitate*: entities should not be multiplied beyond necessity.

Regression Methods

Automatic variable selection

- \square Assume linear model $E(y) = X\beta$
- \square 2^p possible subsets of columns of X, plus transformations, ...
- \square Example: p = 17 gives 131072 possible subsets of variables
- ☐ Fast algorithms (e.g., branch and bound, leaps in R) exist visit them all or just subsets (e.g., stepwise), but we need criteria for comparing models.
- ☐ Many proposals for model comparison
 - cross-validation,
 - information criteria (AIC, AIC_c, BIC, NIC, TIC, ...)
 - Mallow's C_p ,
 - ...
- ☐ Most involve minimising estimated prediction error for future data *like those observed*!

Regression Methods

Autumn 2024 - slide 66

Prediction error

True model $y \sim (\mu, \sigma^2 I_n)$, we assume (perhaps incorrectly) that $\mu = X\beta$, fit $X_{n \times p}$ and obtain fitted value

$$X\widehat{\beta} = Hy \sim (H\mu, \sigma^2 H).$$

- □ Terminology
 - the true model has $\mu = X\beta$ and all $\beta_r \neq 0$;
 - a correct model has $\mu = X\beta$ but some $\beta_r = 0$;
 - a wrong model has $\mu \notin \operatorname{span}(X)$;

so $(I_n - H)\mu = 0$ if the model is true or correct, and $(I_n - H)\mu \neq 0$ if it is wrong.

 \Box The **prediction error** for an independent dataset y_+ with mean vector μ is

$$\Delta = n^{-1} \mathbf{E} \left\{ (y_{+} - X \widehat{\beta})^{\mathrm{T}} (y_{+} - X \widehat{\beta}) \right\} = \begin{cases} n^{-1} \mu^{\mathrm{T}} (I - H) \mu + (1 + p/n) \sigma^{2}, & \text{wrong,} \\ (1 + q/n) \sigma^{2}, & \text{true,} \\ (1 + p/n) \sigma^{2}, & \text{correct,} \end{cases}$$

where $E(\cdot)$ is over both y_+ and y and $p \ge q = \#\{\beta_r : \beta_r \ne 0\}$ when $\mu \in \operatorname{span}(X)$.

 \square In principle we should write $\Delta \equiv \Delta(X)$.

Regression Methods

Note: Computation of Δ

Let $y \sim (\mu, \sigma^2 I)$ and fit $X\beta$, obtaining fitted value

$$X\widehat{\beta} = Hy \sim (H\mu, \sigma^2 H),$$

where $H\mu=\mu$, i.e., $(I-H)\mu=0$ if $\mu\in \mathrm{span}(X)$, but otherwise $(I-H)\mu\neq 0$.

We have a new data set $y_+ \sim (\mu, \sigma^2 I)$, and we compute the average error in predicting y_+ using $X\widehat{\beta}$, i.e.,

$$\Delta = n^{-1} \mathbf{E} \left\{ (y_+ - X\widehat{\beta})^{\mathrm{T}} (y_+ - X\widehat{\beta}) \right\}.$$

Let $e_+=y_+-X\widehat{eta}$ and note that as the trace of a scalar is the scalar and trace is a linear operator,

$$E(e_{+}^{T}e_{+}) = E\{tr(e_{+}^{T}e_{+})\} = E\{tr(e_{+}e_{+}^{T})\} = tr\{E(e_{+}e_{+}^{T})\} = tr\{var(e_{+}) + E(e_{+})E(e_{+})^{T}\}\}$$

Now as y_+ and y are independent and $\mathrm{var}(X\widehat{\beta}) = \sigma^2 H$, we have

$$y_+ - X\widehat{\beta} \sim (\mu - H\mu, \sigma^2 I + \sigma^2 H),$$

so the computation above gives

$$E\left\{ (y_{+} - X\widehat{\beta})^{T}(y_{+} - X\widehat{\beta}) \right\} = tr\{\sigma^{2}(I + H) + (I - H)\mu\mu^{T}(I - H)\} = \sigma^{2}(n + p) + \mu^{T}(I - H)\mu,$$

because tr(I + H) = n + p and I - H is symmetric and idempotent, giving

$$\Delta = \begin{cases} n^{-1}\mu^{\mathrm{\scriptscriptstyle T}}(I-H)\mu + (1+p/n)\sigma^2, & \text{wrong model,} \\ (1+q/n)\sigma^2, & \text{true model,} \\ (1+p/n)\sigma^2, & \text{correct model.} \end{cases}$$

Regression Methods

Autumn 2024 - note 1 of slide 67

Bias/variance trade-off

- \square Minimising Δ involves balancing the
 - bias $n^{-1}\mu^{\rm T}(I-H)\mu$, which is reduced by including useful terms in X, and
 - variance $(1+p/n)\sigma^2$, which is increased by including useless terms in X.
- \square We would like to minimise Δ , but it depends on the unknown μ and σ .
- \square The **cross-validation** estimator of Δ splits the data into X', y' and X^*, y^* , then
 - for each possible subset S of columns of X^* :
 - \triangleright compute $\widehat{\beta}_{\mathcal{S}}^*$ by regressing y^* on $X_{\mathcal{S}}^*$;
 - \triangleright use $\widehat{\beta}_{\mathcal{S}}^*$ to estimate the prediction error for \mathcal{S} by

$$\widehat{\Delta}_{\mathcal{S}} = (n')^{-1} (y' - X_{\mathcal{S}}' \widehat{\beta}_{\mathcal{S}}^*)^{\mathrm{T}} (y' - X_{\mathcal{S}}' \widehat{\beta}_{\mathcal{S}}^*);$$

- finally choose the set of columns S for which $\widehat{\Delta}_S$ is minimised.
- \square This estimator depends on the split, and since $X' \neq X^*$ in general, $\widehat{\Delta}_{\mathcal{S}}$ does not estimate $\Delta_{\mathcal{S}}$, so it would be preferable to use the entire dataset . . .

Regression Methods

Leave-one-out cross-validation

☐ Simplest way to use entire dataset is leave-one-out cross-validation (CV), minimising

$$n\widehat{\Delta}_{\mathrm{CV}} = \mathrm{CV} = \sum_{j=1}^{n} (y_j - x_j^{\mathrm{T}}\widehat{\beta}_{-j})^2,$$

where $\widehat{\beta}_{-j}$ is estimate computed without (x_j, y_j) .

 \Box This seems to require n fits, but the lemma below implies that with one fit we have

$$CV = \sum_{j=1}^{n} \frac{(y_j - x_j^{\mathrm{T}} \widehat{\beta})^2}{(1 - h_{jj})^2}.$$

Lemma 14 For a fit $\widehat{y} = Hy$ where H has jth diagonal element h_{jj} and $\widehat{y}_{j,-j}$ is the fitted value for y_j obtained when (x_j, y_j) is dropped,

$$y_j - \widehat{y}_{j,-j} = \frac{y_j - \widehat{y}_j}{1 - h_{jj}},$$

and therefore

$$\sum_{j=1}^{n} (y_j - \widehat{y}_{j,-j})^2 = \sum_{j=1}^{n} \frac{(y_j - \widehat{y}_j)^2}{(1 - h_{jj})^2}.$$

Regression Methods

Autumn 2024 - slide 69

Note to Lemma 14

- \square Consider any linear fit $\widehat{y} = Hy$, and note that $\widehat{y}_j = \sum_{i=1}^n h_{ji} y_i$.
- \square Now suppose we leave out (x_i, y_i) and compute the corresponding (penalized) estimate

$$\widehat{\beta}_{-j} = \operatorname{argmin}_{\beta} \quad \sum_{i \neq j} (y_i - x_i^{\mathrm{T}} \beta)^2 + \lambda p(\beta),$$

and fitted value $y_j^* = \widehat{y}_{j,-j} = x_j^{\mathrm{T}} \widehat{\beta}_{-j}$ corresponding to x_j .

Inserting (x_j, y_j^*) back into the dataset used to compute $\widehat{\beta}_{-j}$ changes nothing, because $(y^* - x_j^{\mathrm{T}} \widehat{\beta}_{-j})^2 = 0$ and $p(\beta)$ does not depend on the data. For this new dataset,

$$y_j^* = \sum_{i \neq j} h_{ji} y_i + h_{jj} y_j^* = \sum_{i=1}^n h_{ji} y_i + h_{jj} (y_j^* - y_j) = \widehat{y}_j + h_{jj} (y_j^* - y_j)$$

SO

$$y_j - y_j^* = y_j - \widehat{y}_j + h_{jj}(y_j - y_j^*),$$

leading to

$$y_j - y_j^* = y_j - \widehat{y}_{j,-j} = \frac{y_j - \widehat{y}_j}{1 - h_{jj}},$$

and thus to the given formula.

Regression Methods

Autumn 2024 - note 1 of slide 69

Generalized cross-validation

- \square Leave-one-out CV can be unstable if some of the h_{jj} are large.
- \square Generalised cross-validation (GCV) replaces all the h_{ij} by their average $\mathrm{tr}(H)/n = p/n$, giving

GCV =
$$\sum_{j=1}^{n} \frac{(y_j - x_j^{\mathrm{T}} \widehat{\beta})^2}{(1 - p/n)^2}$$
,

and hence

$$E(GCV) = \mu^{T}(I - H)\mu/(1 - p/n)^{2} + n\sigma^{2}/(1 - p/n) \approx n\Delta.$$

- \square Often choose the model that minimises GCV or CV.
- □ Note that these only require the second-order assumptions.

Regression Methods

Autumn 2024 - slide 70

Note: Properties of GCV

We have $(1-p/n)^2 \text{GCV} = e^{\mathrm{T}} e$ where $e=y-X \widehat{\beta} = (I-H)y \sim ((I-H)\mu, (I-H)\sigma^2)$, and

$$E(e^{T}e) = E\{tr(ee^{T})\} = tr\{E(e)E(e)^{T} + var(e)\} = \mu^{T}(I - H)\mu + \sigma^{2}tr(I - H).$$

Now note that ${\rm tr}(I-H)=n-p$ and divide by $(1-p/n)^2$ to give (almost) the required result, for which we need also $(1-p/n)^{-1}\approx 1+p/n$, for $p\ll n$.

Regression Methods

Autumn 2024 - note 1 of slide 70

Akaike information criterion

☐ The above arguments apply only to least squares estimators. More generally, we could aim to minimise the Kullback–Leibler discrepancy

$$D(f_{\theta}, g) = \int \log \left\{ \frac{g(y)}{f(y; \theta)} \right\} g(y) \, \mathrm{d}y \ge 0,$$

between candidate model $f_{\theta} \equiv f(y; \theta)$ and true model g, based on $Y_1, \dots, Y_n \overset{\text{iid}}{\sim} g$.

- \square Suppose that θ_g minimises $D(f_\theta,g)$ within the family of candidate models, and is estimated by the MLE $\widehat{\theta}$, with log likelihood $\widehat{\ell}$.
- \square We suppose there is an independent sample $Y_1^+,\ldots,Y_n^+\stackrel{\mathrm{iid}}{\sim} g$ and aim to estimate

$$E_g\left(E_g^+\left[\sum_{j=1}^n \log\left\{\frac{g(Y_j^+)}{f(Y_j^+;\widehat{\theta})}\right\}\right]\right) = nE_g\left\{D(f_{\widehat{\theta}},g)\right\};$$
(2)

the outer expectation is over the distribution of $\widehat{\theta}$, which is independent of Y^+ .

 $\ \square$ After tedious expansions we end up trying to minimise the **Akaike information criterion**

$$\label{eq:AIC} \mathrm{AIC} = -2\widehat{\ell} + 2p \qquad (\equiv n\log\mathrm{RSS} + 2p \text{ in linear model}).$$

Regression Methods

Note: Derivation of AIC

 \square Taylor series expansion shows that $\log f(y;\widehat{\theta})$ approximately equals

$$\log f(y; \theta_g) + (\widehat{\theta} - \theta_g)^{\mathrm{T}} \frac{\partial \log f(y; \theta_g)}{\partial \theta} + \frac{1}{2} (\widehat{\theta} - \theta_g)^{\mathrm{T}} \frac{\partial^2 \log f(y; \theta_g)}{\partial \theta \partial \theta^{\mathrm{T}}} (\widehat{\theta} - \theta_g),$$

and as θ_q minimizes $D(f_\theta, g)$,

$$\int \frac{\partial \log f(y; \theta_g)}{\partial \theta} g(y) \, dy = 0.$$

Hence taking expectation over Y_1^+, \dots, Y_n^+ , we get

$$nD(f_{\widehat{\theta}}, g) = n \int \log \left\{ \frac{g(y)}{f(y; \widehat{\theta})} \right\} g(y) \, dy \doteq nD(f_{\theta_g}, g) + \frac{1}{2} \operatorname{tr} \left\{ (\widehat{\theta} - \theta_g) (\widehat{\theta} - \theta_g)^{\mathrm{T}} I_g(\theta_g) \right\},$$

where we have used the fact that the trace of a scalar is itself.

 \square Expectation over the distribution of $\widehat{\theta}$ gives its variance matrix, $I_g(\theta_g)^{-1}K(\theta_g)I_g(\theta_g)^{-1}$, and hence

$$nE_g\left\{D(f_{\widehat{\theta}},g)\right\} \doteq nD(f_{\theta_g},g) + \frac{1}{2}\operatorname{tr}\left\{I_g(\theta_g)^{-1}K(\theta_g)\right\},\tag{3}$$

where the second term penalizes the dimension p of θ . The first term here is O(n) but the second is O(p). When $f_{\theta g}=g$, $I_g(\theta_g)=K(\theta_g)$ so $\operatorname{tr}\left\{I_g(\theta_g)^{-1}K(\theta_g)\right\}=p$.

□ To build an estimator, note that $\int \log g(y) g(y) dy$ is constant and can be ignored. Now $\ell(\widehat{\theta}) = \ell(\theta_a) + \{\ell(\widehat{\theta}) - \ell(\theta_a)\}$, so

$$\begin{split} \mathbf{E}_g \left\{ -\ell(\widehat{\theta}) \right\} &= -\mathbf{E}_g \left\{ \ell(\theta_g) + \frac{1}{2}W(\theta_g) \right\} \\ &\doteq nD(f_{\theta_g}, g) - \frac{1}{2} \mathrm{tr} \left\{ I(\theta_g)^{-1} K(\theta_g) \right\} - n \int \log g(y) \, g(y) \, dy, \end{split}$$

where we have used the fact that under the wrong model, the likelihood ratio statistic $W(\theta_g)$ has mean approximately $\operatorname{tr} \left\{ I(\theta_g)^{-1} K(\theta_g) \right\}$. Hence $-\ell(\widehat{\theta})$ tends to underestimate $nD(f_{\theta_g},g) - n \int \log g(y) \, g(y) \, dy$. On reflection this is obvious, because $\ell(\widehat{\theta}) \geq \ell(\theta_g)$ by definition of $\widehat{\theta}$. As p increases, so will the extent of overestimation.

AIC =
$$2\{-\ell(\widehat{\theta}) + p\}$$
, NIC = $2\{-\ell(\widehat{\theta}) + \operatorname{tr}(\widehat{J}^{-1}\widehat{K})\}$; (4)

another possibility is $BIC = -2\ell(\widehat{\theta}) + p \log n$.

- \square The model is chosen to minimize AIC, say, with the factor 2 putting differences of AIC on the same scale as likelihood ratio statistics. Such criteria are used far beyond random samples, and even in cases where the theory above doesn't work.
- $\hfill\Box$ In particular, the maximised log-likelihood for a normal-theory linear model with residual sum of squares RSS can be shown to be

$$-\frac{n}{2}\log(2\pi\widehat{\sigma}) - \frac{n}{2} \equiv -\frac{n}{2}\log RSS + \text{constants},$$

which leads to the formula given on the slide.

Regression Methods

Autumn 2024 - note 1 of slide 71

Other model selection criteria

☐ **'Corrected' AIC** for (normal-theory) regression problems:

$$AIC_{c} \equiv n \log \hat{\sigma}^{2} + n \frac{1 + p/n}{1 - (p+2)/n}.$$

□ Bayes' information criterion

$$\mathrm{BIC} = -2\widehat{\ell} + p\log n.$$

 \square Mallows C_p :

$$C_p = \frac{SS_p}{s^2} + 2p - n,$$

where SS_p is RSS for fitted model and s^2 estimates σ^2 .

- \square When the true model is a candidate and $n \to \infty$,
 - AIC is inconsistent it will not choose the true model with probability one, but tends to pick a more complex model;
 - $\quad AIC_c$ is also inconsistent but gives better results in finite samples;
 - BIC is **consistent** it chooses the true model with probability $\rightarrow 1$.

These results suppose that the models are fixed, but in practice we also have $p\to\infty$ when $n\to\infty$, because we fit ever more complex models when we have more data.

Regression Methods

Autumn 2024 - slide 72

Simulation experiment

Number of times models were selected using various model selection criteria in 50 repetitions using simulated normal data for each of 20 design matrices. The true model has p=3.

n		Number of covariates						
		1	2	3	4	5	6	7
10	C_p		131	504	91	63	83	128
	BIC		72	373	97	83	109	266
	AIC		52	329	97	91	125	306
	$\mathrm{AIC}_{\mathrm{c}}$	15	398	565	18	4		
20	C_p		4	673	121	88	61	53
	BIC		6	781	104	52	30	27
	AIC		2	577	144	104	76	97
	$\mathrm{AIC_c}$		8	859	94	30	8	1
40	C_p			712	107	73	66	42
	BIC			904	56	20	15	5
	AIC			673	114	90	69	54
	$\mathrm{AIC_c}$			786	105	52	41	16

Regression Methods

Ste	epwise methods
	In principle we might wish to fit all 2^p possible choices of covariates, but in practice this is possible only for 'modest' p , using leaps or similar methods (or approximations).
	When p is too large for exhaustive searches, we instead consider subsets of the models, using the methods below (or variants).
	Forward selection: starting from the model with a constant only,
	1. add each remaining term separately to the current model;
	2. if none of these terms improves the fit, stop; otherwise
	3. update the current model to include the most useful new term; go to 1
	Backward elimination: starting from the model with all terms,
	1. if all terms are 'useful', stop; otherwise
	2. update current model by dropping the 'least useful' term; go to 1
	Stepwise: starting from an arbitrary model,
	 consider three options—add a term, delete a term, swap a term in the model for one not in the model;
	2. if model unchanged, stop; otherwise go to 1
	'Useful' means 'reduces the AIC' (but in the past meant 'is significant using an F test').

Regression Methods Autumn 2024 – slide 74

Stepwise methods: Comments

□ Original formulation of stepwise used F tests (or even arbitrary numbers!) to assess significance, but this finds spurious models.
 □ Systematic search minimising AIC or similar over all possible models is preferable, but is often infeasible.
 □ Compare AICs for different models at each step—i.e., use AIC (or AIC_c) as objective function.

☐ Important not to fixate on a specific model, or assume that there is a single 'best' model, but to consider any models whose AIC is within (say) 2 of the minimum — especially if the interpretations of competing models differ.

Regression Methods

Example: Nuclear power stations > nuclear cost date t1 t2 cap pr ne ct bw cum.n pt 1 460.05 68.58 14 46 687 0 1 0 0 2 452.99 67.33 10 73 1065 0 0 1 0 1 0 3 443.22 67.33 10 85 1065 1 0 1 0 1 0 4 652.32 68.00 11 67 1065 0 1 1 0 12 0 5 642.23 68.00 11 78 1065 1 1 1 0 6 345.39 67.92 13 51 514 0 1 1 0 3 0 7 272.37 68.17 12 50 822 0 0 0 0 5 0 8 317.21 68.42 14 59 457 0 0 0 0 1 0 9 457.12 68.42 15 55 822 1 0 0 0 5 0 10 690.19 68.33 12 71 792 0 1 1 1 32 270.71 67.83 7 80 886 1 0 0 1 11 1

Regression Methods

Autumn 2024 - slide 76

Example:	Nuclear power stations
	Full model

	Full mode	·	Backward		Forward		
•	Est (SE)	t	Est (SE)	t	Est (SE)	t	
Constant	-14.24 (4.229)	-3.37	-13.26 (3.140)	-4.22	-7.627 (2.875)	-2.66	
date	0.209 (0.065)	3.21	$0.212 \ (0.043)$	4.91	$0.136 \ (0.040)$	3.38	
log(T1)	0.092 (0.244)	0.38					
log(T2)	$0.290 \ (0.273)$	1.05					
log(cap)	$0.694 \ (0.136)$	5.10	$0.723 \ (0.119)$	6.09	$0.671 \ (0.141)$	4.75	
PR	-0.092 (0.077)	-1.20					
ΝE	$0.258 \ (0.077)$	3.35	0.249 (0.074)	3.36			
CT	$0.120 \ (0.066)$	1.82	0.140 (0.060)	2.32			
ВW	$0.033 \ (0.101)$	0.33					
log(N)	-0.080 (0.046)	-1.74	-0.088 (0.042)	-2.11			
PΤ	$-0.224 \ (0.123)$	-1.83	-0.226 (0.114)	-1.99	$-0.490 \ (0.103)$	-4.77	
s (df)	0.164 (21)		0.159 (25)	0.159 (25)		0.195 (28)	

Regression Methods

M-estimation

- \Box The least squares estimates are linear in y and therefore very sensitive to outliers.
- \square When $y_i \mapsto y_i + c$,

$$\widehat{\beta} = \sum_{j=1}^{n} (X^{\mathrm{T}} X)^{-1} x_j y_j \mapsto \sum_{j=1}^{n} (X^{\mathrm{T}} X)^{-1} x_j y_j + (X^{\mathrm{T}} X)^{-1} x_i c = \widehat{\beta} + (X^{\mathrm{T}} X)^{-1} x_i c,$$

which could be arbitrarily far from $\widehat{\beta}$.

 \Box Try and fix this by replacing

$$\min_{\beta} \sum_{j=1}^{n} (y_j - x_j^{\mathrm{T}} \beta)^2 \qquad \text{by} \qquad \min_{\beta} \sum_{j=1}^{n} \rho \left\{ (y_j - x_j^{\mathrm{T}} \beta) / \sigma \right\},$$

for function $\rho(\cdot)$ that will give a more robust **M**(aximum likelihood-like)-estimator, or equivalently solving the $p \times 1$ system of estimating equations

$$\frac{1}{\sigma} \sum_{j=1}^{n} x_j \rho' \left\{ (y_j - x_j^{\mathrm{T}} \beta) / \sigma \right\} = X^{\mathrm{T}} \rho' = 0$$

say, where $\rho'_{n\times 1}$ has $j\mathrm{th}$ element $\mathrm{d}\rho(u)/\mathrm{d}u$ for $u=(y_j-x_j^\mathrm{T}\beta)/\sigma$.

Regression Methods

Autumn 2024 - slide 79

Choice of ρ

 \square Choose $\rho(u)$ to have desirable properties, e.g., to downweight outliers:

 $\rho(u) = u^2/2$ (normal errors),

 $ho(u) \ = \ |u|$ (Laplace errors),

 $\rho(u) = \nu \log(1 + u^2/\nu)/2$ (t_{ν} errors),

 $\rho(u) = \begin{cases} u^2/2, & |u| < c, \\ c(2|u|-c)/2, & \text{otherwise,} \end{cases}$ (Huber function).

- \square The function $\rho'(u)$ is also called the **influence function** of the estimator, as its value determines what influence an observation at u has on the estimator:
 - Huber $\rho'(u)$ is bounded,
 - t_{ν} function is bounded and redescending, as $\lim_{u\to\pm\infty} \rho'(u) = 0$;
 - Tukey's biweight

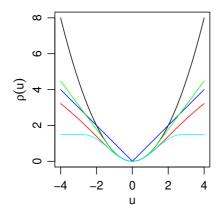
$$\rho'(u) = u \left\{ 1 - (u/c)^2 \right\}^2 I(|u| < c),$$

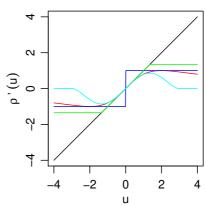
which gives $\rho'(u)=0$ when |u|>c, is also redescending, giving no weight to observations outside $\pm c$.

Regression Methods

 ρ and ρ'

Functions ρ and ρ' for least squares (black), t_5 (red), Laplace (blue), Huber (green) and biweight (cyan) estimators.





Regression Methods

Autumn 2024 - slide 81

Estimation

 \square We need to solve

$$X^{\mathrm{T}}\rho'=0,$$

where ρ' has jth element

$$\sigma^{-1}\rho'\{(y_j - x_j^{\mathrm{T}}\beta)/\sigma\} \propto \frac{\rho'\{(y_j - x_j^{\mathrm{T}}\beta)/\sigma\}}{y_j - x_j^{\mathrm{T}}\beta} \times (y_j - x_j^{\mathrm{T}}\beta) = w_j(\beta, \sigma)(y_j - x_j^{\mathrm{T}}\beta),$$

say, so we write the estimating equation as

$$X^{\mathrm{T}}W(y - X\beta) = 0,$$

with $W = \operatorname{diag}\{w_1(\beta, \sigma), \dots, w_n(\beta, \sigma)\}.$

- \square We use **iterative weighted least squares**: choose some initial $\tilde{\beta}$ and σ , then iterate to convergence the steps
 - compute W using the current $\tilde{\beta}$,
 - compute the weighted least squares estimate,

$$\tilde{\beta} = (X^{\mathrm{T}}WX)^{-1}X^{\mathrm{T}}Wy.$$

 \square Estimate σ using median absolute deviation of residuals $y_j - x_j^{\mathrm{T}} \tilde{\beta}$ at each iteration, or similar robust scale estimate.

Regression Methods

M-estimator variance

 \square Estimator $\tilde{\beta}$ is solution to $p \times 1$ system of equations

$$g(y;\beta) = X^{\mathrm{T}} \rho' = 0.$$

 $\ \square$ Can show that if the estimating function g is **unbiased**, i.e.

$$E\{g(Y;\beta);\beta\}=0,$$
 for any β ,

then under mild regularity conditions

$$\tilde{\beta} \sim \mathcal{N}_p \left(\beta, \mathrm{E} \left\{ -\frac{\partial g(Y; \beta)}{\partial \beta^{\mathrm{T}}} \right\}^{-1} \mathrm{var} \left\{ g(Y; \beta) \right\} \mathrm{E} \left\{ -\frac{\partial g(Y; \beta)}{\partial \beta^{\mathrm{T}}} \right\}^{-1} \right).$$

This is another sandwich variance matrix, with

$$\mathrm{E}\left\{-\frac{\partial g(Y;\beta)}{\partial \beta^{\mathrm{T}}}\right\} = X^{\mathrm{T}}W_{1}X, \quad \operatorname{var}\left\{g(Y;\beta)\right\} = X^{\mathrm{T}}W_{2}X,$$

so if $W_1=A(\sigma)I_n$, $W_2=\sigma^2B(\sigma)I_n$, then

$$\operatorname{var}(\tilde{\beta}) \doteq \sigma^2 (X^{\mathrm{T}} X)^{-1} \times B(\sigma) / A(\sigma)^2.$$

Regression Methods

Note: Sandwich matrix I

 \square The $p \times 1$ estimating function is

$$g(y; \beta) = \sum_{j=1}^{n} x_j \rho' \left(\frac{y_j - x_j^{\mathrm{T}} \beta}{\sigma} \right),$$

and unbiasedness implies that if the individual densities are $\sigma^{-1}f\{(y_j-x_j^{\mathrm{T}}\beta)/\sigma\}$, then

$$0 = \mathrm{E}\left\{g(y;\beta)\right\} = \sum_{j=1}^{n} x_{j} \int \rho' \left(\frac{y_{j} - x_{j}^{\mathrm{T}}\beta}{\sigma}\right) \sigma^{-1} f\left(\frac{y_{j} - x_{j}^{\mathrm{T}}\beta}{\sigma}\right) \, \mathrm{d}y_{j} = X^{\mathrm{T}} a_{n \times 1},$$

say, where a_j is the jth integral above, and setting $u=(y_j-x_j^{\mathrm{T}}\beta)/\sigma$ shows that all the a_j equal

$$\int \rho'(u) f(u) du = 0; \tag{5}$$

this is true by symmetry if the error distribution and ho' are symmetric around the origin. Now

$$\frac{\partial g(y;\beta)}{\partial \beta^{\mathrm{T}}} = -\frac{1}{\sigma} \sum_{j=1}^{n} x_j x_j^{\mathrm{T}} \rho'' \left(\frac{y_j - x_j^{\mathrm{T}} \beta}{\sigma} \right),$$

whose expectation is (using the same transformation)

$$\mathbb{E}\left\{\frac{\partial g(y;\beta)}{\partial \beta^{\mathrm{T}}}\right\} = -\frac{1}{\sigma} \sum_{j=1}^{n} x_{j} x_{j}^{\mathrm{T}} \mathbb{E}\left\{\rho''\left(\frac{Y_{j} - x_{j}^{\mathrm{T}}\beta}{\sigma}\right)\right\} \\
= -\frac{1}{\sigma} \sum_{j=1}^{n} x_{j} x_{j}^{\mathrm{T}} \int \rho''(u) f(u) \, \mathrm{d}u = -\frac{1}{\sigma} X^{\mathrm{T}} X A(\sigma),$$

say.

 \square The components of these sums are independent, so

$$\operatorname{var}\left\{g(Y;\beta)\right\} = \operatorname{var}\left\{\sum_{j=1}^{n} x_{j} \rho'\left(\frac{Y_{j} - x_{j}^{\mathrm{T}} \beta}{\sigma}\right),\right\} = \sum_{j=1}^{n} x_{j} x_{j}^{\mathrm{T}} \operatorname{var}\left\{\rho'\left(\frac{Y_{j} - x_{j}^{\mathrm{T}} \beta}{\sigma}\right)\right\},$$

where the substitution $u=(y_j-x_j^{\mathrm{T}}\beta)/\sigma$ and (??) show that the variance term can be written as

$$\operatorname{var}\left\{\rho'\left(\frac{Y_j - x_j^{\mathrm{T}}\beta}{\sigma}\right)\right\} = \int \rho'(u)^2 f(u) \, \mathrm{d}u = B(\sigma).$$

☐ The sandwich variance formula is therefore

$$\left\{-\frac{1}{\sigma}X^{\mathrm{\scriptscriptstyle T}}XA(\sigma)\right\}^{-1}X^{\mathrm{\scriptscriptstyle T}}XB(\sigma)\left\{-\frac{1}{\sigma}X^{\mathrm{\scriptscriptstyle T}}XA(\sigma)\right\}^{-1} = (X^{\mathrm{\scriptscriptstyle T}}X)^{-1}\times\frac{\sigma^2B(\sigma)}{A(\sigma)^2}.$$

The variance of the LSE is $var(Y_j)(X^TX)^{-1}$, so the asymptotic relative efficiency of the M-estimator based on ρ and the LSE is

$$\frac{\operatorname{var}(Y_j)}{\sigma^2} \times \frac{A(\sigma)^2}{B(\sigma)}.$$

Regression Methods

Autumn 2024 - note 1 of slide 83

Note: Sandwich matrix II

- As a check on this, note that for the normal distribution $\rho'(u) = u$, $f(u) = (2\pi)^{-1}e^{-u^2/2}$, so $A(\sigma) = B(\sigma) = 1$, which gives ARE of 1. If we take $\rho'(u) = \mathrm{sign}(u)$ with the normal density, we have $B(\sigma) = 1$, $A(\sigma) = -2/(2\pi)^{1/2}$, so the sandwich variance formula gives $\sigma^2(X^\mathrm{T}X)^{-1}\pi/2$. So using the ρ -function corresponding to the Laplace distribution when the data are in fact normally distributed leads to an estimator which is $\pi/2 \approx 1.57$ times more variable than would be the case if the appropriate ρ -function were used.
- \square If we take the ho-function ho'(u)=u corresponding to the normal density, and the errors are in fact Laplace, $g(u)=(1/2)e^{-|u|}$, we have

$$A(\sigma) = \int (-1)f(u) du = 1, \quad B(\sigma) = \int u^2 f(u) du = 2$$

and the asymptotic relative efficiency is 1/2.

Regression Methods

Autumn 2024 - note 2 of slide 83

Efficiency

 \square Efficiency of M-estimators of β relative to LSEs of β is

$$\frac{\operatorname{var}(Y_j)}{\sigma^2} \times \frac{A(\sigma)^2}{B(\sigma)};$$

for example, the Huber estimator is 95% efficient if c = 1.345.

- \Box In practice need to balance robustness and efficiency, increasing the latter by increasing c.
- ☐ High numbers of outliers can wreck M-estimators.
- ☐ Highly robust least trimmed squares estimators obtained by minimising

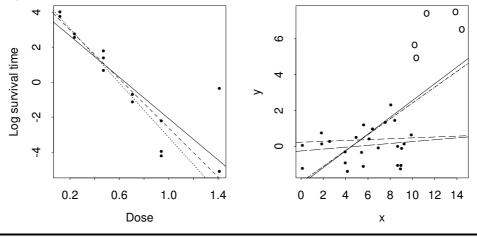
$$\sum_{j=1}^{q} (y_j - x_j^{\mathrm{T}} \beta)_{(j)}^2,$$

where q = |n/2| + |(p+1)/2|.

Regression Methods

Example: Survival data

Left: log survival proportions for rats given doses of radiation, with lines fitted by least squares with (solid) and without (dots) the outlier, and a Huber fit for the entire data (dashes). Right: simulated data with a batch of outliers (circles), and fits by least squares to all data (solid), least squares to good data only (large dash), Huber (dot-dash), biweight (dashes), and least trimmed squares (medium dash). The Huber and biweight fits are the same to plotting accuracy.



Regression Methods

Autumn 2024 - slide 85

Simulation (right-hand panel on slide 85)

Table 1: Bias (standard deviation) of estimators of slope in sample of 25 good data and k outliers, estimated from 200 replications.

\overline{k}	Least squares		M-esti	M-estimation		
	No outliers	With outliers	Huber	Biweight	squares	
1	0.00 (0.07)	0.17 (0.06)	0.07 (0.07)	0.01 (0.07)	-0.01 (0.13)	
2	0.00 (0.07)	0.26 (0.06)	0.13 (0.07)	0.02 (0.09)	0.01 (0.14)	
5	0.00 (0.07)	0.41 (0.05)	0.38 (0.06)	0.19 (0.19)	0.01 (0.14)	
10	0.00 (0.06)	0.48 (0.04)	0.48 (0.04)	0.46 (0.12)	0.05 (0.20)	

Good strategy is initial fit using least trimmed squares, then robust fit using this as starting point.

Regression Methods

Quantile regression

☐ The Laplace distribution has

$$\rho(u) = |u| = uI(u \ge 0) - uI(u < 0),$$

and for continuous Y, the solution to $E\{\rho'(Y-\theta)\}=0$ is the median of Y. Hence

$$\operatorname{argmin} \sum_{j=1}^{n} \rho(y_j - x_j^{\mathrm{T}} \beta)$$

estimates the median of y as a linear function of $X\beta$.

 \square Quantile regression takes $\tau \in (0,1)$ and uses the check function

$$\rho_{\tau}(u) = \tau u I(u \ge 0) - (1 - \tau) u I(u < 0);$$

then

$$\tilde{\beta}_{\tau} = \operatorname{argmin} \sum_{j=1}^{n} \rho_{\tau} (y_j - x_j^{\mathrm{T}} \beta)$$

estimates the τ quantile of y as a linear function of $X\beta$.

- \supset For numerical purposes it may be better to round the cusp of ho.
- \square Note that $ho''_ au(u)=0$, so it's better to bootstrap to find $\mathrm{var}(ildeeta_ au)$.

Regression Methods

Autumn 2024 - slide 87

Expectile regression

- ☐ Quantile regression can be used to estimate value-at-risk in finance settings, but it has the drawback of just counting how many residuals are above/below the quantile.
- ☐ Expectile regression extends the LSE in the same way, taking

$$\rho_{\tau}(y-\theta) = \eta_{\tau}(y-\theta) - \eta_{\tau}(y), \quad \eta_{\tau}(u) = |I(u \le 0) - \tau|u^{2},$$

so $\tau=1/2$ gives the LSE, while taking $\tau>1/2$ leads to a more general form of LSE, with good properties for risk estimation in finance applications (coherent elicitable risk measure).

Regression Methods

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.

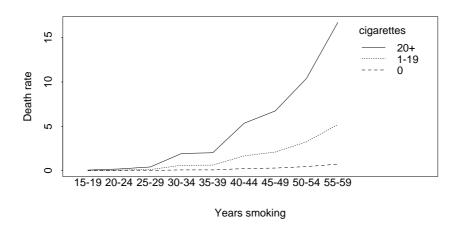
Years of smoking t	Daily cigarette consumption \boldsymbol{d}						
_	Nonsmokers	1–9	10–14	15–19	20–24	25-34	35+
15–19	10366/1	3121	3577	4317	5683	3042	670
20-24	8162	2937	3286/1	4214	6385/1	4050/1	1166
25-29	5969	2288	2546/1	3185	5483/1	4290/4	1482
30-34	4496	2015	2219/2	2560/4	4687/6	4268/9	1580/4
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

_		
Smo	kıng	data
•	5	

- \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
- \square 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

- - Might work as an approximation, but usually extrapolates really badly.
- $\ \square$ 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

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

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

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

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

- \square 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 $Y_j \stackrel{\text{ind}}{\sim} \mathcal{N}(\mu_j, \sigma^2)$ 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 (3), 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 (??) 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 (??) to obtain $\widehat{\beta}$, then setting β equal to $\widehat{\beta}$, and iterating (??) 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 (??), 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.

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 (??), 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 $j{\rm th}$ row of the matrix $X=\partial\eta/\partial\beta^{\rm T}$ is $(\partial\eta_j/\partial\beta_0,\ldots,\partial\eta_j/\partial\beta_{p-1})$, 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 (??) 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 (??) 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 $\widehat{\sigma}^2 = SS(\widehat{\beta})/n$. Although $S^2 = SS(\widehat{\beta})/(n-p)$ 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 $\tilde{\eta}^{\mathrm{T}} = (\tilde{\eta}_1, \dots, \tilde{\eta}_n)$. 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 a < 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 th case (x_j, y_j) is dropped, and r_{Pj} is standardized Pearson residual (see below).

☐ 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 (??) becomes $\widehat{\beta} \doteq (X^{\mathrm{T}}X)^{-1}(X\beta + \tau^2 u)$, where the jth element of u is $\tau^{-1}[1 - \exp\{-(y_j - \eta_j)/\tau\}]$.

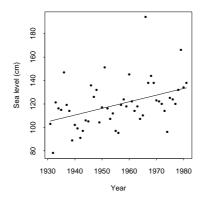
- \square Here it is simplest to fix τ , to obtain $\widehat{\beta}$ by iterating (??) for each fixed value of τ , and then to repeat this over a range of values of τ , giving the profile log likelihood $\ell_p(\tau)$ and hence confidence intervals for τ . Confidence intervals for β_0 and β_1 are obtained from the information matrix.
- With starting value chosen to be the least squares estimates of β , and with $\tau=5$, 19 iterations of (??) 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.

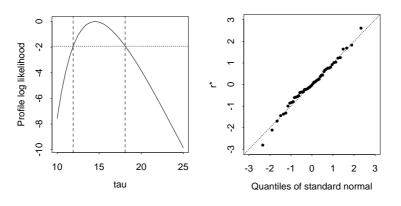


Regression Methods

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.

- $\ \square$ Standard likelihood theory is used for confidence intervals and model comparison.
- ☐ 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

Autumn 2024 - slide 115

2.3 Generalized Linear Models

slide 116

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

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 .
- \square 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,
- ho 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

Autumn 2024 - slide 118

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

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

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

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

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

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 (4) 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.

 \Box If canonical link is used then $\theta_j=x_j^{ \mathrm{\scriptscriptstyle T} }\beta$, 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 (??).

- \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

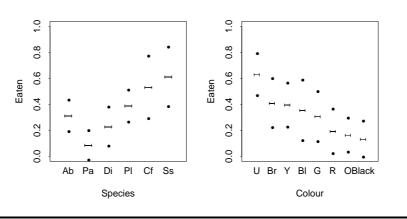
[†] includes *Philaethria dido* also.

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?
- \square 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.

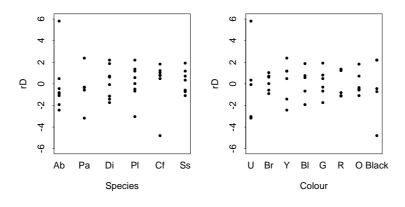
	F	ull data	With	Without outliers				
Terms	df Deviance		df	Deviance				
1	43	134.24	35	73.68				
1+Species	38	114.59	31	46.04				
1+Colour	36	108.46	28	63.20				
1+Species+Colour	31	67.28	24	28.02				

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	stelen	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	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.
- \square 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	Term	df	Deviance
		reduction			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

 \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,

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

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.

m	r	age	stage	grade	xray	acid
6	5	0	1	1	1	1
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
	1	1	1	0	0	0
3 3 3	0	1	0	0	0	1
	0	1	0	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
_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:

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

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

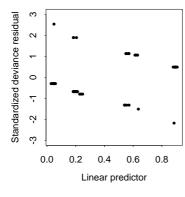
- ☐ Interpretation of this model:
 - for an individual with stage, xray and acid at their lowest levels, the fitted probability of nodal involvement is $e^{-3.05}/(1+e^{-3.05}) \doteq 0.045$ (though there are no such people in the data, so this involves extrapolation);
 - for someone with only $I_{{\tt stage}}=1$, 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.
- ☐ Likewise problems with residuals . . .

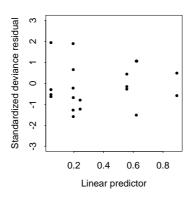
Regression Methods

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 2×2 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.

2.5 Count Data slide 142

Regression Methods

Regression Methods

Autumn 2024 - slide 143

Poisson and multinomial distributions

 \square $Y \sim Pois(\mu)$ 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, $x^{\mathrm{T}}\beta = \eta = \log \mu$.
- \square 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 $\eta = x^{\mathrm{T}}\beta + \log T$
 - offset $\log T$ is fixed part of linear predictor η
- \square If $Y_r \stackrel{\mathrm{ind}}{\sim} \mathrm{Pois}(\mu_r)$, $r=1,\ldots,d$, then the joint distribution of Y_1,\ldots,Y_d given $Y_1+\cdots+Y_d=m$ 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

Autumn 2024 - slide 144

Log-linear and logistic regressions

 \square Special case: if d=2, then

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

 \square If $\mu_1 = \exp(\gamma + x_1^{\mathrm{T}}\beta)$, $\mu_2 = \exp(\gamma + x_2^{\mathrm{T}}\beta)$, 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

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

Regression Methods

Autumn 2024 - slide 147

Premier League data

- □ 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

Premier League data: Analysis of deviance

Poi	sson n	nodel		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.

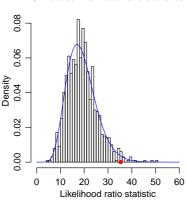
Regression Methods

Autumn 2024 - slide 149

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.

Simulated likelihood ratio statistics

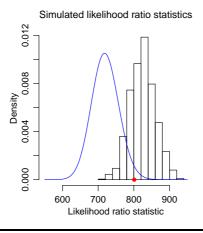


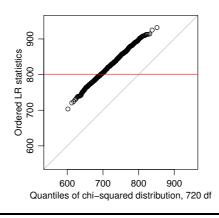
Ougeted LB statistics

Regression Methods

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

Autumn 2024 - slide 151

Premier League data: Estimates

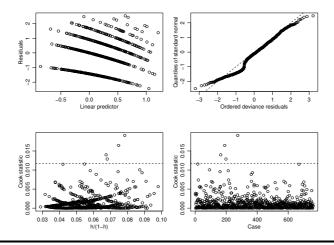
	Overall (δ)	Offensive (α)	Defensive (β)
Manchester United	0.39	0.22	0.15
Liverpool	0.13	0.12	-0.08
Arsenal		0.04	
Chelsea	-0.09	0.08	-0.22
Leeds	-0.10	0.02	-0.17
lpswich	-0.16	-0.10	-0.13
Sunderland	-0.33	-0.31	-0.10
Aston Villa	-0.48	-0.31	-0.15
West Ham	-0.53	-0.33	-0.30
Middlesborough	-0.53	-0.35	-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

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

Autumn 2024 - slide 153

2.7 Contingency Tables

slide 154

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 $0, 1, \ldots, 14$.

- \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

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

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

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 $\mu_{rc} = \exp(\gamma_r + x_{rc}^{\mathrm{T}}\beta)$:
 - 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.

☐ The resulting multinomial log likelihood is

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

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

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:
 - $\triangleright \quad \widehat{\beta}$ 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

Autumn 2024 - slide 157

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

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

so

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

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

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

Autumn 2024 - slide 158

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 * 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

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.
- \square 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

Autumn 2024 - slide 160

2.8 Ordinal Responses

slide 161

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

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

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

- 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.

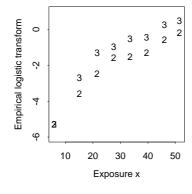
- \square Another standard tolerance distribution is $F(u) = 1 \exp\{-\exp(u)\}$.
- \square To fit, we just apply IWLS to the multinomial likelihood $\prod_{i=1}^n \pi_{I_i}$.

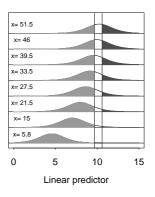
Regression Methods

Autumn 2024 - slide 163

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

Co	mments 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

Autumn 2024 - slide 165

2.9 Overdispersion

slide 166

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

UK monthly AIDS reports 1983-1992											
	Diagnosis period			Reporting-delay interval (quarters):							
Year	·		1	2	3	4	5	6		>14	reports to end of 1992
Teal	Quarter	0^{\dagger}			3	4	5	0		<u> </u>	01 1992
	:	:	:	:	:	:	:	:	:	÷	÷
1988	1	31	80	16	9	3	2	8		6	174
	2	26	99	27	9	8	11	3		3	211
	3	31	95	35	13	18	4	6		3	224
	4	36	77	20	26	11	3	8		2	205
1989	1	32	92	32	10	12	19	12		2	224
	2	15	92	14	27	22	21	12		1	219
	3	34	104	29	31	18	8	6			253
	4	38	101	34	18	9	15	6			233
1990	1	31	124	47	24	11	15	8			281
	2	32	132	36	10	9	7	6			245
	3	49	107	51	17	15	8	9			260
	4	44	153	41	16	11	6	5			285
1991	1	41	137	29	33	7	11	6			271
	2	56	124	39	14	12	7	10			263
	3	53	175	35	17	13	11	2			306
	4	63	135	24	23	12	1				258
1992	1	71	161	48	25	5					310
	2	95	178	39	6						318
	3	76	181	16							273
	4	67	66								133

Regression Methods

Autumn 2024 - slide 168

AIDS data

- □ 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:

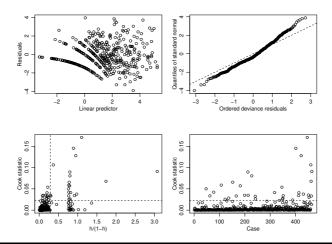
Model	df	Deviance reduction	df	Deviance
			464	14184.3
Time (rows)	37	6114.8	427	8069.5
Delay (cols)	14	7353.0	413	716.5

- \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_i .
- ☐ Residuals on next page show no obvious problems, just generic overdispersion.

Regression Methods

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).



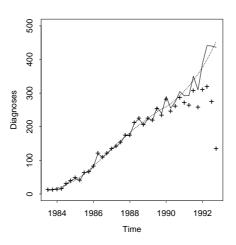
Regression Methods

Autumn 2024 - slide 170

AIDS data

 \square 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

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

Autumn 2024 - slide 172

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

$$\mathrm{E}(Y) = \mathrm{E}_{\varepsilon} \left\{ \mathrm{E}(Y \mid \varepsilon) \right\}, \quad \mathrm{var}(Y) = \mathrm{var}_{\varepsilon} \left\{ \mathrm{E}(Y \mid \varepsilon) \right\} + \mathrm{E}_{\varepsilon} \left\{ \mathrm{var}(Y \mid \varepsilon) \right\},$$

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)

- \Box MLE $\hat{\xi}_Q = 22.7 \ (5.5)$
- \Box Analysis of Deviance (with $\widehat{\xi}_Q$ 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?

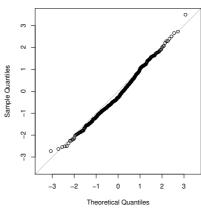
Regression Methods

Autumn 2024 - slide 173

AIDS data: Deviance residuals for NB model

Clear improvement over previous plots, even if not perfect.

Normal Q-Q Plot



Regression Methods

Quasi-likelihood

- ☐ Recall two basic assumptions for the linear model:
 - the responses are uncorrelated with means $\mu_j = x_i^{\mathrm{T}} \beta$ and equal variances σ^2 ;
 - in addition to this, the responses are normally distributed.
- $\hfill\Box$ 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 $\mathrm{E}(Y_j)=\mu_j$, so $\mathrm{E}\{h(\beta;Y)\}=0$, and under mild conditions $\tilde{\beta}$ is consistent (but maybe not efficient) as $n\to\infty$.

Regression Methods

Autumn 2024 - slide 175

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

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

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 $\hat{\phi}^{1/2}$, 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

□ 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 $B_j(\mu_j) = g'(\mu_j)\phi_j V(\mu_j)$. 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 $\text{var}(Y_j) = \phi_jV(\mu_j)$, 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 ${\cal B}$ 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)
> 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

3.1 Basic Notions slide 181

Ta	Il and wide regressions
	So far we have supposed that we have a tall regression :
	- the number of units n exceeds the number of variables p ,
	– the design matrix X has rank p .
	In many 'modern' settings we instead have a wide regression:
	- n and p are comparable, $p > n$, maybe even $p \gg n$;
	– in genomics, for example (typically) $n=O(10^2,10^3)$, $p=O(10^5,10^6)$;
	- hence $\operatorname{rank}(X) = \min(n, p) = n$.
	Even tall X may be 'almost singular', making β 'almost inestimable'.
	Solutions:
	- subset selection (drop certain columns of X);
	 seek different good explanations of response variation, not single model;
	 regularisation (often with prediction in mind).
	Certain regularisation methods (e.g., lasso) also perform subset selection.

Regression Methods Autumn 2024 – slide 182

Different good explanations \square With p > n, perhaps $p \gg n$, X is rank-deficient and many β may give $X\beta = y$. ☐ To find important variables we include intrinsic variables (gender, ...) in all models, and then choose some k (preferably ≤ 15) such that k < n and suppose that $p < k^a$ (let a = 3 for easy visualisation); – assign each variable to a cell of a hyper-cube with coordinates $\{1,\ldots,k\}^a$; - fit a linear model containing each set of k variables corresponding to the ak^{a-1} rows, columns, \dots of the cube, so each variable appears in a distinct models; - for each such model, retain the two variables that are most significant. ☐ Iterate the above procedure, retaining only the most significant variables at each stage, aiming for a final set of 10-20 variables, for which a careful analysis is performed, perhaps leading to several different good explanations of the response variation. ☐ Some cells of the hyper-cube may be empty, and important variables might be assigned to several The above design is a form of **balanced incomplete block design (BIBD)** (with k^a treatments and ak^{a-1} blocks). ☐ See Cox and Battey (2017, PNAS)

Regression Methods

Collinearity

- \square Columns of X collinear if there exists a non-zero $v_{p\times 1}$ such that Xv=0, i.e., $\mathrm{rank}(X)< p$, so there is no unique $\widehat{\beta}$ minimising $\|y-X\beta\|^2$.
- \square Software deals with this by dropping columns of X, but it may be better to write $X\beta = XC\gamma$, where XC is full rank and γ has a clear interpretation.
- \square If X is nearly collinear, its SVD $U_{n\times n}D_{n\times p}V_{p\times p}^{\mathrm{T}}$, with $d_1\geq \cdots \geq d_p\geq 0$, gives

$$\widehat{\beta} = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y = VD_{-}^{\mathrm{T}}U^{\mathrm{T}}y = \sum_{r=1}^{p} (u_r^{\mathrm{T}}y/d_r)v_r,$$

so $\widehat{\beta}$ is a linear combination of the vectors v_r with coefficients $u_r^{\mathrm{\scriptscriptstyle T}} y/d_r$. As $\mathrm{var}(U^{\mathrm{\scriptscriptstyle T}} y) = \sigma^2 I_n$,

$$\operatorname{var}(\widehat{\beta}) = \sigma^2 V D_{-}^{\mathrm{\scriptscriptstyle T}} D_{-} V^{\mathrm{\scriptscriptstyle T}} = \sigma^2 \sum_{r=1}^p d_r^{-2} v_r v_r^{\mathrm{\scriptscriptstyle T}},$$

i.e., $\widehat{\beta}$ is unstable in the directions corresponding to the v_r with small singular values d_r .

 \square In numerical analysis, collinearity often measured using **condition number** $(d_1/d_p)^{1/2}$, but its statistical meaning is unclear.

Regression Methods

Autumn 2024 - slide 184

Regularisation

Stop $\widehat{\beta}$ from fluctuating too wildly in directions with small eigenvalues d_r , by adding a non-negative penalty $p_{\lambda}(\beta)$ and choosing β to minimise the **penalised sum of squares**

$$||y - X\beta||^2 + p_{\lambda}(\beta). \tag{16}$$

- \square The strength of the penalty depends on a positive parameter λ that constrains β more as λ increases.
- \square Often $p_{\lambda}(\beta) = \lambda p(\beta)$, where, for example,
 - $p(\beta) = \|\beta\|_2^2 = \sum_{r=1}^p \beta_r^2$ gives ridge regression (aka Tikhonov regularisation);
 - $p(\beta) = \|\beta\|_1 = \sum_{r=1}^p |\beta_r|$ gives the lasso (aka L_1 regularisation);
 - $p(\beta) = (1 \alpha) \|\beta\|_2^2 + \alpha \|\beta\|_1$ for $0 \le \alpha \le 1$ gives the elastic net;
 - $p(\beta) = \sum_{g=1}^G p_g^{1/2} \|\beta_g\|_2$, with β_g being $p_g \times 1$ sub-vectors of β , gives the **grouped lasso**, which penalises factors with parameters β_g .
- □ It is useful to see regularisation through the lens of Bayesian inference, with the regularising term equivalent to the prior density.

Regression Methods

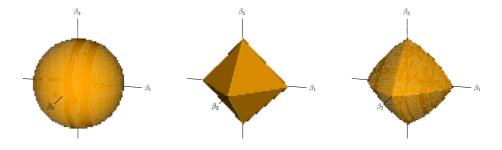
Bound form

☐ Equivalently we can take the **bound form** of the minimisation problem, i.e.,

minimise_{$$\beta$$} $||y - X\beta||_2^2$ subject to $p(\beta) \le t$,

for some $t \ge 0$, where setting $t = \infty$ just gives the least squares estimates.

Below: constraint balls for ridge (left), lasso (centre) and elastic-net (right) regularisation. The sharp corners of the last two allow for variable selection as well as shrinkage.



Regression Methods

Autumn 2024 - slide 186

Bayesian setting

- ☐ Treat all unknowns as random variables, and compute conditional distribution of unobserved unknowns conditional on observed unknowns.
- \square Requires prior density on β , and if σ^2 is known, then a simple combination of **data model** and **prior model** is

$$y \mid \beta, \sigma^2 \sim \mathcal{N}(X\beta, \sigma^2 I_n), \quad \beta \mid \sigma^2 \sim \mathcal{N}_p(\beta_*, \sigma^2 V_*),$$
 (17)

where the prior model is determined by β_* and V_* .

- \Box Full specification would require prior on σ^2 , but we don't need this.
- \Box Let \equiv mean we have dropped additive constants not involving the argument of a density.
- ☐ The log multivariate normal density is

$$\log f(x \mid \mu, \Omega) = -\frac{m}{2} \log 2\pi - \frac{1}{2} \log |\Omega| - \frac{1}{2} (x - \mu)^{\mathrm{T}} \Omega^{-1} (x - \mu)$$

$$\equiv x^{\mathrm{T}} \Omega^{-1} \mu - \frac{1}{2} x^{\mathrm{T}} \Omega^{-1} x$$

$$\equiv Q(x) = x^{\mathrm{T}} a - \frac{1}{2} x^{\mathrm{T}} B x,$$

say, and as $\exp Q(x)$ is proportional to a unique probability density function,

$$\mathrm{E}(X) = \mu = B^{-1}a, \quad \mathrm{var}(X) = \Omega = B^{-1}, \quad \text{where } B \text{ is the precision matrix}.$$

Regression Methods

Bayesian linear model I

☐ The model (6) gives

$$\log f(\beta \mid y, \sigma^{2}) = \log \left\{ \frac{f(y \mid \beta, \sigma^{2}) f(\beta \mid \sigma^{2})}{f(y \mid \sigma^{2})} \right\}$$

$$\equiv \log f(y \mid \beta, \sigma^{2}) + \log f(\beta \mid \sigma^{2})$$

$$\equiv -\frac{(y - X\beta)^{\mathrm{T}} (y - X\beta)}{2\sigma^{2}} - \frac{(\beta - \beta_{*})^{\mathrm{T}} V_{*}^{-1} (\beta - \beta_{*})}{2\sigma^{2}}$$

$$\propto \|y - X\beta\|_{2}^{2} + (\beta - \beta_{*})^{\mathrm{T}} V_{*}^{-1} (\beta - \beta_{*}).$$

- Comparison with (5) shows that $p_{\lambda}(\beta)$ represents prior beliefs about the likely values of β : before seeing the data, the most plausible value is β_* , with precision V_*^{-1} .
- ☐ Dropping more constants,

$$\log f(\beta \mid y, \sigma^{2}) \equiv \frac{1}{\sigma^{2}} \left\{ \beta^{\mathrm{T}} X^{\mathrm{T}} y - \beta^{\mathrm{T}} (X^{\mathrm{T}} X) \beta / 2 + \beta^{\mathrm{T}} V_{*}^{-1} \beta_{*} - \beta^{\mathrm{T}} V_{*}^{-1} \beta / 2 \right\}$$

$$= \frac{1}{2\sigma^{2}} \left\{ 2\beta^{\mathrm{T}} (X^{\mathrm{T}} y + V_{*}^{-1} \beta_{*}) - \beta^{\mathrm{T}} (X^{\mathrm{T}} X + V_{*}^{-1}) \beta \right\},$$
(18)

which is Q(x) with x, a and B replaced by β , $(X^{\mathrm{T}}y+V_*^{-1}\beta_*)/\sigma^2$ and $(X^{\mathrm{T}}X+V_*^{-1})/\sigma^2$.

 \square Hence $f(\beta \mid y, \sigma^2)$ is multivariate normal with mean vector and variance matrix

$$E(\beta \mid y, \sigma^2) = (X^{\mathrm{T}}X + V_*^{-1})^{-1}(X^{\mathrm{T}}y + V_*^{-1}\beta_*), \quad \text{var}(\beta \mid y, \sigma^2) = \sigma^2(X^{\mathrm{T}}X + V_*^{-1})^{-1}.$$

Regression Methods

Autumn 2024 - slide 188

Bayesian linear model II

- The maximum a posteriori (MAP) estimator of β is $E(\beta \mid y, \sigma^2)$, and the MAP estimator of $A_{q \times p} \beta$ is $AE(\beta \mid y, \sigma^2)$, which has a posterior normal density.
- \square When $X^{\mathrm{T}}X$ is invertible,

$$\tilde{\beta} = \mathcal{E}(\beta \mid y, \sigma^2) = (X^{\mathrm{T}}X + V_*^{-1})^{-1}(X^{\mathrm{T}}X\hat{\beta} + V_*^{-1}\beta_*)$$

is an average of $\widehat{\beta}$ and β_* , weighted by $X^{\mathrm{T}}X$ and V_*^{-1} .

☐ The posterior precision matrix

$$var(\beta \mid y, \sigma^2)^{-1} = X^{\mathrm{T}} X / \sigma^2 + V_*^{-1} / \sigma^2$$

adds the Fisher information and the prior precision matrix, $V_{*}^{-1}/\sigma^{2}.$

- $\hfill\Box$ High precision corresponds to small variance, and conversely:
 - letting $V_*^{-1} \to 0$ yields an improper prior density; and
 - for large V_{\ast}^{-1} the posterior precision is essentially determined by the prior precision.

Thus the prior density regularises $\widehat{\beta}$ by including β_* and V_* .

Regression Methods

Improper prior density

 \square We only need V_* to add information in directions corresponding to small singular values of X, so we might use an **improper prior** in which V_* is singular:

$$f(\beta \mid \sigma^2) = \frac{1}{(2\pi)^{p/2} |V_*|_+^{1/2}} \exp\left\{-(\beta - \beta_*)^{\mathrm{T}} V_*^{-} (\beta - \beta_*) / (2\sigma^2)\right\},\tag{19}$$

where V_* has spectral decomposition ED_*E^{T} ,

- $|V_*|_+$ denotes the product of the non-zero elements of D_* , and
- $V_*^- = \sum_{r:d_{*r}>0} e_r e_r^{\mathrm{T}}/d_{*r}$ is a generalized inverse of V_* .
- \square Below we write V_*^- even when V_* is invertible.
- \square (8) is improper because it is not integrable in the directions of the columns of E for which the corresponding d_r^* equal zero, but we need only that the posterior density of β be proper, i.e., that the posterior precision matrix

$$var(\beta \mid y, \sigma^2)^{-1} = X^{T}X/\sigma^2 + V_*^{-}/\sigma^2$$

is invertible.

Regression Methods

Autumn 2024 - slide 190

Empirical Bayes

- \square Use the data to estimate the prior: construct estimators using Bayesian arguments, but assess their properties using classical criteria (bias, MSE, ...)
- \square The estimator $\tilde{\beta} = \mathrm{E}(\beta \mid y, \sigma^2)$ has mean and variance

$$E(\tilde{\beta} \mid \beta) = (X^{T}X + V_{*}^{-})^{-1}(X^{T}X\beta + V_{*}^{-}\beta_{*})$$

$$= \beta + (X^{T}X + V_{*}^{-})^{-1}V_{*}^{-}(\beta_{*} - \beta),$$

$$var(\tilde{\beta} \mid \beta) = \sigma^{2}(X^{T}X + V_{*}^{-})^{-1}X^{T}X(X^{T}X + V_{*}^{-})^{-1}.$$
(20)

- \square Hence $\tilde{\beta}$
 - is biased unless $\beta_* = \beta$,
 - has smaller variance than $\widehat{\beta}$,

so maybe there is a bias-variance tradeoff when estimating $A\beta$.

 \square If we write $\mu = E(\tilde{\beta} \mid \beta)$, then the MSE is

$$\begin{split} \mathbf{E} \left(\| A \tilde{\beta} - A \beta \|^2 \mid \beta \right) &= \mathbf{E} \{ (\tilde{\beta} - \beta)^{\mathrm{T}} A^{\mathrm{T}} A (\tilde{\beta} - \beta) \mid \beta \} \\ &= \mathbf{E} \left[\mathrm{tr} \left\{ A (\tilde{\beta} - \beta) (\tilde{\beta} - \beta)^{\mathrm{T}} A^{\mathrm{T}} \right\} \mid \beta \right] \\ &= \mathbf{tr} \left[\mathbf{E} \left\{ A (\tilde{\beta} - \mu + \mu - \beta) (\tilde{\beta} - \mu + \mu - \beta)^{\mathrm{T}} A^{\mathrm{T}} \mid \beta \right\} \right]. \end{split}$$

Regression Methods

Empirical Bayes II

 \square The expectation above is

$$A\left\{ \mathrm{var}(\tilde{\beta} \mid \beta) + (X^{\mathrm{\scriptscriptstyle T}}X + V_*^-)^{-1}V_*^-(\beta - \beta_*)(\beta - \beta_*)^{\mathrm{\scriptscriptstyle T}}V_*^-(X^{\mathrm{\scriptscriptstyle T}}X + V_*^-)^{-1} \right\} A^{\mathrm{\scriptscriptstyle T}},$$

giving the MSE when estimating a fixed β .

 \square Taking expectations over the prior model for β gives

$$\mathrm{E}\left(\|A\tilde{\beta} - A\beta\|^{2}\right) = \sigma^{2} \mathrm{tr}\left\{A(X^{\mathrm{T}}X + V_{*}^{-})^{-1}A^{\mathrm{T}}\right\},\tag{21}$$

which is larger than $A \text{var}(\tilde{\beta} \mid \beta) A^{\text{T}}$ and does not depend on β_* .

- ☐ This computation uses only the mean and variance, so holds under second-order assumptions.
- \square From now on we set $\beta_* = 0$, unless we state otherwise.

Regression Methods

Autumn 2024 - slide 192

Equivalent degrees of freedom

 \Box If we set $\beta_* = 0$, then the fitted values are

$$\tilde{y} = X\tilde{\beta} = X(X^{\mathrm{T}}X + V_{*}^{-})^{-1}X^{\mathrm{T}}y = H_{*}y,$$

say.

☐ We define the **equivalent degrees of freedom** of the fit as

$$edf = tr(H_*) = tr\{X(X^{\mathrm{T}}X + V_*^{-})^{-1}X^{\mathrm{T}}\} = p - tr\{(X^{\mathrm{T}}X + V_*^{-})^{-1}V_*^{-}\},\$$

- \square This is lower than p unless $V_*^-=0$, so regularisation reduces the degrees of freedom by an amount that depends on V_* .
- ☐ The penalised estimate is a linear function of the unpenalised one (if it exists), as we can write

$$\tilde{\beta} = (X^{\mathrm{T}}X + V_*^-)^{-1}X^{\mathrm{T}}X\widehat{\beta} = P_*\widehat{\beta},$$

say. As

$$edf = tr(H_*) = tr(P_*),$$

this gives an alternative formula useful in complex models.

Regression Methods

How much penalisation?

- \square Often V_*^- depends on some $\lambda > 0$ that must be chosen, as well as σ^2 , which is usually estimated by a (penalised) residual sum of squares.
- \square To estimate λ , we compare y_j with its predicted value $\widehat{y}_{\lambda,j} = x_j^{\mathrm{T}} \widehat{\beta}_{\lambda,-j}$, where $\widehat{\beta}_{\lambda,-j}$ is

$$\widehat{\beta}_{\lambda} = (X^{\mathrm{T}}X + V_{*}^{-})^{-1}X^{\mathrm{T}}y$$

computed with the jth rows x_i and y_i of X and y omitted.

☐ Using Lemma 14, the leave-one-out cross-validation sum of squares is then

$$CV_{\lambda} = \sum_{j=1}^{n} (y_j - \widehat{y}_{\lambda,j})^2 = \|y - \widehat{y}_{\lambda}\|^2 = \sum_{j=1}^{n} \frac{(y_j - \widehat{y}_{\lambda,j})^2}{(1 - h_{\lambda,jj})^2},$$

where $\widehat{y}_{\lambda,j}$ is the jth element of the complete-data fitted value $H_{\lambda}y$ and $h_{\lambda,jj}$ is the jth diagonal element of $H_{\lambda}=X(X^{\mathrm{T}}X+V_{*}^{-})^{-1}X^{\mathrm{T}}$ for the overall fit.

☐ More often we use the **generalized cross-validation** criterion

$$GCV_{\lambda} = \sum_{j=1}^{n} \frac{(y_j - \widehat{y}_{\lambda,j})^2}{\{1 - \operatorname{tr}(H_{\lambda})/n\}^2}.$$

 \square Whichever criterion is used, it is typically minimised numerically over a grid of values of λ .

Regression Methods

Autumn 2024 - slide 194

REML

- ☐ Cross-validation makes only second-order assumptions.
- Under normality, the marginal density of y is $\mathcal{N}\{X\beta_*, \sigma^2(I_n + XV_*X^{\mathrm{T}})\}$, so we could estimate β_* , σ^2 and λ by maximising the corresponding likelihood.
- □ If n and p are large, this results in biased estimates of λ and σ^2 , so we prefer to eliminate β_* , resulting in a \log restricted likelihood whose form is given below, with $W_{\lambda}^{-1} = I_n + XV_*X^{\mathrm{T}}$.

Lemma 31 In a model in which $y \sim \mathcal{N}(X\beta, \sigma^2 W_{\lambda}^{-1})$, where W_{λ} depends on a parameter λ , a log restricted likelihood for σ^2 and λ is

$$\ell_{\text{REML}}(\sigma^2, \lambda) \equiv \frac{1}{2} \log(|W_{\lambda}|/|X^{\mathsf{T}}W_{\lambda}X|) - \frac{n-p}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y - \widehat{y}_{\lambda})^{\mathsf{T}}W_{\lambda}(y - \widehat{y}_{\lambda}),$$

where $\widehat{\beta}_{\lambda}=(X^{\mathrm{T}}W_{\lambda}X)^{-1}X^{\mathrm{T}}W_{\lambda}y$ and $\widehat{y}_{\lambda}=X\widehat{\beta}_{\lambda}$. For fixed λ the restricted maximum likelihood estimator of σ^2 is therefore

$$\widehat{\sigma}_{\lambda}^{2} = \frac{1}{n-p} (y - \widehat{y}_{\lambda})^{\mathrm{T}} W_{\lambda} (y - \widehat{y}_{\lambda}),$$

and the resulting profile log restricted likelihood for λ is

$$\ell_{\mathbf{p}}(\lambda) \equiv \frac{1}{2} \log(|W_{\lambda}|/|X^{\mathrm{T}}W_{\lambda}X|) - \frac{(n-p)}{2} \log \widehat{\sigma}_{\lambda}^{2}.$$

Regression Methods

Note on Lemma 31

 \square Suppose that $f(y; \alpha, \beta)$ depends on two parameters, that interest is focused on α , and that for fixed α there is a minimal sufficient statistic s_{α} for β . Then $f(y; \alpha, \beta) = f(y \mid s_{\alpha}; \alpha) f(s_{\alpha}; \alpha, \beta)$, and since the first density on the right is a proper conditional density not depending on β , we can use it for inference on α , in the form

$$\log f(y \mid s_{\alpha}; \alpha) = \log f(y; \alpha, \beta) - \log f(s_{\alpha}; \alpha, \beta).$$

As the left-hand side of this expression does not depend on β , we may be able to simplify the right-hand side by an astute choice of β .

 $\Box \quad \text{In the normal model we take } \alpha = (\sigma^2, \lambda). \text{ If } \alpha \text{ is fixed, then } s_\alpha = \widehat{\beta}_\alpha = (X^{\mathrm{T}} W_\lambda X)^{-1} X^{\mathrm{T}} W_\lambda y \text{ is sufficient for } \beta; \text{ its distribution is } \mathcal{N}_p \{\beta, \sigma^2 (X^{\mathrm{T}} W_\lambda X)^{-1}\}. \text{ Hence }$

$$\ell_{\text{REMI}}(\sigma^2, \lambda) = \log f(y \mid \widehat{\beta}_{\lambda}; \sigma^2, \lambda) = \log f(y; \sigma^2, \lambda, \beta) - \log f(\widehat{\beta}_{\lambda}; \sigma^2, \lambda, \beta)$$

which equals

$$-\frac{n}{2}\log\sigma^{2} + \frac{1}{2}\log|W_{\lambda}| - \frac{1}{2\sigma^{2}}(y - X\beta)^{\mathrm{T}}W_{\lambda}(y - X\beta) + \frac{p}{2}\log\sigma^{2} - \frac{1}{2}\log|X^{\mathrm{T}}W_{\lambda}X| + \frac{1}{2\sigma^{2}}(\widehat{\beta}_{\lambda} - \beta)^{\mathrm{T}}X^{\mathrm{T}}W_{\lambda}X(\widehat{\beta}_{\lambda} - \beta),$$

or equivalently, on setting $\beta=0$ and $\widehat{y}_{\lambda}=X\widehat{\beta}_{\lambda}$

$$\frac{1}{2}\log(|W_{\lambda}|/|X^{\mathrm{\scriptscriptstyle T}}W_{\lambda}X|) - \frac{(n-p)}{2}\log\sigma^2 - \frac{1}{2\sigma^2}\left(y^{\mathrm{\scriptscriptstyle T}}W_{\lambda}y - \widehat{y}_{\lambda}^{\mathrm{\scriptscriptstyle T}}X^{\mathrm{\scriptscriptstyle T}}W_{\lambda}\widehat{y}_{\lambda}\right).$$

- The last term reduces to the given form because $\widehat{y}_{\lambda}^{\mathrm{T}}W_{\lambda}(y-\widehat{y}_{\lambda})=0$, so the term in brackets in the last displayed equation is the residual sum of squares $(y-\widehat{y}_{\lambda})^{\mathrm{T}}W_{\lambda}(y-\widehat{y}_{\lambda})$.
- □ The restricted maximum likelihood estimator $\widehat{\sigma}_{\lambda}^2$ and the profile log restricted likelihood for λ are obtained by maximising $\ell_{\text{REML}}(\sigma^2, \lambda)$, for fixed λ and then dropping constant terms from $\ell_{\text{REML}}(\widehat{\sigma}_{\lambda}^2, \lambda)$.

Regression Methods

Autumn 2024 - note 1 of slide 195

Numerical example from Wood (2011, JRSSB)

The usual methods (AIC, GCV, ...) for choosing λ are available, but we focus on likelihood methods; see below.

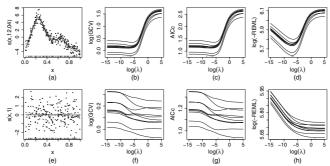


Fig. 1. Example comparison of GCV, AICc and REML criteria: (a) some (x,y)-data modelled as $y_i = f(x_i) + \varepsilon_i$, ε_i independent and identically distributed $N(0, \sigma^2)$ where smooth function f was represented by using a rank 20 thin plate regression spline (Wood, 2003); (b)–(d) various smoothness selection criteria plotted against logarithmic smoothing parameters, for 10 replicates of the data (each generated from the same truth') (note how shallow the GCV and AICc minima are relative to the sampling variability, resulting in rather variable optimal λ -values (which are shown as a rug plot), and a propensity to undersmooth; in contrast the REML optima are much better defined, relative to the sampling variability, resulting in a smaller range of λ -estimates); (e)–(f) are equivalent to (a)–(d), but for data with no signal, so that the appropriate smoothing parameter should tend to ∞ (note GCV's and AICc's occasional multiple minima and undersmoothing in this case, compared with the excellent behaviour of REML; ML (which is not shown) has a similar shape to REML)

Regression Methods

Autumn 2024 - slide 196

3.2 Simple Applications

slide 197

Ridge regression

- \square Used for prediction when X is close to singular.
- \square If the first column of X is 1_n , we set $\beta_*=0$ and $V_*^-=\lambda S=\lambda {
 m diag}(0,I_{p-1})$, giving

$$\widehat{\beta}_{\lambda} = (X^{\mathrm{T}} + \lambda S)^{-1} X^{\mathrm{T}} y, \quad \widehat{y}_{\lambda} = X \widehat{\beta}_{\lambda} = X (X^{\mathrm{T}} + \lambda S)^{-1} X^{\mathrm{T}} y = H_{\lambda} y,$$

and effective degrees of freedom

$$\operatorname{edf}_{\lambda} = \operatorname{tr}(H_{\lambda}) = \operatorname{tr}\{(X^{\mathsf{T}}X + \lambda S)^{-1}X^{\mathsf{T}}X\} = \sum_{r=1}^{p} \frac{1}{1 + \lambda \delta_{r}},$$

where $\delta_p \geq \cdots \geq \delta_2 > \delta_1 = 0$ are the eigenvalues of $(X^{\mathrm{T}}X)^{-1/2}S(X^{\mathrm{T}}X)^{-1/2}$.

- \square As λ increases from zero to infinity, $\operatorname{edf}_{\lambda}$ decreases from $p=\operatorname{rank}(X)$ to 1. The two are equivalent, but $\operatorname{edf}_{\lambda}$ is more easily interpreted, because it is not related to the scale of X.
- \square The inverse exists even if X^TX is singular, but if it is invertible then

$$\widehat{\beta}_{\lambda} = (X^{\mathrm{T}}X + \lambda S)^{-1}(X^{\mathrm{T}}X + \lambda S - \lambda S)(X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y = \widehat{\beta} - \lambda(X^{\mathrm{T}}X + \lambda S)^{-1}S\widehat{\beta},$$

so as $\lambda \to \infty$ all the elements of $\widehat{\beta}_{\lambda}$ tend to zero, other than the first. This corresponds to reducing the prior variance to zero, thereby giving the data themselves less and less influence on the elements of $\widehat{\beta}_{\lambda}$ other than the first, and thus stabilises the estimator.

Regression Methods

Example: Cement data > cement x1 x2 x3 x4 7 26 6 60 78.5 1 29 15 52 74.3 11 56 8 20 104.3 8 47 11 31 87.6 7 52 6 33 95.9 11 55 9 22 109.2 3 71 17 6 102.7 1 31 22 44 72.5 2 54 18 22 93.1 4 26 115.9 10 21 47 11 1 40 23 34 83.8 12 11 66 9 12 113.3 13 10 68 8 12 109.4

Regression Methods

Autumn 2024 - slide 199

Example: Cement data

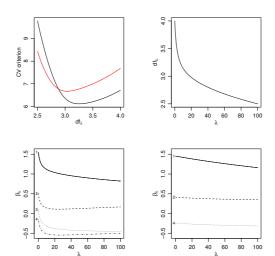
	Fu	II model	Redu	Reduced model			
Parameter	Estimate	Standard error	Estimate	Standard error			
eta_0	62.41	70.07	71.64	14.14			
eta_1	1.55	0.74	1.45	0.12			
eta_2	0.51	0.72	0.42	0.19			
eta_3	0.10	0.75					
eta_4	-0.14	0.71	-0.24	0.17			

- \square The next slide shows results for ridge fits for these models.
- \square Looks like 3 df is optimal for prediction.
- \square Software often preprocesses X and y by either
 - centering both, by subtracting column means, or
 - centering y and centering and scaling X, so the column means are zero and the column variances are unity.
- \square The singular values for the centred X matrix are 78.8, 28.5, 12.2, 1.7, and those for the centred and scaled X matrix are 5.18, 4.35, 1.50, 0.14, so it matters which is used.
- \Box The singular values for the (centred) reduced matrix are 78.8, 19.8 and 9.15.
- \square The shrinkage due to increasing λ occurs more slowly for the reduced model.

Regression Methods

Example: Cement data/Ridge analysis

Top left: CV (black) and GCV (red) as functions of degrees of freedom df_{λ} . Top right: dependence of df_{λ} on λ . Bottom left: $\widehat{\beta}_{\lambda}$ as a function of λ , with all four covariates. Bottom right: $\widehat{\beta}_{\lambda}$ as a function of λ , with x_1 , x_2 , and x_4 only.



Regression Methods Autumn 2024 – slide 201

Comments

☐ The literature on ridge regression is very large and very dispersed, with many variants and many connections to ML techniques.

 \square Be careful with software: any pre-processing of X is not always described.

Regression Methods

Autumn 2024 - slide 202

Semiparametric regression

□ Normal linear model has two main aspects:

- systematic variation, $E(y) = \mu$, and $\mu = X\beta$ with parameters β ;
- stochastic variation, $y \sim \mathcal{N}_n(\mu, \sigma^2 I_n)$.

☐ Can relax the stochastic assumption using other distributions or second-order assumptions, but still have parametric model for the systematic part.

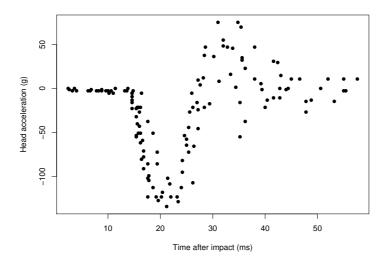
☐ Often want to relax systematic part for more flexible models, for

- exploratory data analysis 'will a linear model be adequate?'
- confirmatory data analysis 'I've fitted a linear model, is it adequate?'
- general modelling 'the data are too complex to expect a simple parametric model to work, so what can I do?'
- semiparametric modelling 'I will use a parametric model for the effects of interest, but can I model nuisance effects more flexibly?'
- ☐ Most basic tool is the **scatterplot smoother**.

Regression Methods

Example: Motorcycle data

Measurements of head acceleration (g) at time after impact (ms) in a simulated motorcycle accident, used to test crash helmets:



Regression Methods

Autumn 2024 - slide 204

Scatterplot smoothing

- \square Have data $(x_1, y_1), \ldots, (x_n, y_n)$, with $x_- \le x_1 < \cdots < x_n \le x_+$ (ahem) and we wish to estimate $\mathrm{E}(y) = \mu(x)$, for $x \in \mathcal{X} = [x_-, x_+]$.
- Suppose that $\mu \in \mathcal{M}$, a function space spanned by n linearly independent basis functions that can be identified by evaluation at x_1, \ldots, x_n , and let $\mu_j = \mu(x_j)$.
- \square Can choose a basis $\{b_1(x), \ldots, b_n(x)\}$ for \mathcal{M} such that $\mu(x) = \sum_{j=1}^n \mu_j b_j(x)$ interpolates $(x_1, \mu_1), \ldots, (x_n, \mu_n)$.
- Suppose that \mathcal{M} contains the linear functions on \mathcal{X} and that the second derivatives of the $b_j(x)$ are not all zero, so functions in \mathcal{M} may also be nonlinear in x.
- \square To estimate μ we minimise a **penalised sum of squares**,

$$\sum_{j=1}^{n} \{y_j - \mu(x_j)\}^2 + \lambda \int_{\mathcal{X}} \{\mu''(x)\}^2 \, \mathrm{d}x,\tag{22}$$

where the **roughness penalty** imposes smoothness: if $\lambda \to 0$, then $\mu(x_j) \to y_j$ and $\widehat{\mu}$ interpolates, but when $\lambda \to \infty$ even tiny wiggles in μ will give a huge penalty, making $\widehat{\mu}$ linear.

 \square The penalty does not affect linear functions, so $\mathcal{M} = \mathcal{L} \bigoplus \mathcal{P}$, where \mathcal{L} and \mathcal{P} are the two-dimensional vector space of linear functions on \mathcal{X} and an (n-2)-dimensional vector space of nonlinear functions on \mathcal{X} , and \bigoplus denotes addition of vector spaces.

Regression Methods

Scatterplot smoothing II

☐ The roughness term is

$$\int_{\mathcal{X}} \{\mu''(x)\}^2 dx = \int_{\mathcal{X}} \left\{ \sum_{j=1}^n \mu_j b_j''(x) \right\}^2 dx = \sum_{i,j=1}^n \mu_i \mu_j \int_{\mathcal{X}} b_i''(x) b_j''(x) dx = \mu^{\mathrm{T}} S \mu,$$

say, where $\mu^{\mathrm{T}} = (\mu_1, \dots, \mu_n)$.

- \square $S_{n\times n}$ has (i,j) element $\int_{\mathcal{X}} b_i''(x)b_j''(x)\,\mathrm{d}x$ and is symmetric and positive semi-definite of rank n-2, because linear functions are unpenalised, so $S1_n=S(x_1,\ldots,x_n)^\mathrm{T}=0$.
- ☐ The penalised sum of squares

$$(y - \mu)^{\mathrm{T}}(y - \mu) + \lambda \mu^{\mathrm{T}} S \mu \equiv -2\mu^{\mathrm{T}} y + \mu^{\mathrm{T}} (I_n + \lambda S) \mu,$$

is minimised by $\widehat{\mu}_{\lambda} = (I_n + \lambda S)^{-1} y$.

- \square As λ increases from zero, the fitted value $\widehat{\mu}_{\lambda}$ shrinks from y towards the straight-line regression fit to y, which is unpenalised.
- The equivalent degrees of freedom are $\mathrm{edf}_{\lambda} = \mathrm{tr}(H_{\lambda}) = \sum_{j=1}^{n} (1+\lambda \delta_{j})^{-1}$, where $\delta_{1} \geq \cdots \geq \delta_{3} > \delta_{2} = \delta_{1} = 0$ are the eigenvalues of S. As λ increases edf_{λ} decreases monotonically from $\mathrm{edf}_{0} = n$ towards $\mathrm{edf}_{\infty} = 2$.

Regression Methods

Autumn 2024 - slide 206

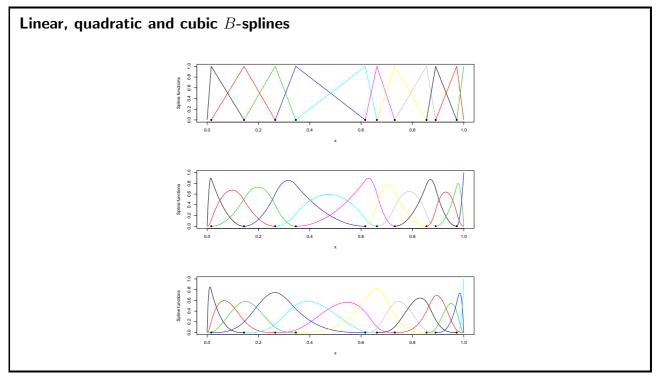
Scatterplot smoothing III

- ☐ In principle we might take any basis functions, but in practice we usually take local polynomials known as **splines** that have good approximation properties.
- ☐ There are many forms of splines, which
 - are often cubic polynomials with finite support between values of x known as **knots**, x_1^*, \ldots, x_K^* , and then S is tri-diagonal,
 - sometimes form a **natural cubic spline**, which has K = n and certain optimality properties,
 - are discussed in more detail later.
- If there is no penalisation ($\lambda = 0$) then we have a standard linear model, and spline basis functions are called **regression splines**.
- \square Under second-order assumptions we choose λ by minimising $CV(\lambda)$ or $GCV(\lambda)$.
- \Box Under normal-theory assumptions we can use REML to estimate σ^2 and λ .
- \square Obvious generalisation allows weight matrix $W = \operatorname{diag}(w_1, \dots, w_n)$.
- If the x_1, \ldots, x_n are not unique, write $E(y) = N_{n \times n'} \mu_{n' \times 1}$ in terms of the means μ at the n' unique elements of x, and minimise

$$(y - N\mu)^{\mathrm{T}}W(y - N\mu) + \lambda\mu^{\mathrm{T}}S\mu.$$

where $S_{n'\times n'}$ arises as before from the roughness penalty on $\mu(x)$.

Regression Methods

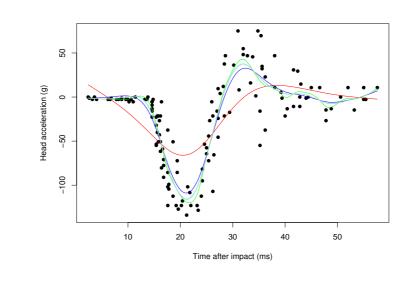


Regression Methods

Autumn 2024 - slide 208

Example: Motorcycle data

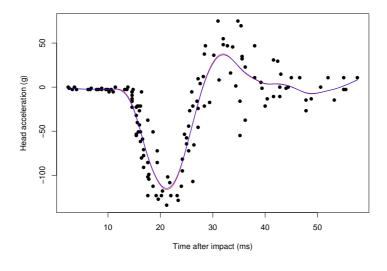
Scatterplot smooths based on natural cubic splines with edf equal to 5 (red), 10 (blue), 20 (green), and chosen by CV (cyan, edf = 12.8) and GCV (pink, edf = 12.26):



Regression Methods

Example: Motorcycle data

Scatterplot smooths based on natural cubic splines with weights 16 when $x \le 12$ and 1 for x > 12, and edf chosen by CV (red, edf = 14.7) and GCV (blue, edf = 13.7):



Regression Methods

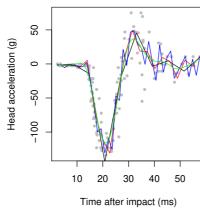
Autumn 2024 - slide 210

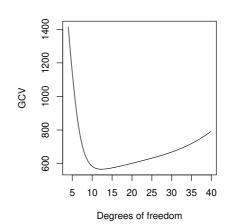
Choosing K and λ

- Above we took K=n basis functions, but for statistical purposes we seek a summary of the data, so we hope that $\mathrm{edf} \ll n$, so we hope that K < n, maybe even $K \ll n$.
- \square Theory suggests that as $n \to \infty$ we need $K = O(n^{1/5})$ or even $O(n^{1/9})$ to get near-optimal estimation of $\mu(x)$, when μ lies in reasonable function classes;
- In practice we take K (more than) large enough to give enough flexibility (increasing it if results are suspect, K = 9 by default in mgcv), and allow λ to determine the smoothness of the curve;
- \square Typically the knots x_k^* are placed at equally-spaced quantiles of x.

Regression Methods

Example: Motorcycle data





Left: linear spline fits with $\lambda=0$ and K=10 (black), 20 (red), 40 (blue), and optimal GCV choice of λ with K=40 (green)

Right: $GCV(\lambda)$ as a function of df_{λ} for K=40.

Regression Methods

Autumn 2024 - slide 212

Comments

 \square We discuss inference (beyond 'point' estimation) and adaptive estimation of weights later ...

☐ Here we are producing point estimates; later we discuss the construction of confidence sets.

An alternative local averaging approach uses locally weighted fits, such as the Nadaraya-Watson estimator

$$\widehat{\mu}(x) = \frac{\sum_{j=1}^{n} K\{(x - x_j)/h\} y_j}{\sum_{j=1}^{n} K\{(x - x_j)/h\}},$$

where

– the kernel function K is something like the Gaussian density, and

- the **bandwidth** h plays a role similar to edf.

This is also a linear smoother, and in fact the spline smoothers have representations in terms of equivalent kernels.

☐ Local averaging can be extended to **local likelihood** fitting of more complex models.

Regression Methods

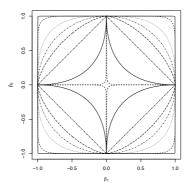
3.3 Lasso slide 214

L_q penalties

 \square The quadratic penalty $\|eta\|_2$ generalises to other L_q penalties

$$\|\beta\|_q = \sum_{r=1}^p |\beta_r|^q,$$

shown below for p=2 and (working inwards) $q=100,\ 10,\ 3,\ 2,\ 1.5,\ 1,\ 0.5,\ 0.2;$ $\|\beta\|_0=\#\{\beta_r\neq 0\}$ counts the number of non-zero parameters.



(Some picture credits here and later: Simon Wood)

Regression Methods Autumn 2024 – slide 215

Basic geometry

 \square If $D(\beta)$ is a sum of squares or negative log likelihood, then

$$\tilde{\beta}_{\lambda} = \operatorname{argmin}_{\beta} \left\{ D(\beta) + \lambda \|\beta\|_{q} \right\},\,$$

- satisfies $\| ilde{eta}_{\lambda}\|_q=t$ for some t, and
- minimises $D(\beta)$ on that contour, i.e.,

$$\tilde{\beta}_{\lambda} = \mathrm{argmin}_{\beta} D(\beta) \quad \text{such that} \quad \|\tilde{\beta}_{\lambda}\|_q = t,$$

because otherwise we could reduce $D(\beta)$ while leaving the penalty unchanged, i.e., $\tilde{\beta}_{\lambda}$ would not be optimal.

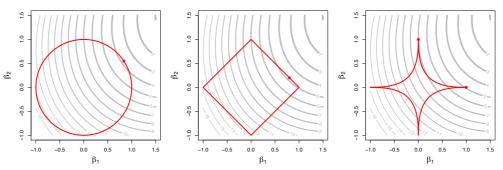
- \square The sets $\|\tilde{\beta}_{\lambda}\|_q = t$
 - have cusps (and thus can set $\beta_r=0$) when $q\leq 1$,
 - are non-convex (and thus may give non-unique solutions) when q < 1,

so there is a unique solution if the contours of $D(\beta)$ and $\|\beta\|_q$ are convex, and both a unique solution and the possibility of choosing variables (sparsity) by setting $\beta_r = 0$ when q = 1.

Regression Methods

Basic geometry II

Penalised solutions (red dots) for q=2, 1, 0.45, with contours of $D(\beta)$ in grey and solution contour for $\|\beta\|_q$ in red.



As $\lambda \to \infty$ the constraint tightens and the red contours shrink around the origin, and as $\lambda \to 0$ the constraint relaxes and the $\tilde{\beta}_{\lambda}$ tends to the unconstrained estimate.

Regression Methods

Autumn 2024 - slide 217

Lasso

☐ The lasso (least absolute shrinkage and selection operator) objective function can be written as

$$L = \frac{1}{2} ||y - X\beta||_2 + \lambda ||\beta||_1,$$

so suppose we have minimised this for some λ_0 , giving active set $A=\{r:\tilde{\beta}_r\neq 0\}$ and

$$L = \frac{1}{2}(y - X_A \tilde{\beta}_A)^{\mathrm{T}}(y - X_A \tilde{\beta}_A) + \lambda \sum_{r \in A} |\tilde{\beta}_r|,$$

and now we aim to decrease λ (i.e., to relax the constraint).

 \square Now $\mathrm{d}|x|/\mathrm{d}x = \mathrm{sign}(x)$, so when

$$\frac{\mathrm{d}L}{\mathrm{d}\tilde{\beta}_A} = X_A^{\mathrm{T}}(X_A\tilde{\beta}_A - y) + \lambda \operatorname{sign}(\tilde{\beta}_A) = 0,$$

we have

$$\tilde{\beta}_A = (X_A^{\mathrm{T}} X_A)^{-1} X_A^{\mathrm{T}} y - \lambda (X_A^{\mathrm{T}} X_A)^{-1} \mathrm{sign}(\tilde{\beta}_A) = b - \lambda a,$$

say, i.e., $\tilde{\beta}_A$ is linear in λ until A changes.

- \square A changes on deleting a column X_r from X_A or on adding one from its complement X_{A^c} .
- $\square \quad \mathrm{sign}(\tilde{\beta}_A)$ only changes when (say) $\tilde{\beta}_r$ passes through zero, but r leaves A when $\tilde{\beta}_r=0$.

Regression Methods

Lasso algorithm

- \square A variable in A is deleted if a component of $\tilde{\beta}_A = b \lambda a$ hits zero as λ decreases from λ_0 , which occurs at $\lambda_- = \max_{\lambda < \lambda_0} b_r/a_r$.
- \square If X_r is the rth column of X, then r will enter A if adding $X_r\beta_r$ decreases L, i.e., if

$$\frac{\mathrm{d}L}{\mathrm{d}\beta_r} = X_r^{\mathrm{\scriptscriptstyle T}}(X\beta - y) + \lambda \mathrm{sign}(\beta_r) \quad \begin{cases} <0, & \beta_r > 0, \\ >0, & \beta_r < 0, \end{cases}$$

so β_r remains inactive if $|X_r^{\mathrm{T}}(y-X\beta)| \leq \lambda$.

 \square Thus as λ decreases, A changes when for some r in the complement A^c of A we have

$$X_r^{\mathrm{T}}(y - X_A \tilde{\beta}_A) = \pm \lambda,$$

or, setting $\tilde{\beta}_A = b - \lambda a$,

$$X_{A^c}^{\mathrm{T}}(y - X_A b) + \lambda (X_{A^c}^{\mathrm{T}} X_A a \pm 1) = 0 \implies c + \lambda (d \pm 1) = 0,$$

say: the next variable is added when $\lambda = \lambda_+ = \max_{\lambda < \lambda_0} \{-c_r/(d_r \pm 1)\}.$

- \square Hence if $s = \operatorname{sign}(\beta)$, the algorithm decreases λ from
 - the highest λ at which the a first variable is active, and defines the A and s, then
 - finds the next λ at which A changes, stores it and the corresponding $\tilde{\beta}$, updating A and s.

Regression Methods

Autumn 2024 - slide 219

Practical matters and thresholding

- □ Usually
 - λ is chosen by dividing the data into training and testing subsets and minimising some measure of prediction error for the test subset,
 - $-\ \ y$ is centered and X has no column of ones, and
 - the columns of X are standardized to have zero mean and unit variance what this means in terms of interpreting the components of β is then unclear!
- We can think of penalised estimators as using different sorts of **thresholding** functions, where $\widehat{\beta}$ is replaced by $\widetilde{\beta} = g_{\lambda}(\widehat{\beta})$ and (conceptually)
 - for the lasso there is soft thresholding,

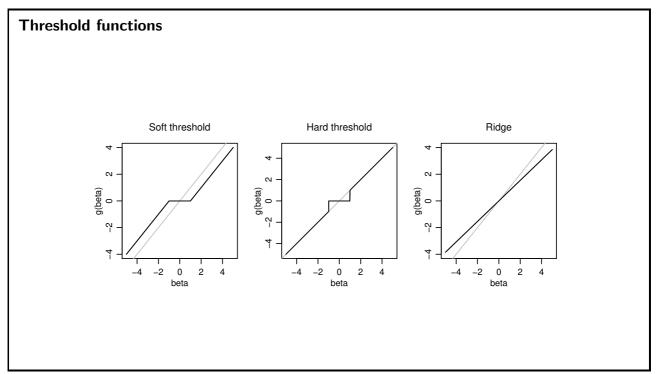
$$g_{\lambda}(u) = \begin{cases} 0, & |u| < \lambda, \\ \operatorname{sign}(u)(|u| - \lambda), & \text{otherwise,} \end{cases}$$

- for variable selection there is hard thresholding,

$$g_{\lambda}(u) = \begin{cases} 0, & |u| < \lambda, \\ u, & \text{otherwise,} \end{cases}$$

- for ridge regression there is shrinkage, $g(u) = u/(1 + \lambda)$.

Regression Methods

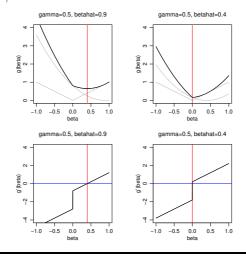


Regression Methods

Autumn 2024 - slide 221

Soft thresholding

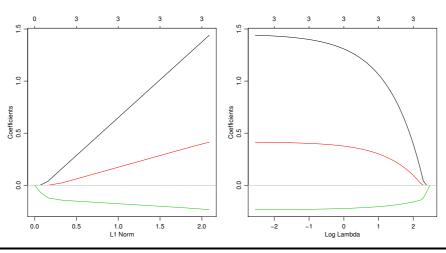
Top panels: the sum $g(\beta)$ of the L_1 penalty and the least squares function (both in grey) is the black line, which has a cusp at $\beta=0$. If the left- and right-hand derivatives of the sum are equal at zero, then the minimiser (at the red vertical line) is non-zero, but not otherwise. Bottom panels: the derivative $g'(\beta)=0$ when $\beta=\tilde{\beta}$.



Regression Methods

Example: cement data

 \square Estimated coefficients for lasso fit against L_1 norm and λ :



Regression Methods

Autumn 2024 - slide 223

Comments

Least angle regression (LAR) is similar to the lasso, and can compute the lasso solution path for all λ in $O(n^3)$ operations (faster than ridge, $O(np^2)$, when $p \gg n$).

Theory: one can ask about the properties of $\tilde{\beta}_{\lambda}$ in suitable settings (e.g., $n,p\to\infty$ with $p/n\to c>0$). Then under certain conditions one can show that lasso variable is consistent (i.e., the probability that the variables with $\beta_r\neq 0$ are selected tends to 1), but that the $\tilde{\beta}_{\lambda}$ themselves are inconsistent (because soft thresholding implies that $|\tilde{\beta}_{\lambda,r}|$ is systematically smaller than $|\beta_r|$).

☐ Many (many!) variants and related procedures exist to overcome such problems.

Computation: lasso and elastic net penalisations available in R package glmnet and extend to generalized linear models and more general regressions (later).

☐ For any regression model we can define the degrees of freedom as

$$\sigma^{-2} \sum_{j=1}^{n} \operatorname{cov}(y_j, \widehat{y}_j) = \operatorname{tr}\{\operatorname{cov}(y, \widehat{y})\} / \sigma^2;$$

this reduces to previous definitions but can be computed in more situations.

When $D(\beta)$ is a general loss function (e.g., a negative log likelihood for a GLM), the exact algorithm above is replaced by a **coordinate descent algorithm** that updates each $\tilde{\beta}_r$ in turn, with the other components fixed. This too is very efficient.

Regression Methods

3.4 Splines slide 225

Basis functions

 \square We seek to estimate a function $\mu(x)$ based on data $(x_1,y_1),\ldots,(x_n,y_n)$.

There are n parameters $\mu_1 = \mu(x_1), \dots, \mu_n = \mu(x_n)$ (plus noise, ...), so we assume that $\mu(x)$ belongs to a suitable class of functions, defined for $x \in \mathcal{X}$.

☐ Simple linear model is

$$\mu_{n\times 1} = B_{n\times p}\beta_{n\times 1}, \quad \operatorname{rank}(B) = p \le n,$$

with the columns of B evaluations at x_1, \ldots, x_n of basis functions.

 \Box The basis functions may be

- global (e.g., polynomials, trigonometric/Fourier functions),
- local (e.g., splines),
- multiscale (e.g., wavelets).
- \square We choose the basis for
 - suitability for the problem at hand (e.g., suitably smooth), and
 - computational reasons—want fast, preferably $\mathcal{O}(n)$, handling of $n \times n$ matrices.
- ☐ Focus on **spline functions**, on which there is a huge literature.

Regression Methods

Autumn 2024 - slide 226

Aside: Polynomial regression

 \square Classical approach is to fit a polynomial of degree p-1, i.e.,

$$\mu(x_j) = \beta_0 + \beta_1 x_j + \dots + \beta_{p-1} x_j^{p-1},$$

and choose $\beta_0,\ldots,\beta_{p-1}$ to minimise the sum of squares

$$\sum_{j=1}^{n} \{y_j - \mu(x_j)\}^2 = \sum_{j=1}^{n} \{y_j - (\beta_0 + \beta_1 x_j + \dots + \beta_{p-1} x_j^{p-1})\}^2,$$

giving $\widehat{\beta}_{p \times 1} = (B^{\mathrm{T}}B)^{-1}B^{\mathrm{T}}y$, where (j,i) element of $n \times p$ matrix B is x_j^{i-1} .

☐ Comments:

- easily copes with missing values/unequally spaced observations;
- use orthogonal polynomials to avoid numerical problems if n, k large;
- sensitivity to observations at extremities of series often leads to poor fit;
- usually doesn't work well because infinite differentiability everywhere is generally unnecessarily restrictive.

Regression Methods

Piecewise linear basis

 \square Place **knots** of a univariate x at $x_1^* < \cdots < x_K^*$, and define **tent functions**

$$b_1(x) = \begin{cases} (x_2^* - x)/(x_2^* - x_1^*), & x_1^* \le x \le x_2^*, \\ 0, & \text{otherwise}, \end{cases}$$

$$b_k(x) = \begin{cases} (x - x_{k-1}^*)/(x_k^* - x_{k-1}^*), & x_{k-1}^* < x \le x_k^*, \\ (x_{k+1}^* - x)/(x_{k+1}^* - x_k^*), & x_k^* < x \le x_{k+1}^*, \end{cases}$$

$$b_K(x) = \begin{cases} (x - x_{K-1}^*)/(x_K^* - x_{K-1}^*), & x_{K-1}^* \le x \le x_K^*, \\ 0, & \text{otherwise} : \end{cases}$$

$$correspondent conductors and take value 1 at x^* .$$

these are non-zero only in (x_{k-1}^*, x_{k+1}^*) (compact support) and take value 1 at x_k^* .

 \square An exact linear interpolant of data y_1,\ldots,y_K at the knots is the function

$$\mu(x) = \sum_{k=1}^{K} b_k(x) y_k = B(x)^{\mathrm{T}} y,$$

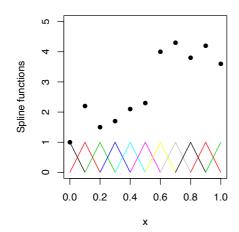
which by construction

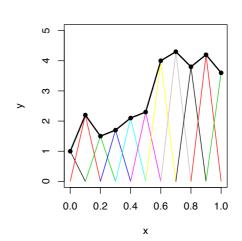
- passes through the points (x_k^st,y_k) and
- is linear between the knots.

Regression Methods

Autumn 2024 – slide 228

Piecewise linear basis II





- \square Left: piecewise linear basis functions $b_k(x)$ and data (x_k^*, y_k) .
- \square Right: functions $b_k(x)y_k$ and linear interpolant (bold).

Regression Methods

Statistical use

- \square Aim for summary of the n observations, so interpolation not useful.
- \square Could use K < n knots, but fit tends to depend heavily on their locations, so better to use high(ish) K and impose structure by penalising roughness of $\mu(x)$:

$$\widehat{\beta}_{\lambda} = \operatorname{argmin}_{\beta} \left\{ \|y - B\beta\|^2 + \lambda \sum_{k=2}^{K-1} \left\{ \mu(x_{k-1}^*) - 2\mu(x_k^*) + \mu(x_{k+1}^*) \right\}^2 \right\}.$$

- \Box The second term sums squared numerical second derivatives at the internal knots, and λ imposes the degree of penalisation:
 - $\lambda = 0$ (no penalty) gives the interpolant,
 - $-\lambda \to \infty$ forces the second derivatives to be zero, so gives a straight-line fit.
- \square On setting $\beta_k = \mu(x_k^*)$ and writing

$$\begin{pmatrix} \beta_1 - 2\beta_2 + \beta_3 \\ \beta_2 - 2\beta_3 + \beta_4 \\ \beta_3 - 2\beta_4 + \beta_5 \\ \vdots \end{pmatrix} = \begin{pmatrix} 1 & -2 & 1 & 0 & 0 & 0 & \cdots \\ 0 & 1 & -2 & 1 & 0 & 0 & \cdots \\ 0 & 0 & 1 & -2 & 1 & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \vdots \end{pmatrix} = D_{(K-2)\times K}\beta_{K\times 1},$$

the penalty is $\sum_{k=2}^{K-1} (\beta_{k-1} - 2\beta_k + \beta_{k+1})^2 = (D\beta)^{\mathrm{T}} D\beta = \beta^{\mathrm{T}} D^{\mathrm{T}} D\beta = \beta^{\mathrm{T}} S\beta$, say.

Regression Methods

Autumn 2024 – slide 230

Penalized fit

 \square The penalty matrix S is of side $K \times K$ but of rank K-2, because

$$S1_K = Sx_{K \times 1}^* = 0_K$$
:

the null space of S consists of all straight lines $\beta_0 1_K + \beta_1 x^*$, which are unpenalised.

☐ Hence (recalling ridge regression),

$$\widehat{\beta}_{\lambda} = \operatorname{argmin}_{\beta} \left\{ \|y - B\beta\|^2 + \lambda \beta^{\mathrm{T}} S\beta \right\} = (B^{\mathrm{T}} B + \lambda S)^{-1} B^{\mathrm{T}} y$$

giving

fitted values
$$\widehat{y} = B\widehat{eta}_{\lambda} = B(B^{\mathrm{T}}B + \lambda S)^{-1}B^{\mathrm{T}}y = H_{\lambda}y,$$

equivalent degrees of freedom
$$df_{\lambda} = tr(H_{\lambda}) = \sum_{k=1}^{K} \frac{1}{1 + \eta_k \lambda},$$

where

- $\eta_1 \leq \cdots \leq \eta_K \in [0,1]$ are the eigenvalues of $(B^{\mathrm{T}}B)^{-1/2}S(B^{\mathrm{T}}B)^{-1/2}$,
- $\eta_1 = \eta_2 = 0$, corresponding to the null space of S, so
- df_{λ} is monotone decreasing in λ , with

$$(\lambda = 0)$$
 $K \ge \mathrm{df}_{\lambda} \ge 2$ $(\lambda \to \infty)$.

Regression Methods

Higher-order splines

 $\ \square$ The pth degree spline basis with knots $x_1^* < \cdots < x_K^*$ is

$$1, x, \ldots, x^p, (x - x_1^*)_+^p, \ldots, (x - x_K^*)_+^p,$$

where $u_+ = \max(u, 0)$ is the **positive part function**.

- \Box The resulting basis matrix B is highly collinear and gives an implausible statistical model.
- \square B-spline bases span the same linear space, but have better numerical properties. They are defined by adding boundary knots x_0^* and x_{K+1}^* and setting up an augmented knot sequence

$$\tau_1 \le \dots \le \tau_M \le x_0^* \le \tau_{M+1} = x_1^* \le \dots \le \tau_{M+K} = x_K^* \le x_{K+1}^* \le \tau_{K+1+M} \le \dots \le \tau_{K+2M};$$

typically the τ_k outside $[x_0^*, x_{K+1}^*]$ are set to the boundary knot values. Then

$$B_{k,1}(x) = I(\tau_k \le x < \tau_{k+1}), \quad k = 1, \dots, K + 2M - 1,$$

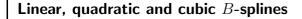
$$B_{k,m}(x) = \frac{x - \tau_k}{\tau_{k+m-1} - \tau_k} B_{k,m-1}(x) + \frac{\tau_{k+m} - x}{\tau_{k+m} - \tau_{k+1}} B_{k+1,m-1}(x), \quad k = 1, \dots, K + 2M - m,$$

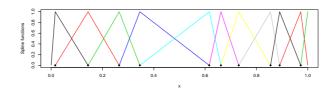
where we set $B_{k,1} \equiv 0$ if $\tau_k = \tau_{k+1}$ (avoiding division by zero).

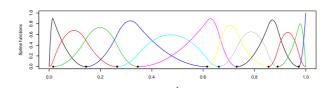
- \square Cubic splines (p=3, M=4) give visually smooth functions.
- \square K=10 on the next slide, with M=2 (linear), M=3 (quadratic) and M=4 (cubic), and the τ_k set to equal the boundary knots.

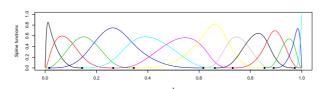
Regression Methods

Autumn 2024 - slide 232









Regression Methods

Natural cubic spline

 \square Suppose the x_j are distinct (no loss of generality) and

$$a < x_1 < \dots < x_n < b, \quad \mathcal{X} = [a, b] \subset \mathbb{R}.$$

- A **natural cubic spline** adds the constraint that the function is linear outside $[x_1, x_n]$, and thus avoids high variance due to quadratic and higher terms outside this interval.
- ☐ A natural cubic spline
 - has K = n knots, at $x_1 < \cdots < x_n$,
 - is a cubic polynomial on each interval between knots,
 - is continuous, with continuous first and second derivatives at each knot, and
 - is linear on $[a, x_1]$ and $[x_n, b]$, with zero second and third derivatives at x_1 and x_n ,
 - has

$$2+4(n-1)+2$$
 parameters $-3n$ linear constraints $=n$

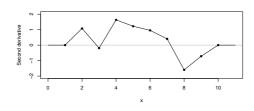
degrees of freedom (df), which can be split into

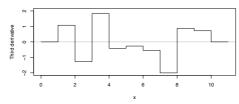
- ≥ 2 df for a linear fit, plus
- \triangleright n-2 df for the second derivatives $\mu''(x_2), \ldots, \mu''(x_{n-1})$.

Regression Methods

Autumn 2024 - slide 234

Natural cubic spline





- A natural cubic spline may be constructed by integrating a linear second derivative function $\mu''(x)$ which is determined by $\mu''(x_2), \dots, \mu''(x_{K-1})$ and because $\mu''(x) \equiv 0$ for $x \notin (x_1, x_K)$.
- On integrating twice we gain two constants: $\mu(x) = \beta_0 + \beta_1 x + \int_0^x \int_0^{x'} \mu''(u) du dx'$.
- \square Above $x_1 = 1, \dots, x_{10} = 10$, so the spline is determined by $\mu''(2), \dots, \mu''(9)$ and the line.

Regression Methods

Optimality of natural cubic splines

- Let $S_2(\mathcal{X})$ denote the set of functions μ differentiable on $\mathcal{X} = [a,b]$ with absolutely continuous first derivative μ' : i.e., there exists an integrable function μ'' such that $\int_a^x \mu''(u) du = \mu'(x) \mu'(a)$ for $x \in \mathcal{X}$.
- \square Clearly any μ with two continuous derivatives on \mathcal{X} lies in $\mathcal{S}_2(\mathcal{X})$.

Theorem 32 Suppose $n \geq 2$, that $a < x_1 < \cdots < x_n < b$, and that μ is the natural cubic spline interpolating y_1, \ldots, y_n at x_1, \ldots, x_n . If $\tilde{\mu} \in \mathcal{S}_2(\mathcal{X})$ also interpolates the y_j , then

$$\int_{\mathcal{X}} \tilde{\mu}''^2 \ge \int_{\mathcal{X}} \mu''^2,$$

with equality iff $\tilde{\mu} \equiv \mu$.

 \square Thus μ minimises the **roughness penalty** $\lambda \int_{\mathcal{X}} \mu''^2$ in a larger class of functions than that to which it belongs, making it a natural choice as an interpolant, because minimising

$$\sum_{j=1}^{n} \{y_j - \tilde{\mu}(x_j)\}^2 + \lambda \int_{\mathcal{X}} \tilde{\mu}''(x)^2 dx$$

for $\tilde{\mu} \in \mathcal{S}_2(\mathcal{X})$ will automatically result in a natural cubic spline μ : if $\tilde{\mu}(x_j) = \mu(x_j)$, then the penalty is reduced by using μ .

Regression Methods

Autumn 2024 - slide 236

Note to Theorem 32

Let $\nu = \tilde{\mu} - \mu \in \mathcal{S}_2(\mathcal{X})$, and note that $\nu(x_j) = 0$ for each j, since $\mu(x_j) = \tilde{\mu}(x_j) = y_j$. The natural boundary conditions imply that $\mu''(a) = \mu''(b) = 0$, so integration by parts yields

$$0 = \left[\mu''(x)\nu'(x) \right]_a^b = \int_{\mathcal{X}} (\mu''\nu')' = \int_{\mathcal{X}} \mu''\nu'' + \int_{\mathcal{X}} \mu'''\nu',$$

and hence the facts that μ''' is piecewise constant and that $\nu(x_i) = 0$ yields

$$\int_{\mathcal{X}} \mu'' \nu'' = -\int_{\mathcal{X}} \mu''' \nu' = -\sum_{j=1}^{n-1} \mu'''(x_j^+) \int_{x_j}^{x_{j+1}} \nu' = -\sum_{j=1}^{n-1} \mu'''(x_j^+) \{\nu(x_{j+1}) - \nu(x_j)\} = 0.$$

Hence

$$\int_{\mathcal{X}} \tilde{\mu}''^2 = \int_{\mathcal{X}} (\mu'' + \nu'')^2 = \int_{\mathcal{X}} \mu''^2 + 2 \int_{\mathcal{X}} \mu'' \nu'' + \int_{\mathcal{X}} \nu''^2 = \int_{\mathcal{X}} \mu''^2 + \int_{\mathcal{X}} \nu''^2 \geq \int_{\mathcal{X}} \mu''^2,$$

wth equality iff $\nu''(x) \equiv 0$. This occurs iff $\nu(x)$ is linear, but since $\nu(x_j) = 0$ at at least two points, $\nu(x) = 0$ for all $x \in \mathcal{X}$.

Regression Methods

Autumn 2024 - note 1 of slide 236

More splines

- □ Sometimes cyclic effects (e.g., seasonality, diurnal variation) must be modelled smoothly, so (e.g.) December joins smoothly onto January. Then the penalty and spline basis must be modified accordingly, to give a cyclic (cubic) spline.
- P-splines are a version of B-splines (usually with equally-spaced knots) in which a difference penalty is applied to the parameters to control the wiggliness of μ , e.g.,

$$\sum_{k=1}^{K-1} w_k (\beta_{k+1} - \beta_k)^2 = \beta^{\mathrm{T}} D^{\mathrm{T}} W D \beta, \quad \text{with} \quad D = \begin{pmatrix} -1 & 1 & 0 & 0 \cdots \\ 0 & -1 & 1 & 0 & \cdots \\ \cdot & \cdot & \cdot & \cdot & \cdots \end{pmatrix},$$

and $W = \operatorname{diag}(w_1, \dots, w_{K-1})$. These are easy to set up and flexible, but messy if the knots are not equi-spaced, and the penalty is less readily interpreted.

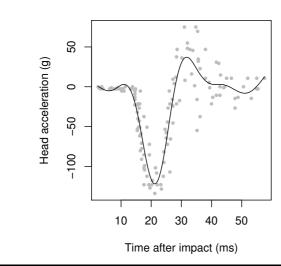
- \square For an adaptive spline we can let $w_k \equiv w_k(x)$ vary with x, for example setting $w(x) = B(x)\lambda_{L\times 1}$ and thus having $D^{\mathrm{T}}WD = \sum_l \lambda_l D^{\mathrm{T}}\mathrm{diag}\{B_l(x)\}D$, where $B_l(x)$ is the lth column of B(x), then estimating the vector λ .
- ☐ Other possibilities include (Wood, 2017, Chapter 5)
 - shape-constrained splines to impose, e.g., monotonicity on the fit;
 - thin-plate, Duchon and tensor product splines used in spatial problems; and
 - soap film splines used when smoothing over complex domains.

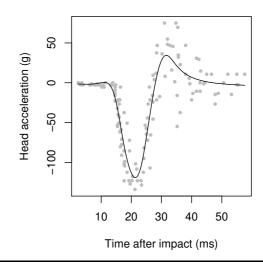
Regression Methods

Autumn 2024 - slide 237

Motorcycle data: adaptive fit

Standard (left) and adaptive (right) spline fits, the latter with K=40 and L=5:





Regression Methods

Generalisations

- \square We've discussed estimation of a single function $\mu(x)$, but in applications we may have
 - covariates to be treated parametrically,
 - several smooth functions,
 - non-normal response variable,
 - random effects (later).
- ☐ To include ordinary covariates and allow for weights, we write

$$y \mid b \sim (B\theta, \sigma^2 W), \quad B\theta = X\beta + Zb,$$

where B=(X,Z) is $n\times d$, $\theta=(\beta^{\rm T},b^{\rm T})^{\rm T}$ is $d\times 1$, d=p+q and

- the $n \times p$ matrix X represents the ordinary covariates, plus any unpenalised columns for smooth components,
- the $p \times 1$ parameter vector β is unpenalized,
- the $n \times q$ matrix Z represents the bases for any smooth functions,
- the $q \times 1$ vector b is penalized,
- the n imes n diagonal matrix $W = \mathrm{diag}(w_1, \ldots, w_n)$ contains positive weights,

and everything 'goes through as before'.

Regression Methods

Autumn 2024 - slide 240

Additivity and identifiability

☐ Consider the additive model

$$E(y) = \mu_1(x) + \mu_2(z),$$

where μ_1 , μ_2 belong to suitable classes of smooth functions; if

$$x \equiv \text{time}, \quad z \equiv \text{space},$$

then μ_1 is defined on $\mathcal{X}_1 \subset \mathbb{R}$ and μ_2 is defined on $\mathcal{X}_2 \subset \mathbb{R}^2$.

☐ There is an identifiability problem, since we could map

$$\mu_1(x) \mapsto \mu_1(x) + a, \quad \mu_2(z) \mapsto \mu_2(z) - a, \quad a \in \mathbb{R},$$

and the fitted values would not change, so we must constrain μ_1 and μ_2 .

 \square As before, we use bases for μ_1 and μ_2 , writing

$$E(y) = Zb = \begin{pmatrix} Z_1(x) & Z_2(z) \end{pmatrix} \begin{pmatrix} b_1 \\ b_2 \end{pmatrix},$$

where we penalise the q_1 elements of b_1 and the q_2 elements of b_2 .

Regression Methods

Ensuring identifiability

☐ The identifiability problem is solved by **centering** the fitted smooth, i.e., enforcing

$$1_n^{\mathrm{T}} Z_{n \times q} b_{q \times 1} = 0$$

for each smooth term.

 \square In general we can use a QR decomposition. If $C_{a\times q}b_{q\times 1}=0_{a\times 1}$, with a< q, write

$$C_{q \times a}^{\mathrm{T}} = Q_{q \times q} R_{q \times a} = \begin{pmatrix} Q_1 & Q_2 \end{pmatrix} \begin{pmatrix} R_1 \\ 0 \end{pmatrix},$$

where Q is orthogonal,

- Q_1 has dimension $q \times a$,
- Q_2 has dimension $q \times (q-a)$, and
- R_1 has dimension $a \times a$ and is upper triangular.

Then if we set $b_{q\times 1}=Q_2b'_{(q-a)\times 1}$, we have

$$Cb = R^{\mathrm{T}}Q^{\mathrm{T}}b = \begin{pmatrix} R_1^{\mathrm{T}} & 0 \end{pmatrix} \begin{pmatrix} Q_1^{\mathrm{T}} \\ Q_2^{\mathrm{T}} \end{pmatrix} Q_2b' = \begin{pmatrix} R_1^{\mathrm{T}} & 0 \end{pmatrix} \begin{pmatrix} 0 \\ I_{q-m} \end{pmatrix}b' = 0.$$

 \square Thus the constraint is satisfied if we replace $Z_{n\times q}$ by $(ZQ_2)_{n\times (q-1)}$; this reduces b to dimension $(q-1)\times 1$.

Regression Methods

Autumn 2024 - slide 242

Penalty formulation

☐ Minimise

$$(y - B\theta)^{\mathrm{T}} W (y - B\theta) + \theta^{\mathrm{T}} S_{\lambda} \theta = (y - X\beta - Zb)^{\mathrm{T}} W (y - X\beta - Zb) + \theta^{\mathrm{T}} S_{\lambda} \theta$$

where S_{λ} is a sum of symmetric positive semi-definite $d \times d$ matrices S_m , such that

$$\theta^{\mathrm{T}} S_{\lambda} \theta = \theta^{\mathrm{T}} \left(\sum_{m=1}^{M} \lambda_m S_m \right) \theta = \sum_{m=1}^{M} \lambda_m b_m^{\mathrm{T}} S_m^* b_m, \quad \lambda_m \ge 0,$$

where S_m^* is the non-zero diagonal block of S_m and b has sub-vectors $b_1,\ldots,b_M.$

 \square With M=2, β , b_1 and b_2 are vectors of respective lengths p, q_1 and q_2 , and S_1^* and S_2^* are square matrices of sides q_1 and q_2 , so

$$\theta = \begin{pmatrix} \beta \\ b_1 \\ b_2 \end{pmatrix}, \quad S_{\lambda} = \lambda_1 S_1 + \lambda_2 S_2 = \lambda_1 \begin{pmatrix} 0 & 0 & 0 \\ 0 & S_1^* & 0 \\ 0 & 0 & 0 \end{pmatrix} + \lambda_2 \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & S_2^* \end{pmatrix},$$

with S_1 and S_2 partitioned conformably with θ .

- \square Let S_{λ}^* denote the $q \times q$ corner of S_{λ} corresponding to b; here $S_{\lambda}^* = \operatorname{diag}(\lambda_1 S_1^*, \lambda_2 S_2^*)$.
- \square Note that $|S_{\lambda}|_{+} = |S_{\lambda}^{*}|_{+}$.

Regression Methods

Estimation

 \Box For fixed λ , the minimiser and fitted values for

$$(y - B\theta)^{\mathrm{T}}W(y - B\theta) + \theta^{\mathrm{T}}S_{\lambda}\theta$$

are

$$\widehat{\theta}_{\lambda} = (B^{\mathsf{T}}WB + S_{\lambda})^{-1}B^{\mathsf{T}}Wy, \quad \widehat{y}_{\lambda} = B\widehat{\theta}_{\lambda} = B(B^{\mathsf{T}}WB + S_{\lambda})^{-1}B^{\mathsf{T}}Wy = H_{\lambda}y.$$

 \Box If the unpenalized least squares estimator $\widehat{ heta}=(B^{ ext{ iny T}}WB)^{-1}B^{ ext{ iny T}}Wy$ exists, then

$$\widehat{\theta}_{\lambda} = (B^{\mathrm{T}}WB + S_{\lambda})^{-1}B^{\mathrm{T}}WB\widehat{\theta} = \widehat{\theta} - (B^{\mathrm{T}}WB + S_{\lambda})^{-1}S_{\lambda}\widehat{\theta} = P_{\lambda}\widehat{\theta},$$

and if \widehat{y} is the unpenalised fitted value, then

$$\widehat{y}_{\lambda} = \widehat{y} - B(B^{\mathsf{T}}WB + S_{\lambda})^{-1}S_{\lambda}\widehat{\theta}.$$

- ☐ Now we must decide
 - how many degrees of freedom for each smooth?
 - how to select the smoothing parameters?

Regression Methods

Autumn 2024 - slide 244

Amount of smoothing

☐ We write

$$\widehat{\theta}_{\lambda} = P_{\lambda}\widehat{\theta},$$

say, where P_{λ} shows how penalisation shrinks $\widehat{\theta}$ towards $\widehat{\theta}_{\infty} = (\widehat{\beta}^{\mathrm{T}}, 0^{\mathrm{T}})^{\mathrm{T}}$.

- $\Box \quad \text{If } \lambda \approx 0 \text{, then } P_{\lambda} \approx I_{p+q} \text{ and the degrees of freedom of the two fits are both } \approx p+q \text{, but as } \\ \lambda \to \infty \text{, } P_{\lambda} \text{ tends to the projection matrix onto the column space of } X_{n \times p}.$
- \square On slide 193 with just one smooth term we defined

$$\operatorname{edf}_{\lambda} = \operatorname{tr}(H_{\lambda}) = \operatorname{tr}(P_{\lambda}) = \sum_{r=1}^{p+q} P_{\lambda,rr} \in (p, p+q),$$

which gives the usual definition for a linear model.

- If $\theta^{\mathrm{T}} = (\beta^{\mathrm{T}}, b_1^{\mathrm{T}}, \dots, b_M^{\mathrm{T}})$, we define the **effective degrees of freedom** edf_{λ_m} associated to the mth smooth as being the sum of those $P_{\lambda,rr}$ that correspond to the elements of b_m in θ .
- \square To choose the vector λ we use either
 - $CV(\lambda)$ or $GCV(\lambda)$ (second-order assumptions),
 - REML (normal-theory assumptions).
- \square Must optimise over (log) λ , e.g., by grid search (CV/GCV) or other methods (REML).

Regression Methods

Inference

- \square So far we have discussed only 'point estimation' of a smooth function $\mu(x)$, but in applications we also want
 - pointwise confidence intervals for smooth functions,
 - overall confidence bands for (say) $\{\mu(x):x\in\mathcal{S}\}$, where \mathcal{S} is some subset of \mathcal{X} , and
 - tests of hypotheses such as 'is the spline part needed?' and 'is the curve monotonic?'
- ☐ Under the normal model we have the Bayesian interpretation from slide 191,

$$\theta \mid y, \sigma^2, \lambda \sim \mathcal{N}_d\left(\widehat{\theta}_{\lambda}, V_{\lambda}\right), \quad V_{\lambda} = \sigma^2 (B^{\mathrm{T}} W B + S_{\lambda})^{-1},$$

from which we can simulate to find bounds for any function $A(\theta)$.

 \Box If $A(\theta) = A_{m \times d}\theta$, then

$$A\theta \mid y, \sigma^2, \lambda \sim \mathcal{N}_m(A\widehat{\theta}_{\lambda}, AV_{\lambda}A^{\mathrm{T}}),$$

and generalisation of (10) gives that its mean square error is

$$MSE = E\left(\|A\widehat{\theta}_{\lambda} - A\theta\|^{2}\right) = tr(AV_{\lambda}A^{T}),$$

which takes into account both estimation error and prior uncertainty about θ .

Regression Methods

Autumn 2024 - slide 246

Average coverage probabilities

- \square Bayesian credible intervals have good frequentist properties, averaged over the domain of x.
- \square Let the random index variable J choose the m rows a_j^{T} of A with equal probabilities, and aim to choose constants d and c_j such that the average coverage probability

$$ACP = P\left\{ |a_J^T \widehat{\theta}_{\lambda} - a_J^T \theta| \le dc_J \right\} = 1 - \alpha;$$

i.e., ACP has a desired value averaged over y, θ and J.

☐ The random variable

$$a_J^{\mathrm{T}}(\widehat{\theta}_{\lambda} - \theta)/c_J = a_J^{\mathrm{T}}\{\widehat{\theta}_{\lambda} - \mathrm{E}(\widehat{\theta}_{\lambda})\}/c_J + a_J^{\mathrm{T}}\{\mathrm{E}(\widehat{\theta}_{\lambda}) - \theta\}/c_J = S + T,$$

say, has a mixture of normal distributions, where

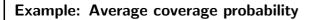
- S is approximately normal and E(S) = 0,
- T is random (because of J) with $E(T) \approx 0$, but $var(T) \ll var(S)$.
- \square We now choose $C = \operatorname{diag}(c_1, \ldots, c_m) = \operatorname{diag}(AV_{\lambda}A^{\mathrm{T}})^{1/2}$, so that

$$\operatorname{var}(S+T) \approx m^{-1} \operatorname{E} \left\{ \|C^{-1} A(\widehat{\theta}_{\lambda} - \theta)\|^{2} \right\} = m^{-1} \operatorname{tr} \left(C^{-1} A V_{\lambda} A^{\mathrm{T}} C^{-1} \right) = 1,$$

and then setting $d=z_{1-\alpha/2}$ gives the required value for ACP.

 \square This ignores estimation error for σ^2 and λ .

Regression Methods



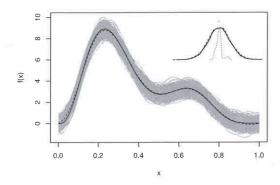


Figure 6.7 The Nychka (1988) idea. The main black curve shows a true function f(x), while the grey curves show 500 replicate spline estimates $\hat{f}(x)$. The dashed curve is $\mathbb{E}\hat{f}(x)$. Inset at top right are scaled kernel smooth estimates of the distributions of the sampling error, $\hat{f} - \mathbb{E}\hat{f}$ (continuous black); the bias, $\mathbb{E}\hat{f} - f$, evaluated at a random x (dotted) and $\hat{f} - f$ evaluated at a random x (dashed). In grey is the normal approximation to the dashed curve. Evaluation at a random x turns the bias into a random variable, which has substantially lower variance than the approximately normal $\hat{f} - \mathbb{E}\hat{f}$. Hence the sum of the randomized bias and sampling error is approximately normally distributed. The variance of this sum turns out to be well approximated by the Bayesian posterior covariance for \hat{f} .

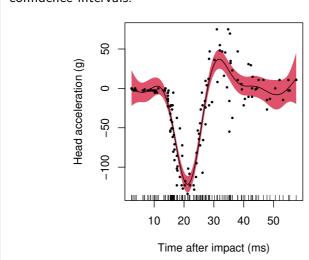
(Wood, 2017)

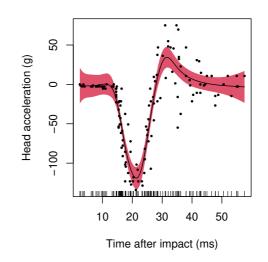
Regression Methods

Autumn 2024 - slide 248

Example: Motorcycle data

Standard (left) and adaptive (right) spline fits, the latter with K=40 and L=5, and 95% pointwise confidence intervals:





Regression Methods

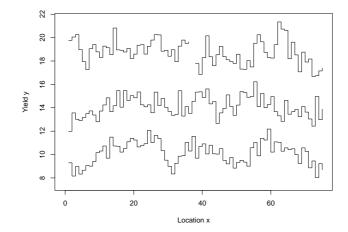
Plot yield at harvest for 75 varieties of spring barley sown in 3 blocks each of 75 plots:

Location x	Block 1		Blo	ck 2	Bloo	Block 3		
	Variety	Yield y	Variety	Yield y	Variety	Yield y		
1	57	9.29	49	7.99	63	11.77		
2	39	8.16	18	9.56	38	12.05		
3	3	8.97	8	9.02	14	12.25		
4	48	8.33	69	8.91	71	10.96		
5	75	8.66	29	9.17	22	9.94		
6	21	9.05	59	9.49	46	9.27		
7	66	9.01	19	9.73	6	11.05		
8	12	9.40	39	9.38	30	11.40		
9	30	10.16	67	8.80	16	10.78		
10	32	10.30	57	9.72	24	10.30		
11	59	10.73	37	10.24	40	11.27		
12	50	9.69	26	10.85	64	11.13		
13	5	11.49	16	9.67	8	10.55		
14	23	10.73	6	10.17	56	12.82		
15	14	10.71	47	11.46	32	10.95		
16	68	10.21	36	10.05	48	10.92		
17	41	10.52	64	11.47	54	10.77		
18	1	11.09	63	10.63	37	11.08		
:	:	:	:	:	:	<u>:</u>		

Regression Methods

Autumn 2024 - slide 250

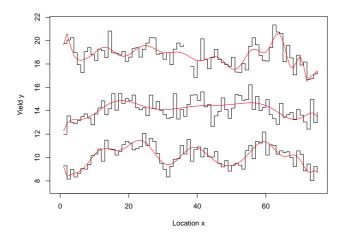
Example: Spring barley data



Yield as a function of location for the three blocks, with yields for blocks 2 and 3 offset by the addition of 4 and of 7 respectively. Value 37 in block 3 is missing.

Regression Methods

Spring barley data and polynomial fits



Yield as a function of location for the three blocks, with yields for blocks 2 and 3 offset by the addition of 4 and of 7 respectively, with fitted polynomials of degrees 20, 10 and 50.

Regression Methods

Autumn 2024 - slide 252

Example: Spring barley data

☐ We fit a model with parametric variety effects and smooth effects for the fertility patterns in the blocks,

$$y_{n\times 1} \sim (X_{n\times 75}\beta_{75\times 1} + Z_1b_1 + Z_2b_2 + Z_3b_3, \sigma^2 I_n),$$

where

- n = 224, as one of the responses is missing,
- X is a matrix of indicators (0/1) of which variety is in which plot in each block,
- $-\beta$ are the variety effects, with the model parametrized without an overall mean,
- Z_m of dimension $n \times (p_m + q_m)$ corresponds to the basis functions for the smooth in block m, and
- b_m are of dimensions $(p_m+q_m) imes 1$, for m=1,2,3, corresponding to the smooth effects, and
- $p_m + q_m = 9$ by default (after centering) when using gam in R package mgcv.
- Taking $p_m=2$ would correspond to null smooth $\beta_0+\beta_1x$ for each block (i.e., linear fertility pattern), but the identifiability constraints impose $\beta_0=0$. Hence in fact $p_m=1$ for a linear baseline smooth and the degrees of freedom for the smooth terms lie in [1,9] (see slide 255).

Regression Methods

```
Example: Spring barley data

library(SMPracticals)
data(barley)

library(mgcv)

# ML fit of variety as fixed effect, with GCV estimation of lambdas,
# with splines for fertility gradients within each block

fit.gcv <- gam(y~Variety-1+s(Location,by=Block),data=barley)

# fit of variety as fixed effect, with REML estimation of lambdas,
# with splines for fertility gradients within each block

fit <- gam(y~Variety-1+s(Location,by=Block),method="REML",data=barley)

# REML fit with variety as a random effect and splines for fertilities

fit.re <- gam(y~s(Variety,bs="re")+s(Location,by=Block),method="REML",data=barley)</pre>
```

Regression Methods

Autumn 2024 - slide 254

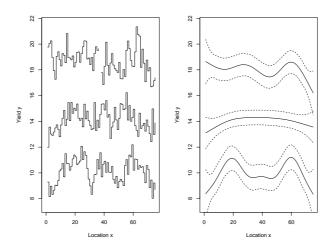
Example: Spring barley data

- Using GCV the smooths have $df_{\lambda}=8.3$, 6.8, 6.3, with $\widehat{\sigma}=0.65$ and AIC=513.1, the residual degrees of freedom is $224-75-8.3-6.8-6.3\approx 130.6$, with SEs around 0.4 for the estimated variety effects (0.54 for variety 27).
- Using REML the smooths have $df_{\lambda} = 7.2$, 3, 6.1, with $\hat{\sigma} = 0.66$ and AIC = 518.3, the residual degrees of freedom is 132.7, with SEs around 0.4 for the estimated variety effects (0.53 for variety 27).
- \square The estimated smoothing parameters are $\widehat{\lambda}_1=0.0029,\ \widehat{\lambda}_2=0.18$ and $\widehat{\lambda}_3=0.0078$.
- \Box The effective degrees of freedom for the smooth terms, with the totals:

Block					$P_{\lambda,rr}$					Total	
1	1.00	1.07	0.90	0.7	0.65	0.17	0.38	1.31	1	7.18	
2	0.61	0.21	0.12	-0.2	0.03	-0.26	0.01	1.49	1	3.00	
3	0.99	1.04	0.76	0.4	0.41	-0.18	0.18	1.47	1	6.07	

- \square The $P_{\lambda,rr}$ need not be positive, though their total for each smooth is positive.
- \square In applications it would be wise to check whether increasing q_m would lead to very different fits.

Regression Methods



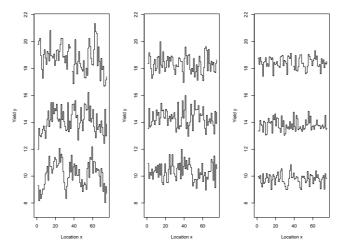
Left: data (offset by adding 4 and 8 to blocks 2 and 3).

Right: estimated fertility patterns (with estimated df 7.2, 3, 6.1) and 95% unconditional pointwise confidence intervals, fitted using REML. The intervals are wider for blocks 1 and 3.

Regression Methods

Autumn 2024 - slide 256

Example: Spring barley data

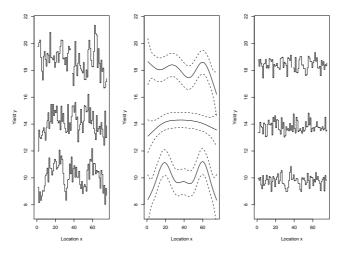


Left: data (offset by adding 4 and 8 to blocks 2 and 3).

Center: Estimated variety effects (also offset)

Right: residuals (also offset, and showing serial autocorrelation?)

Regression Methods



Left: data (offset by adding 4 and 8 to blocks 2 and 3). Center: estimated fertility patterns (REML), also offset.

Right: residuals.

Regression Methods

Autumn 2024 - slide 258

Example: Spring barley data

☐ Should the varieties be treated as randomly selected from a population of varieties?

 \square If so, we use the same basis matrix X as in the previous model, but add a penalty matrix $\lambda_{\beta}S_{\beta}$ and minimise the penalised sum of squares

$$(y - B\theta)^{\mathrm{T}}(y - B\theta) + \theta^{\mathrm{T}}S_{\lambda}\theta,$$

where

$$S_{\lambda} = \lambda_{\beta} S_{\beta} + \lambda_1 S_1 + \lambda_2 S_2 + \lambda_3 S_3,$$

where $S_{\beta} = \operatorname{diag}(I_{75}, 0)$.

 \Box The effective degrees of freedom are then 44.8 for β and 7.5, 3.9 and 6.4 for the splines.

 \Box The optimal smoothing parameters are $\widehat{\lambda}_{\beta}=1.76,~\widehat{\lambda}_{1}=0.0027,~\widehat{\lambda}_{2}=0.073$ and $\widehat{\lambda}_{3}=0.0070.$

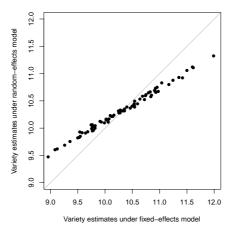
The fixed-effects model has 75 degrees of freedom for β , so this is substantial shrinkage; the estimated standard deviation drops from 0.65 to 0.39.

☐ The estimates under the random-effects model have standard errors around 0.31 (0.36 for variety 27), compared to 0.41 (0.54 for variety 27) for the fixed-effects model.

☐ The next slide compares the estimates.

Regression Methods

Comparison of estimated variety effects under fixed-effects and random-effects models:



Regression Methods

Autumn 2024 - slide 260

Comments

- ☐ Penalised estimation extends the basic smoothers to include
 - parametric terms in models,
 - several smooth terms,
 - spatial and more complex smoothing,
 - 'random effect' parameters,

and extends to generalized additive models in a natural way.

- \square The baseline variance σ^2 and smoothing parameter(s) λ are estimated using cross-validation under second-order assumptions or REML under normality.
- The empirical Bayes formulation allows inference on parameters and smooth functions in a unified way usually ignoring the uncertainty for σ^2 and λ is not too critical.
- \square In practice n and d may be very big, so direct matrix inversion is computationally painful, and then indirect methods (e.g., based on the Woodbury formula) are needed to compute $\widehat{\theta}_{\lambda}$ and V_{λ} .

Regression Methods

Background and motivation

- ☐ All the models so far have involved just one level of randomness, corresponding to 'measurement error' on individual responses.
- ☐ Complex layering of randomness can arise in applications, and then conclusions may depend on how it is dealt with.
- ☐ Two conceptually different set-ups (which may give the same models):
 - observational/experimental setup generates several layers of randomness;
 - we find it useful to treat the parameters of some model as drawn from a distribution.

The first concerns logical properties of the data, whereas the second is a modelling assumption.

Regression Methods

Autumn 2024 - slide 263

Example: Blood pressure

- \square Blood pressure data: P=25 patients each made V=16 visits to a clinic, and on each occasion their systolic and diastolic blood pressures were measured twice.
- ☐ Consider just the diastolic pressure. We expect there to be variation
 - between patients,
 - between visits within patients, and
 - between measurements within visits,

which we could model as

$$y_{pvm} = \mu + b_p + e_{pv} + \varepsilon_{pvm}, \quad p = 1, \dots, P, v = 1, \dots, V, m = 1, \dots, M,$$

where

- μ is the population mean diastolic blood pressure (DBP),
- b_p is the difference between the patient and population mean DBP,
- $-\ e_{pv}$ is the difference between this and the mean DBP on the vth visit, and
- ε_{pvm} is the difference between the mean DBP for the pth patient at the vth visit and the mth measurement on that visit.

and

$$b_p \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_b^2) \perp \!\!\!\perp e_{pv} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_e^2) \perp \!\!\!\perp \varepsilon_{pvm} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2).$$

Regression Methods

Example: Blood pressure patno patient visno dbp1 dbp2 sbp1 sbp2

Regression Methods

Autumn 2024 - slide 265

Fixed and random effects

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

- ☐ Times (min) for four chimpanzees to learn each of ten words.
- ☐ A possible model for log time is

$$y_{cw} \mid \alpha_c, \beta_w \stackrel{\text{ind}}{\sim} \mathcal{N}(\mu + \alpha_c + \beta_w, \sigma^2), \quad c = 1, \dots, C = 4, w = 1, \dots, W = 10.$$

- \square The α_c and/or the β_w would be considered as constant **fixed effects** if we were interested in the relative linguistic abilities of these particular chimps and/or if we planned further tests with these particular words.
- \square Either (or both) of the α_c and β_w might be considered to be **random effects** if they were thought to be sampled from a larger population whose variation is of interest.

Regression Methods

Two distinctions

- ☐ We distinguish **fixed** and **random** effects (above).
- ☐ We distinguish **nested** and **crossed** effects:
 - in the blood pressure data, replicate measurements at each visit are **nested** within visit, because there is no logical connection between $y_{p,v_1,1}$ and $y_{p,v_2,1}$ (we could permute the final index m within each patient/visit combination without changing the data structure). Likewise if we ignore any possible time effects between visits, we could consider that visits are nested within patients;
 - in the chimp data, the effects are **crossed**, because permuting chimps or words would entail permuting entire rows or columns of the data table: there is a logical connection between y_{c_1w} and y_{c_2w} , and between y_{cw_1} and y_{cw_2} ;
- □ In R syntax, with patient and visit number declared as factors, for nested effects we write

y ~ patient/visno

read as 'separate effects for visit number within the levels of patient' and for crossed effects with chimp and word declared as factors we write

y ~ chimp + word

Regression Methods

Autumn 2024 - slide 267

Nested model ANOVA

☐ For the nested model

$$y_{pvm} = \mu + b_p + e_{pv} + \varepsilon_{pvm}, \quad p = 1, \dots, P, v = 1, \dots, V, m = 1, \dots, M,$$

and with a dot and bar denoting averaging over that index, we write

$$y_{pvm} - \overline{y}_{...} = y_{pvm} - \overline{y}_{pv.} + \overline{y}_{pv.} - \overline{y}_{p...} + \overline{y}_{p...} - \overline{y}_{...}$$

and note that

$$\begin{array}{rcl} y_{pvm} - \overline{y}_{pv} & = & \varepsilon_{pvm} - \overline{\varepsilon}_{pv}, \\ \overline{y}_{pv} - \overline{y}_{p..} & = & e_{pv} + \overline{\varepsilon}_{pv} - (\overline{e}_{p.} + \overline{\varepsilon}_{p..}), \\ \overline{y}_{p..} - \overline{y}_{...} & = & b_{p} + \overline{e}_{p.} + \overline{\varepsilon}_{p..} - (\overline{b}. + \overline{e}... + \overline{\varepsilon}...), \end{array}$$

so the overall sum of squares is

$$\sum_{p,v,m} (y_{pvm} - \overline{y}_{...})^{2} = \sum_{p,v,m} (y_{pvm} - \overline{y}_{pv.})^{2} + \sum_{p,v,m} (\overline{y}_{pv.} - \overline{y}_{p..})^{2} + VM \sum_{p} (\overline{y}_{pv.} - \overline{y}_{...})^{2},$$

where these terms are independent sums of squares for variables that are

$$\mathcal{N}(0,\sigma^2)$$
, $\mathcal{N}(0,\sigma_e^2+\sigma^2/M)$, $\mathcal{N}\{0,\sigma_b^2+\sigma_e^2/V+\sigma^2/(VM)\}$.

Regression Methods

Nested model ANOVA II

☐ Hence

$$\begin{split} & \sum_{p,v,m} (y_{pvm} - \overline{y}_{pv\cdot})^2 & \sim & \sigma^2 \chi_{PV(M-1)}^2, \\ & \sum_{p,v,m} (\overline{y}_{pv\cdot} - \overline{y}_{p\cdot\cdot})^2 & \sim & M(\sigma_e^2 + \sigma^2/M) \chi_{P(V-1)}^2 \overset{\mathrm{D}}{=} (M\sigma_e^2 + \sigma^2) \chi_{P(V-1)}^2, \\ & \sum_{p,v,m} (\overline{y}_{p\cdot\cdot} - \overline{y}_{\cdot\cdot\cdot})^2 & \sim & VM\left(\sigma_b^2 + \frac{\sigma_e^2}{V} + \frac{\sigma^2}{VM}\right) \chi_{P-1}^2 \overset{\mathrm{D}}{=} (VM\sigma_b^2 + M\sigma_e^2 + \sigma^2) \chi_{P-1}^2, \end{split}$$

and we can estimate the components of variance σ^2 , σ_e^2 and σ_b^2 from the ANOVA table.

 \Box The interpretation of the ANOVA depends on whether we regard $\delta_b^2=\sum_p(b_p-\overline{b}.)^2$ and $\delta_e^2=\sum_{p,v}(e_{pv}-\overline{e}_{p\cdot})^2$ as random or fixed:

Term	df	Sum of squares	$\mathrm{E}(Mean\ square)$ when terms below random				
			ε	ε, e	arepsilon, e, b		
Between patients	P-1	$\sum (\overline{y}_{p\cdots} - \overline{y}_{\cdots})^2$	$VM\delta_b^2 + M\delta_e^2 \\ + \sigma^2$	$VM\delta_b^2 + M\sigma_e^2 \\ + \sigma^2$	$VM\sigma_b^2 + M\sigma_e^2 + \sigma^2$		
Between visits within patients	P(V-1)	$\sum (\overline{y}_{pv\cdot} - \overline{y}_{p\cdot\cdot})^2$	$M\delta_e^2 + \sigma^2$	$M\sigma_e^2 + \sigma^2$	$M\sigma_e^2 + \sigma^2$		
Between measures within visits	PV(M-1)	$\sum (y_{pvm} - \overline{y}_{pv})^2$	σ^2	σ^2	σ^2		

Regression Methods

Autumn 2024 - slide 269

Nested and crossed ANOVA

☐ Nested analysis of the blood pressure data:

☐ Likewise, crossed analysis of the chimpanzee data:

There are C-1 degrees of freedom for chimps, W-1 for words, and (C-1)(W-1) for the residual.

 \square In both cases, we can use the ANOVA table to estimate the variance components and then perform synthesis of variance: e.g., how large would W need to be to distinguish the learning abilities of two chimps with probability 0.95?

Regression Methods

Example: Blood pressure

 \square Solving the equations

$$\sigma^2 = 7.7$$
, $M\sigma_e^2 + \sigma^2 = 104.2$, $VM\sigma_b^2 + M\sigma_e^2 + \sigma^2 = 960.8$,

gives (in units of millimeters of mercury, mmHg)

$$\widehat{\sigma} = 2.8, \quad \widehat{\sigma}_e = 6.9, \quad \widehat{\sigma}_b = 5.2,$$

so the largest variation is between different visits within patients, while that between measurements on a single visit is smallest.

- ☐ Different comparisons require appropriate baseline variances:
 - if we are interested in how patient p's response varies from visit to visit, we use

$$\overline{y}_{pv_1} - \overline{y}_{pv_2} = \mu + b_p + e_{pv_1} + \overline{\varepsilon}_{pv_1} - (\mu + b_p + e_{pv_2} + \overline{\varepsilon}_{pv_2}) \sim \mathcal{N}(0, 2\sigma_e^2 + 2\sigma^2/M),$$

as a basis for a test of a significant difference, whereas to compare average blood pressures for two different patients we use

$$\overline{y}_{p_1..} - \overline{y}_{p_2..} = b_{p_1} + \overline{e}_{p_1}. + \overline{e}_{p_1}. + \overline{e}_{p_2}.. - (b_{p_2} + \overline{e}_{p_2}.. + \overline{e}_{p_2}..) \sim \mathcal{N}\{0, 2\sigma_b^2 + 2\sigma_e^2/V + 2\sigma^2/(VM)\}.$$

Split-unit designs are set up to make the most important comparisons within units (here patients) and less important ones between units, and the ANOVA reflects this.

Regression Methods

Autumn 2024 - slide 271

General form

☐ We could have written the nested model above as

$$y = 1_n \mu + X_b b + X_e e + \varepsilon,$$

with design matrices X_b and X_e for the patient and visit-within-patient effects.

- ☐ Then if
 - b and e are treated as fixed (ordinary parameters),

$$y \sim \mathcal{N}_n(1_n\mu + X_bb + X_ee, \sigma^2 I_n),$$

- b is treated as fixed but $e \sim \mathcal{N}_{PV}(0, \sigma_e^2 I_{PV})$, then

$$y \sim \mathcal{N}_n(1_n \mu + X_b b, \sigma_e^2 X_e X_e^{\mathrm{T}} + \sigma^2 I_n),$$

– and if $b \sim \mathcal{N}_P(0, \sigma_b^2 I_P)$ independent of $e \sim \mathcal{N}_{PV}(0, \sigma_e^2 I_{PV})$, then

$$y \sim \mathcal{N}_n(1_n \mu, \sigma_b^2 X_b X_b^{\mathrm{T}} + \sigma_e^2 X_e X_e^{\mathrm{T}} + \sigma^2 I_n).$$

 \Box Hence random e or b give patterned covariance matrices depending on their variances.

Regression Methods

Summary

- ☐ Components of variance ANOVA is easily performed directly for balanced data.
- ☐ Standard ANOVA tables have different interpretations, depending on which components of variance are taken to be random or fixed.
- Extensions are needed to deal with more complex settings, with unbalanced data, or with non-linear or non-normal errors hence **mixed models**, i.e., models with both random and fixed parts, arising in many different settings (and with different names):
 - components of variance (as above),
 - classical experimental design (split-plot designs, ...),
 - repeated measures,
 - longitudinal models,
 - multi-level models,
 - hierarchical models.
- ☐ Can subsume linear versions into the **linear mixed model**, which can be extended to nonlinear models, GLMs, . . .

Regression Methods

Autumn 2024 - slide 273

3.7 Linear Mixed Model

slide 274

Linear mixed model

☐ The linear mixed model may be written as

$$y_{n\times 1} = X_{n\times p}\beta_{p\times 1} + Z_{n\times q}b_{q\times 1} + \varepsilon_{n\times 1}, \quad b \sim N_q(0,\Omega_b), \quad \varepsilon \sim \mathcal{N}_n(0,\Omega),$$

where

- β represents the **fixed effects**,
- b represents the random effects, and
- usually $\Omega = \sigma^2 I_n$.
- \Box This has the same structure as when smoothing, with the columns of Z giving the structure of the random effects.
- ☐ Equivalently,

$$y \mid b \sim \mathcal{N}_n(X\beta + Zb, \Omega), \quad b \sim \mathcal{N}_a(0, \Omega_b),$$

which gives marginal response distribution

$$y \sim \mathcal{N}_n(X\beta, Z\Omega_b Z^{\mathrm{T}} + \Omega), \quad Z\Omega_b Z^{\mathrm{T}} + \Omega = \sigma^2 \Delta^{-1}(\psi),$$

say, with ψ the vector of distinct variance ratios appearing in Δ^{-1} (e.g., $\sigma_b^2/\sigma^2,\ldots$).

 \square Although Ω is often diagonal, $Z\Omega_bZ^{\mathrm{T}}$ is not, so inverting $Z\Omega_bZ^{\mathrm{T}} + \Omega$ involves $O(n^3)$ flops in general, and we should avoid working with Δ .

Regression Methods

Maximum likelihood estimation

 \square Let \tilde{b} denote the MLE of b for fixed β (and ψ). Then

$$\begin{split} f(y;\beta,\sigma^{2},\psi) &= \int f(y\mid b;\beta,\sigma^{2},\psi) f(b;\sigma^{2},\psi) \,\mathrm{d}b \\ &= f(y,\tilde{b};\beta,\sigma^{2},\psi) \times \frac{(2\pi)^{q/2}}{|Z^{\mathrm{T}}\Omega^{-1}Z + \Omega_{b}^{-1}|^{1/2}} \\ &\propto \frac{f(y\mid \tilde{b};\beta,\sigma^{2},\psi) f(\tilde{b}\mid \sigma^{2},\psi)}{|Z^{\mathrm{T}}\Omega^{-1}Z + \Omega_{b}^{-1}|^{1/2}}, \end{split}$$

so (apart from additive constants) $-2\log f(y;\beta,\sigma^2,\psi)$ equals

$$(y - X\beta - Z\tilde{b})^{\mathrm{T}}\Omega^{-1}(y - X\beta - Z\tilde{b}) + \tilde{b}^{\mathrm{T}}\Omega_{b}^{-1}\tilde{b} + \log\{|\Omega||\Omega_{b}||Z^{\mathrm{T}}\Omega^{-1}Z + \Omega_{b}^{-1}|\}.$$

- The first two (quadratic) terms here depend on β and b, so given ψ and σ^2 we can find $\widehat{\beta}_{\psi}$ and $\widetilde{b}(\widehat{\beta},\psi)$ explicitly, and thus obtain $\ell_{\mathrm{p}}(\psi)$.
- \square By noting that

$$f(b \mid y; \beta, \sigma^2, \psi) = f(y \mid b; \beta, \sigma^2, \psi) f(b; \sigma^2, \psi) / f(y; \beta, \sigma^2, \psi)$$

and taking logs, we obtain

$$b \mid y \sim \mathcal{N}_q \left\{ \tilde{b}, (Z^{\mathrm{T}} \Omega^{-1} Z + \Omega_b^{-1})^{-1} \right\}, \quad \tilde{b} = \left(Z^{\mathrm{T}} \Omega^{-1} Z + \Omega_b^{-1} \right)^{-1} Z^{\mathrm{T}} \Omega^{-1} \left(y - X \beta \right).$$

Regression Methods

Autumn 2024 - slide 276

Note on maximum likelihood estimation

 \Box Suppressing the parameters β , σ^2 and ψ for now, we write the log integrand in

$$f(y) = \int f(y,b) db = \int f(y \mid b) f(b) db$$

in the form

$$\log f(y,b) = \log f(y,\tilde{b}) - \frac{1}{2}(b-\tilde{b})^{\mathrm{T}}H(\tilde{b})(b-\tilde{b}),$$

where the linear term of the Taylor series equals zero, because it is evaluated at the maximising value \tilde{b} , and the given Taylor series is exact because the log likelihood is quadratic.

 \Box On ignoring terms not involving b we have

$$-2\log f(y,b) = -2\log f(y \mid b) - 2\log f(b) \equiv (y - X\beta - Zb)^{\mathrm{T}}\Omega^{-1}(y - X\beta - Zb) + b^{\mathrm{T}}\Omega_b^{-1}b,$$

SO

$$H(b) \equiv H = Z^{\mathrm{T}} \Omega^{-1} Z + \Omega_b^{-1}$$

does not depend on b, and thus

$$f(y) = f(y, \tilde{b}) \int \exp\left\{-\frac{1}{2}(b - \tilde{b})^{\mathrm{T}}H(b - \tilde{b})\right\} db$$
$$= f(y, \tilde{b}) \times (2\pi)^{q/2}|H|^{-1/2} = f(y, \tilde{b}) \times \frac{(2\pi)^{q/2}}{|Z^{\mathrm{T}}\Omega^{-1}Z + \Omega_b^{-1}|^{1/2}},$$

as announced; the integral equals the normalising constant for a $\mathcal{N}_q(ilde{b},H^{-1})$ density.

Regression Methods

Inference on β

 \square Since

$$y \sim \mathcal{N}_n(X\beta, Z\Omega_b Z^{\mathrm{T}} + \Omega),$$

weighted least squares gives

$$\widehat{\beta} = \{ X^{\mathrm{T}} (Z\Omega_b Z^{\mathrm{T}} + \Omega)^{-1} X \}^{-1} X^{\mathrm{T}} (Z\Omega_b Z^{\mathrm{T}} + \Omega)^{-1} y,$$

with

$$\widehat{\beta} \sim \mathcal{N}_p \left[\beta, \{ X^{\mathrm{T}} (Z\Omega_b Z^{\mathrm{T}} + \Omega)^{-1} X \}^{-1} \right],$$

where in general we need $O(n^3)$ flops to invert the $n \times n$ matrix $Z\Omega_b Z^{\mathrm{T}} + \Omega$.

 \square For cheaper calculation of $var(\widehat{\beta})$, we use the inversion formulae and obtain

$$\begin{pmatrix} \operatorname{var}(\widehat{\beta})_{p \times p} & \cdot \\ \cdot & \cdot \end{pmatrix} = \begin{pmatrix} X^{\mathsf{T}} \Omega^{-1} X & X^{\mathsf{T}} \Omega^{-1} Z \\ Z^{\mathsf{T}} \Omega^{-1} X & Z^{\mathsf{T}} \Omega^{-1} Z + \Omega_b^{-1} \end{pmatrix}_{d \times d}^{-1},$$

where d=p+q, which involves only $O\{nd^2\}$ flops, as Ω is usually diagonal.

- $\hfill \square$ Note that $\mathrm{var}(b\mid y)=(Z^{\scriptscriptstyle \mathrm{T}}\Omega^{-1}Z+\Omega_b^{-1})^{-1}$ can be obtained as a by-product.
- \square In practice these formulae are evaluated at the MLEs $\widehat{\sigma}^2$ and $\widehat{\psi}$ and used to compute confidence intervals etc. for elements of β .

Regression Methods

Autumn 2024 - slide 277

Inference on random effects

- \square Conventional terminology: we estimate parameters β and predict random variables b.
- \square To find the best predictor $\tilde{b}(y)$ of b we minimise

$$\mathrm{E}_{b,y}\left[\left\{\tilde{b}(y)-b\right\}^{\mathrm{T}}\left\{\tilde{b}(y)-b\right\}\right],$$

which gives $\tilde{b}(y) = E(b \mid y)$, with (Woodbury formula):

$$E(b \mid y) = (Z^{T}\Omega^{-1}Z + \Omega_{b}^{-1})^{-1}Z^{T}\Omega^{-1}(y - X\beta),$$

$$var(b \mid y) = (Z^{T}\Omega^{-1}Z + \Omega_{b}^{-1})^{-1}.$$

- \square Replace parameters β , σ^2 , ψ by estimates to get **best linear unbiased predictor (BLUP)** \tilde{b} and its estimated variance.
- ☐ Residuals

$$y - X\widehat{\beta} = Z\widetilde{b} + y - X\widehat{\beta} - Z\widetilde{b}$$

= $Z\widetilde{b} + \left\{ I_n - Z \left(Z^{\mathsf{T}}\widehat{\Omega}^{-1}Z + \widehat{\Omega}_b^{-1} \right)^{-1} Z^{\mathsf{T}}\widehat{\Omega}^{-1} \right\} \left(y - X\widehat{\beta} \right),$

split into two parts, with $Z\tilde{b}$ attributable to random effects, and the second the usual residual $y-X\widehat{\beta}$ shrunk towards zero; this estimates $\varepsilon.$

Regression Methods

Note on conditional mean and variance

☐ First we write

$$\tilde{b}(y) - b = \tilde{b}(y) - \mathcal{E}(b \mid y) + \mathcal{E}(b \mid y) - b,$$

expand $\{\tilde{b}(y)-b\}^{\mathrm{T}}\{\tilde{b}(y)-b\}$ and take expectation over b conditional on y to get

$$\mathbf{E}\left[\left\{\tilde{b}(y) - b\right\}^{\mathrm{T}}\left\{\tilde{b}(y) - b\right\} \mid y\right] = \left\{\tilde{b}(y) - \mathbf{E}(b \mid y)\right\}^{\mathrm{T}}\left\{\tilde{b}(y) - \mathbf{E}(b \mid y)\right\} + \mathbf{var}(b \mid y),$$

which is minimised when $\tilde{b}(y) = \mathrm{E}(b \mid y)$. Any other choice will give a larger expectation when we take E_{y} , so this is optimal.

 \square To obtain $E(b \mid y)$, we note that

$$\begin{pmatrix} y \\ b \end{pmatrix} \sim \mathcal{N}_{n+q} \left\{ \begin{pmatrix} X\beta \\ 0 \end{pmatrix}, \begin{pmatrix} \Omega + Z\Omega_b Z^{\mathrm{T}} & Z\Omega_b \\ \Omega_b Z^{\mathrm{T}} & \Omega_b \end{pmatrix} \right\},\,$$

so using standard formulae for conditional normal distributions, we have

$$E(b \mid y) = \Omega_b Z^{\mathrm{T}} (\Omega + Z \Omega_b Z^{\mathrm{T}})^{-1} (y - X \beta),$$

$$var(b \mid y) = \Omega_b - \Omega_b Z^{\mathrm{T}} (\Omega + Z \Omega_b Z^{\mathrm{T}})^{-1} Z \Omega_b.$$

☐ The Woodbury formula applied to the conditional variance gives

$$var(b | y) = (Z^{T}\Omega^{-1}Z + \Omega_{b}^{-1})^{-1}$$

as required.

 \Box For the conditional mean we apply the Woodbury formula to $(\Omega+Z\Omega_bZ^{\scriptscriptstyle {\rm T}})^{-1}$ and get

$$\begin{split} \mathbf{E}(b \mid y) &= \Omega_b Z^{\mathsf{T}} \left\{ \Omega^{-1} - \Omega^{-1} Z \left(\Omega_b^{-1} + Z^{\mathsf{T}} \Omega^{-1} Z \right)^{-1} Z^{\mathsf{T}} \Omega^{-1} \right\} (y - X \beta) \\ &= \Omega_b \left\{ I_q - Z^{\mathsf{T}} \Omega^{-1} Z \left(\Omega_b^{-1} + Z^{\mathsf{T}} \Omega^{-1} Z \right)^{-1} \right\} Z^{\mathsf{T}} \Omega^{-1} (y - X \beta) \\ &= \Omega_b \left\{ \Omega_b^{-1} \left(\Omega_b^{-1} + Z^{\mathsf{T}} \Omega^{-1} Z \right)^{-1} \right\} Z^{\mathsf{T}} \Omega^{-1} (y - X \beta), \end{split}$$

as required, where we wrote the term in braces in the second line as $I-B(A+B)^{-1}=A(A+B)^{-1}$, with $A=\Omega_b^{-1}$ and $B=Z^{\rm T}\Omega^{-1}Z$

Regression Methods

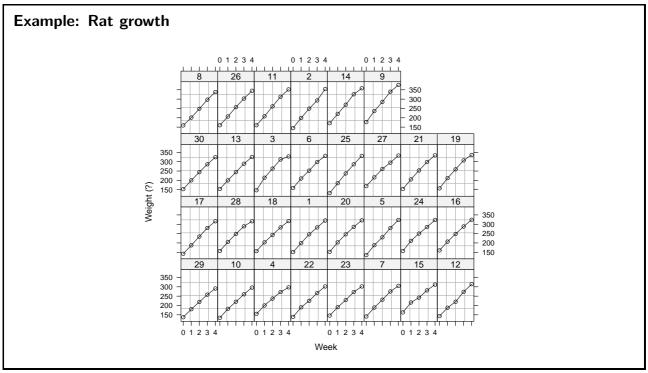
Example: Rat growth

Weights (units unknown) of 30 young rats over a five-week period

			Week				Week							
	1	2	3	4	5		1	2	3	4	5			
1	151	199	246	283	320	16	160	207	248	288	324			
2	145	199	249	293	354	17	142	187	234	280	316			
3	147	214	263	312	328	18	156	203	243	283	317			
4	155	200	237	272	297	19	157	212	259	307	336			
5	135	188	230	280	323	20	152	203	246	286	321			
6	159	210	252	298	331	21	154	205	253	298	334			
7	141	189	231	275	305	22	139	190	225	267	302			
8	159	201	248	297	338	23	146	191	229	272	302			
9	177	236	285	340	376	24	157	211	250	285	323			
10	134	182	220	260	296	25	132	185	237	286	331			
11	160	208	261	313	352	26	160	207	257	303	345			
12	143	188	220	273	314	27	169	216	261	295	333			
13	154	200	244	289	325	28	157	205	248	289	316			
14	171	221	270	326	358	29	137	180	219	258	291			
15	163	216	242	281	312	30	153	200	244	286	324			

Regression Methods

Autumn 2024 - slide 279



Regression Methods

Example: Rat growth

Example 33 (Rat growth data)

□ Write

$$y_{jt} = \beta_0 + b_{j0} + (\beta_1 + b_{j1})x_{jt} + \varepsilon_{jt}, \quad t = 1, \dots 5, j = 1, \dots, 30,$$

where the random variables (b_{j0}, b_{j1}) have a joint normal distribution with mean vector zero and unknown variance matrix and the $\varepsilon_{jt} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$. In matrix terms,

$$\begin{pmatrix} y_{j1} \\ \vdots \\ y_{j5} \end{pmatrix} = \begin{pmatrix} 1 & x_{j1} \\ \vdots & \vdots \\ 1 & x_{j5} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \begin{pmatrix} 1 & x_{j1} \\ \vdots & \vdots \\ 1 & x_{j5} \end{pmatrix} \begin{pmatrix} b_{j0} \\ b_{j1} \end{pmatrix} + \begin{pmatrix} \varepsilon_{j1} \\ \vdots \\ \varepsilon_{j5} \end{pmatrix}, \quad j = 1, \dots, 30;$$

the overall model with n=150 is obtained by stacking these expressions.

- \square We set $(x_{j1},\ldots,x_{j5})=(0,\ldots,4)$, so that β_0 is the mean weight in week 1.
- \square p=2 parameters; q=60 since two random variables per rat.

Regression Methods

```
Example: Rat growth
> rat.growth
   rat week y
    1 0 151
2
         1 199
    1
3
    1 2 246
    1 3 283
    1 4 320
    2 0 145
> fit.reml <- lme(fixed= y~week, random=~week|rat, data=rat.growth)</pre>
> summary(fit.reml)
Linear mixed-effects model fit by REML
Data: rat.growth
     AIC
             BIC logLik
 1096.58 1114.563 -542.2899
Random effects:
Formula: ~week | rat
Structure: General positive-definite, Log-Cholesky parametrization
           StdDev
                  Corr
(Intercept) 10.932986 (Intr)
week
          3.534747 0.184
Residual
           5.817426
Fixed effects: y ~ week
              Value Std.Error DF t-value p-value
(Intercept) 156.05333 2.1589786 119 72.28109
           43.26667 0.7275228 119 59.47122
Correlation:
    (Intr)
week 0.007
```

Example: Rat growth

Results from fit of mixed model to rat growth data, using REML. Values in parentheses are for ML fit. In each case $\hat{\sigma}^2=5.82^2$.

Parameter		Fixed	Random					
	Estimate Standard error		Variance	Correlation				
Intercept	156.05	2.16 (2.13)	$10.93^2 \ (10.71^2)$					
Slope	43.27	0.73 (0.72)	$3.53^2 \ (3.46^2)$	0.18(0.19)				

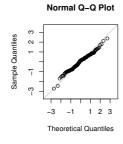
- \square REML estimates of Ω_b slightly larger than ML estimates, but effect is small since p=2.
- \Box Estimated mean weight in week 1 is 156, but SD of individual rats around this is 11.
- ☐ Correlation between slope and intercept is small but positive: initially heavier rats tend to gain weight faster.
- \square Variation around individual slopes is given by $\widehat{\sigma}$, smaller than for the intercept variance.
- ☐ Shrinkage of intercept estimates, shown on next page, is small in this case.
- \square Residuals look acceptably normal.

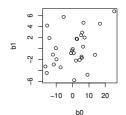
Regression Methods

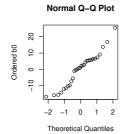
Autumn 2024 - slide 283

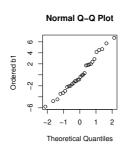


Residuals and random effects









Regression Methods

Comments

Testing for non-zero variance components involves tests on the boundary of the parameter space, which have nasty asymptotic properties: if $\psi=0$, then a likelihood ratio statistic for testing $\psi=0$ satisfies $W \stackrel{.}{\sim} \frac{1}{2}\chi_0^2 + \frac{1}{2}\chi_1^2$ as $n\to\infty$, meaning that

$$P_0(W=0) = \frac{1}{2}, \quad P_0(W>w) = \frac{1}{2}P(\chi_1^2>w), \quad w>0.$$

Unfortunately,

- $P_0(W=0)$ can be very different from $\frac{1}{2}$ even in large samples, and
- in more complex problems, the limiting distribution can be much more complex.
- ☐ Sometimes clearer to write a mixed model in multi-level model form

$$y = X\beta + Z_L b_L + \dots + Z_0 b_0,$$

where the $q_l \times 1$ vectors b_l are all mutually independent with means zero and variance matrices Ω_l , so $Y \sim \mathcal{N}_n(X\beta, \sum_{l=0}^L Z_l \Omega_l Z_l^\mathrm{T})$, where $Z_0 = I_n$, $b_0 = \varepsilon$ and $\Omega_0 = \sigma^2 I_n$.

☐ The same basic approaches apply in **nonlinear mixed models** and **generalized linear mixed models** (**GLMMs**), but integrals appear everywhere and have to be approximated numerically, leading to nastier computations.

Regression Methods

Autumn 2024 - slide 285

3.8 Generalized Additive Models

slide 286

Generalized additive model

☐ Now we write

$$E(y) = \mu$$
, $g(\mu) = \eta = B\theta = X\beta + Zb$,

where

- y follows a GLM (or more general) distribution,
- $g(\cdot)$ is a link function,
- the rest is as before . . .

giving a generalized additive model (GAM).

☐ For a general treatment, suppose we have a penalized log likelihood,

$$\ell_{\lambda}(\theta) = \ell(\theta) - \frac{1}{2}\theta^{\mathrm{T}} S_{\lambda} \theta = \sum_{j=1}^{n} \ell_{j} \{ \eta_{j}(\theta) \} - \frac{1}{2}\theta^{\mathrm{T}} S_{\lambda} \theta,$$

where $\theta_{d\times 1}$ (with d=p+q) contains $\beta_{p\times 1}$ and $b_{q\times 1}$, the latter penalized using a symmetric positive semidefinite $d\times d$ matrix S_{λ} , and the underlying observations y_1,\ldots,y_n giving likelihood contributions ℓ_1,\ldots,ℓ_n are assumed to be independent.

Now we apply the argument leading to the IWLS algorithm to ℓ_{λ} , leading to the **penalized** iterative weighted least squares (PIWLS) algorithm.

Regression Methods

PIWLS

 \Box For fixed λ , we apply (ridge regression) iterative weighted least squares with update step

$$\widehat{\theta}_{\lambda} = (B^{\mathrm{T}}WB + S_{\lambda})^{-1}B^{\mathrm{T}}Wz,$$

where S_{λ} is the penalty matrix, and

 $B_{n \times d} = \partial \eta / \partial \theta^{\mathrm{T}}, \text{ (design matrix)}$

 $W_{n\times n} = \operatorname{diag}(w_1, \dots, w_n), \quad w_j = \{\operatorname{E}(-\partial^2 \ell_j/\partial \eta_j^2)\}, \quad \text{(weights)}$

 $u_{n\times 1} = \partial \ell/\partial \eta$, (score vector),

 $z_{n \times 1} = B\theta + W^{-1}u$, (adjusted dependent variable).

It is easier (but less stable) to use the (random) $-\partial^2\ell_j/\partial\eta_i^2$ in place of $E(-\partial^2\ell_j/\partial\eta_i^2)$.

 \square Thus to obtain (penalized) MLEs $\widehat{\theta}_{\lambda}$ we use the **PIWLS algorithm**:

 \Box fix λ and take an initial $\widehat{\theta}_{\lambda}$. Repeat

- compute η, B, W, u, z ;

- compute new $\widehat{\theta}_{\lambda} = (B^{\mathrm{T}}WB + S_{\lambda})^{-1}B^{\mathrm{T}}Wz$;

until changes in $\ell_{\lambda}(\widehat{\theta}_{\lambda})$ (or $\widehat{\theta}_{\lambda}$, or both) are lower than some tolerance.

 \square We may add a line search: if $\ell_{\lambda}(\widehat{\theta}_{\lambda,\mathrm{new}}) < \ell_{\lambda}(\widehat{\theta}_{\lambda,\mathrm{old}})$, halve the step length and try again.

Regression Methods

Autumn 2024 - slide 288

Note: Derivation of PIWLS algorithm

 \Box To find the estimate $\widehat{\theta}_{\lambda}$ starting from a trial value θ , we make a Taylor series expansion in the score equation

$$0 = \frac{\partial \ell_{\lambda}(\widehat{\theta}_{\lambda})}{\partial \theta} \doteq \frac{\partial \ell_{\lambda}(\theta)}{\partial \theta} + \frac{\partial^{2} \ell_{\lambda}(\theta)}{\partial \theta \partial \theta^{T}} (\widehat{\theta}_{\lambda} - \theta),$$

where

$$\frac{\partial \ell_{\lambda}(\theta)}{\partial \theta} = B^{\mathrm{T}} u(\theta) - S_{\lambda} \theta, \quad \frac{\partial^{2} \ell_{\lambda}(\theta)}{\partial \theta_{r} \partial \theta_{s}} = \sum_{j=1}^{n} \frac{\partial \eta_{j}(\theta)}{\partial \theta_{r}} \frac{\partial^{2} \ell_{j}(\theta)}{\partial \eta_{j}^{2}} \frac{\partial \eta_{j}(\theta)}{\partial \theta_{s}} + \sum_{j=1}^{n} \frac{\partial^{2} \eta_{j}(\theta)}{\partial \theta_{r} \partial \theta_{s}} u_{j}(\theta) + S_{\lambda,r,s},$$

where $B \equiv B(\theta) = \partial \eta / \partial \theta^{\scriptscriptstyle {\rm T}}.$ If we use the approximation

$$-\frac{\partial^{2} \ell_{\lambda}(\theta)}{\partial \theta \partial \theta^{T}} \doteq B^{T} W B + S_{\lambda}, \quad W = \operatorname{diag} \left\{ -\operatorname{E} \left(\partial^{2} \ell_{j} / \partial \eta_{j}^{2} \right) \right\},$$

where the diagonal matrix of second derivatives is replaced by its expectation, then

$$0 \doteq B^{\mathrm{T}}u(\theta) - S_{\lambda}\theta - (B^{\mathrm{T}}WB + S_{\lambda})(\widehat{\theta}_{\lambda} - \theta)$$
$$= B^{\mathrm{T}}u(\theta) + B^{\mathrm{T}}WB\theta - (B^{\mathrm{T}}WB + S_{\lambda})\widehat{\theta}_{\lambda}.$$

If $B^{\mathrm{T}}WB + S_{\lambda}$ is invertible, this gives

$$\widehat{\theta}_{\lambda} \doteq (B^{\mathsf{\scriptscriptstyle T}}WB + S_{\lambda})^{-1}B^{\mathsf{\scriptscriptstyle T}}(u + WB\theta) = (B^{\mathsf{\scriptscriptstyle T}}WB + S_{\lambda})^{-1}B^{\mathsf{\scriptscriptstyle T}}Wz,$$

where $z = B\theta + W^{-1}u$, as required.

Regression Methods

Relation with least squares

 \square With fixed λ , the penalized MLE

$$\widehat{\theta}_{\lambda} = (B^{\mathrm{T}}WB + S_{\lambda})^{-1}B^{\mathrm{T}}Wz$$

results from fixing θ , and then iteratively solving the minimization problem

$$\min_{\theta} \left\| {W^{1/2}z \choose 0}_{(n+d)\times 1} - {W^{1/2}B \choose Q_{\lambda}}_{(n+d)\times d} \theta_{d\times 1} \right\|^2,$$

where Q_{λ} is a matrix square root of S_{λ} , i.e., $Q_{\lambda}^{\mathrm{T}}Q_{\lambda}=S_{\lambda}$.

☐ The corresponding smoothing matrix is taken to be

$$H_{\lambda} = B(B^{\mathrm{T}}WB + S_{\lambda})^{-1}B^{\mathrm{T}}W,$$

and the effective degrees of freedom for a smooth component are defined as the sum of the corresponding diagonal elements of

$$P_{\lambda} = (B^{\mathrm{T}}WB + S_{\lambda})^{-1}B^{\mathrm{T}}WB,$$

with both H_{λ} and P_{λ} evaluated at the final step of the iteration.

Regression Methods

Autumn 2024 - slide 289

Approaches to iteration

- \square Having chosen how to choose λ for fixed θ , there are two main algorithms:
 - **performance iteration** repeat $\{$ fix λ , update θ with one step of PIWLS, update λ $\}$ to convergence;
 - **outer iteration** repeat $\{$ fix λ , iterate PIWLS to convergence, update λ $\}$ to convergence.
- ☐ Performance iteration
 - can be faster,
 - but since the objective function for θ changes at each step, it may not converge—especially in the context of **concurvity** (collinearity for curves ...), when two or more smooth functions are (almost) confounded.
- ☐ Outer iteration
 - is computationally more burdensome,
 - but will converge to a (local) optimum.

Regression Methods

Choice of λ

 \Box The choice of λ can be based on the marginal density of y,

$$f(y; \beta, \lambda) = \int f(y \mid b; \beta) f(b; \lambda) db,$$

which has no closed form in general (but is Gaussian if both fs are Gaussian).

- ☐ Various ways to approximate the integral:
 - quadrature (doesn't work well when $\dim(b)$ is high);
 - simulation (e.g., importance sampling, same problems as quadrature);
 - Laplace approximation;
 - use the EM algorithm to avoid approximating the integral.
- \square We focus on Laplace approximation.

Regression Methods

Autumn 2024 - slide 291

Laplace approximation

Lemma 34 Let h(u) be a smooth convex function defined for $u \in \mathbb{R}^d$, with a minimum at $u = \tilde{u}$, where $\partial h(\tilde{u})/\partial u = 0$ and the matrix of partial derivatives $h_2 \equiv \partial^2 h(\tilde{u})/\partial u \partial u^{\mathrm{T}}$ is positive definite, and let

$$I_n = \int_{\mathbb{R}^d} e^{-nh(u)} \, \mathrm{d}u.$$

Then $I_n = \tilde{I}_n \left\{ 1 + O(n^{-1}) \right\}$, and its Laplace approximation is

$$\tilde{I}_n = \frac{(2\pi)^{d/2}}{|nh_2|^{1/2}} e^{-nh(\tilde{u})}.$$

 $\ \, \Box \quad \text{For marginal density approximation we let } \theta = (\beta_{p \times 1}^{\mathrm{\scriptscriptstyle T}}, b_{q \times 1}^{\mathrm{\scriptscriptstyle T}})^{\mathrm{\scriptscriptstyle T}} \sim \mathcal{N}_d(0, S_{\lambda}^-) \text{, and write}$

$$f(y; \beta, \lambda) = \int f(y; \theta) f(\theta; \lambda) d\theta = \frac{|S_{\lambda}|_{+}^{1/2}}{(2\pi)^{d/2}} \int \exp \{\ell_{\lambda}(\theta)\} d\theta,$$

where β is unpenalised, $|S_{\lambda}|_{+}$ is the product of the non-negative eigenvalues of S_{λ} , and

$$\ell_{\lambda}(\theta) = \ell(\theta) - \frac{1}{2}\theta^{\mathrm{T}}S_{\lambda}\theta = O(n);$$

the assumptions of Lemma 34 should be satisfied by $h(u) \equiv -n^{-1}\ell_{\lambda}(\theta).$

Regression Methods

Note on Lemma 34

 \Box Close to \tilde{u} a Taylor series expansion gives

$$h(u) \doteq h(\tilde{u}) + h'(\tilde{u})^{\mathrm{T}}(u - \tilde{u}) + \frac{1}{2}(u - \tilde{u})^{\mathrm{T}}h''(\tilde{u})(u - \tilde{u}) = h(\tilde{u}) + \frac{1}{2}(u - \tilde{u})^{\mathrm{T}}h_2(u - \tilde{u})$$

so if we set $z=(nh_2)^{1/2}(u-\tilde{u})$ then $u=\tilde{u}+(nh_2)^{1/2}z$, $\mathrm{d}u/\mathrm{d}z=(nh_2)^{-1/2}$, and arguing heuristically (ignoring the third and higher terms),

$$I_n \doteq e^{-nh(\tilde{u})} \int e^{-n(u-\tilde{u})^{\mathrm{T}} h_2(u-\tilde{u})/2} \, \mathrm{d}u$$

$$= e^{-nh(\tilde{u})} \int_{-\infty}^{\infty} e^{-z^2/2} \frac{\mathrm{d}u}{\mathrm{d}z} \, \mathrm{d}z$$

$$= \left(\frac{(2\pi^d)}{|nh_2|}\right)^{1/2} e^{-nh(\tilde{u})},$$

because the d-dimensional normal density has unit integral.

 \square A more detailed accounting is needed to get the error term. Take the scalar case (d=1) for simplicity. We start by writing

$$nh(u) \doteq nh(\tilde{u}) + \frac{1}{2}nh_2(u - \tilde{u})^2 + \frac{1}{6}nh_3(u - \tilde{u})^3 + \frac{1}{24}nh_4(u - \tilde{u})^4 + \cdots$$

$$= nh(\tilde{u}) + \frac{1}{2}z^2 + \frac{1}{6}\frac{h_3/h_2^{3/2}}{n^{1/2}}z^3 + \frac{1}{24}\frac{h_4/h_2^2}{n}z^4 + O(n^{-3/2})$$

$$= nh(\tilde{u}) + \frac{1}{2}z^2 + \frac{A}{n^{1/2}}z^3 + \frac{B}{n}z^4 + O(n^{-3/2})$$

say. Hence

$$\begin{split} e^{-nh(u)} &= e^{-nh(\tilde{u}) - \frac{1}{2}z^2} \left\{ 1 - \frac{A}{n^{1/2}}z^3 - \frac{B}{n}z^4 + \frac{1}{2} \left(-\frac{A}{n^{1/2}}z^3 - \frac{B}{n}z^4 \right)^2 + O(n^{-3/2}) \right\} \\ &= e^{-nh(\tilde{u}) - \frac{1}{2}z^2} \left\{ 1 - \frac{A}{n^{1/2}}z^3 - \frac{B}{n}z^4 + \frac{1}{2}\frac{A^2}{n}z^6 + O(n^{-3/2}) \right\}. \end{split}$$

 \square As the odd moments of the normal density are zero, integration with respect to z leaves only the n^{-1} term and the next remaining term is $O(n^{-2})$. The fourth and sixth moments of the standard normal distribution are respectively 3 and 15, and

$$15A^{2}/2 - 3B = 15(h_{3}/h_{2}^{3/2}/6)^{2}/2 - 3\{h_{4}/(24h_{2})\} = \frac{15h_{3}^{2}}{72h_{2}^{3}} - \frac{h_{4}}{8h_{2}^{2}} = \frac{5h_{3}^{2}}{24h_{2}^{3}} - \frac{h_{4}}{8h_{2}^{2}},$$

as required. The same argument works for m > 1, but it is more of a bloodbath.

Regression Methods

Comments on Laplace approximations

- \square The O(1/n) error is relative, so the approximation is often surprisingly accurate;
- \square since the odd moments of the normal density are all zero, the expansion has only terms whose orders are even powers of $n^{-1/2}$, i.e., n^{-1}, n^{-2}, \ldots ;
- \square \tilde{I}_n involves only h and the hessian matrix h_2 at \tilde{u}_n , so is easily found, numerically if necessary;
- the series is asymptotic, so the partial sums may not converge, and including additional terms may not be useful;
- \square as most of the normal probability lies within ± 3 standard deviations of the mean, the limits of the integral are almost irrelevant provided they are far enough away from \tilde{u} ;

□ if

$$I_n = \int_{-\infty}^{\infty} e^{-nh(u)} du, \quad J_n = \int_{-\infty}^{\infty} e^{-nh^*(u)} du,$$

where $h^*(u) = h(u) + O(n^{-1})$, then

$$(I_n/J_n) \div (\tilde{I}_n/\tilde{J}_n) = 1 + O(n^{-2}),$$

so two Laplace approximations can be better than one.

Regression Methods

Autumn 2024 - slide 293

Approximate REML

☐ Laplace approximation gives the approximate restricted log likelihood

$$\ell_{\mathbf{p}}(\lambda) \equiv \frac{1}{2} \log |S_{\lambda}|_{+} - \frac{1}{2} \log |B^{\mathsf{T}}WB^{\mathsf{T}} + S_{\lambda}| + \ell(\widehat{\theta}_{\lambda}) - \frac{1}{2} \widehat{\theta}_{\lambda}^{\mathsf{T}} S_{\lambda} \widehat{\theta}_{\lambda} + O_{p}(n^{-1}),$$

where $O_p(n^{-1})$ is a (random) term of order n^{-1} and

$$\widehat{\theta}_{\lambda} = (B^{\mathrm{T}}WB + S_{\lambda})^{-1}B^{\mathrm{T}}Wz$$

results from iterating PIWLS to convergence for fixed λ and satisfies $\partial \ell_{\lambda}(\widehat{\theta}_{\lambda})/\partial \theta = 0$.

 \Box The expression for $\widehat{\theta}_{\lambda}$ contains

$$B \equiv B(\widehat{\theta}_{\lambda}), \quad W \equiv W(\widehat{\theta}_{\lambda}), \quad z = B(\widehat{\theta}_{\lambda})\widehat{\theta}_{\lambda} + W^{-1}(\widehat{\theta}_{\lambda})u(\widehat{\theta}_{\lambda}),$$

which involve the first two derivatives of the log likelihood contributions ℓ_j .

 \square Newton–Raphson maximization of $\ell_p(\lambda)$ requires its first two derivatives, so we need

$$\frac{\partial \widehat{\theta}_{\lambda}}{\partial \lambda}, \quad \frac{\partial^2 \widehat{\theta}_{\lambda}}{\partial \lambda \partial \lambda^{\mathrm{T}}},$$

which will involve the third and fourth derivatives of the ℓ_j ... could be painful.

A version of this is implemented in mgcv.

Regression Methods

UK monthly AIDS reports 1983–1992											
	Diagnosis period		Reporting-delay interval (quarters):								
Year	Quarter	0^{\dagger}	1	2	3	4	5	6		≥14	to end of 1992
	÷	:	:	:	÷	÷	÷	:	:	:	:
1988	1	31	80	16	9	3	2	8		6	174
	2	26	99	27	9	8	11	3	• • •	3	211
	3	31	95	35	13	18	4	6	• • •	3	224
	4	36	77	20	26	11	3	8	• • •	2	205
1989	1	32	92	32	10	12	19	12		2	224
	2	15	92	14	27	22	21	12		1	219
	3	34	104	29	31	18	8	6			253
	4	38	101	34	18	9	15	6			233
1990	1	31	124	47	24	11	15	8	• • •		281
	2	32	132	36	10	9	7	6	• • •		245
	3	49	107	51	17	15	8	9			260
	4	44	153	41	16	11	6	5			285
1991	1	41	137	29	33	7	11	6			271
	2	56	124	39	14	12	7	10			263
	3	53	175	35	17	13	11	2			306
	4	63	135	24	23	12	1				258
1992	1	71	161	48	25	5					310
	2	95	178	39	6						318
	3	76	181	16							273
	4	67	66								133

Autumn 2024 - slide 295

AIDS data

 \square Chain-ladder model: number of reports in row j and column k is Poisson, with mean

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

but

- why should there be different parameters α_j and β_k for every row and column?
- Wouldn't smooth variation be more plausible?
- \square Better models (maybe?):

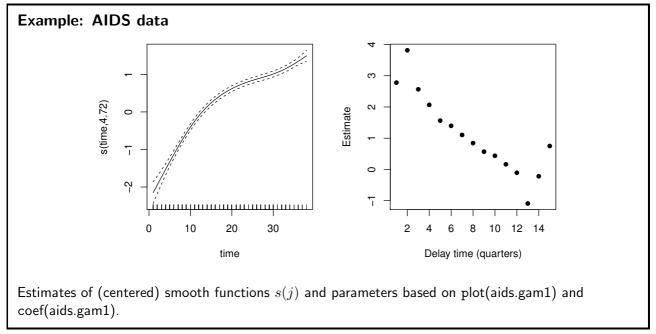
$$\mu_{jk} = \exp\{s(j) + \beta_k\}, \quad \mu_{jk} = \exp\{s(j) + s(k)\},$$

where the time effect s(j) and the delay effect s(k) vary smoothly.

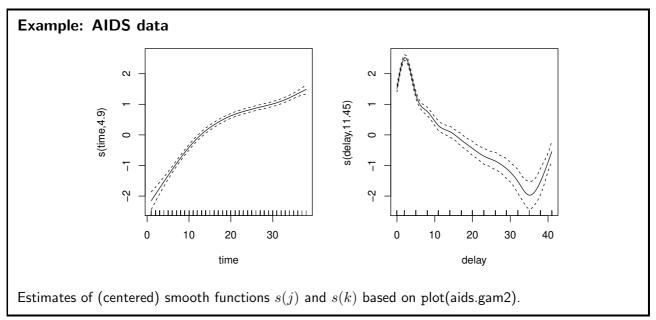
 \square Should also account for the overdispersion . . .

Regression Methods

```
Example: AIDS data
library(mgcv); library(boot)
data(aids)
aids.in <- aids[c(1:570)[as.logical(1-aids$dud)],] # these are elements in the two-way table
aids.glm <- glm(y~factor(time)+factor(delay),family=quasipoisson,data=aids.in)</pre>
aids.gam1 <- mgcv::gam(y~s(time,k=20)+factor(delay)-1,family=quasipoisson,data=aids.in)
plot(aids.gam1,page=1)
> anova(aids.gam1)
Formula:
y \sim s(time, k = 20) + factor(delay)
Parametric Terms:
              df
                      F p-value
factor(delay) 14 261.6 <2e-16
Approximate significance of smooth terms: # Ref.df can be ignored
          edf Ref.df
                          F p-value
s(time) 4.891 6.129 189.1 <2e-16
aids.gam2 <- mgcv::gam(y~s(time,k=20)+s(delay,k=15),family=quasipoisson,data=aids.in)
> anova(aids.gam2)
Formula:
y \sim s(time, k = 20) + s(delay, k = 15)
Approximate significance of smooth terms:
            edf Ref.df
                            F p-value
s(time)
          4.896 6.134 189.0 <2e-16
s(delay) 11.453 12.754 285.5 <2e-16
The fits are very similar, but aids.gam2 has slightly lower AIC of 792.0 compared to 792.1 — these
are so similar that the choice should be based on interpretability rather than on AIC.
```

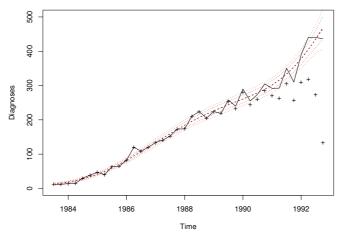


Autumn 2024 – slide 298



Regression Methods

Example: AIDS data



Numbers of recorded deaths (+), with estimated mean deaths per quarter based on chain-ladder model (solid) and on Poisson (black dashes) and quasi-likelihood GAMs with Poisson variance function $V(\mu)=\mu$ (red dashes). The last two estimates have 95% pointwise confidence intervals (dots) based on the fit (treating the smoothing parameters as fixed). To make these I had to compute the fitted means for the missing lower right triangle of the data table.

Regression Methods

Autumn 2024 - slide 300

Closing

- ☐ The basic ideas of regression, dependence of a response on explanatory variables, extend far beyond the linear model, to
 - non-linear dependence on explanatory variables;
 - general response distributions (Poisson, binomial, ...);
 - random effects models—some parameters treated as random, and others as fixed;
 - smooth curve fitting by basis function methods in (generalized) additive models.
- ☐ Unifying themes are:
 - (semi-)parametric modelling using basis functions;
 - maximum likelihood inference;
 - estimation using iterative weighted least squares algorithms;
 - penalized fitting to allow for random effects/basis functions;
 - analysis of deviance;
 - residuals and other diagnostics.

Regression Methods