

# STAT6001

## 2: Linear Models

### 2.0 P

#### 2.1 Inference

General  
OF

2.1.1 Estimation of params

2.1.2 Est  $\sigma^2$

2.1.3 Weighted least squares

2.1.4 CIs & tests

Multiple  
Regression

#### 2.2 Multiple lin reg

2.2.1 Basic results

2.2.2 Interpretation

2.2.3 Model adequacy

Robust  
Regression

#### 2.3 Robust

2.3.0

2.3.1 M-estimation

2.3.2 S-estimator

2.3.3 MM-estimator

Model  
Selection

#### 2.4 Var select

2.4.1 t-tests, F-tests

2.4.2 Best subset

2.4.3 Criteria for models

2.4.4 Stepwise methods

2.4.5 Lasso

#### 2.5 ANOVA

General  
Theory

## 3: Generalized linear models

### 3.0 Exp family

#### 3.0.1 Properties

### 3.2 GLM theory

#### 3.2.1 Estimation

3.2.2 CIs for params

3.2.3 Hypothesis test

Logistic  
Binomial  
Data

### 3.3 Binomial Data, logistic reg

3.3.1 Param interpretation

3.3.2 Basic results

3.3.3 Model adequacy

### 3.3.4 Model selection

3.3.5 Analysis of deviance

### 3.4 Contingency tables

3.4.1 Background

3.4.2 3-way

3.4.3 Fitting models

3.4.4 Goodness of fit, checks

### 3.5 GAMs

3.5.1 Intro

3.5.2 Structure

3.5.3 Param est

3.5.4 Inference

Assumptions: Linearity  
 Normality (of  $e_i$ )  
 Homogeneity ( $e_i \perp\!\!\!\perp x$ )  
 Independence (of  $e_i$ )

## Linear Models

$$\cdot \underline{y} = \underline{X} \underline{\beta} + \underline{e} \quad \left( \begin{array}{l} e_i \stackrel{iid}{\sim} N(0, \sigma^2) \\ y_i \stackrel{iid}{\sim} N(\mu_i, \sigma^2), \mu_i = \underline{x}_i^T \underline{\beta} \end{array} \right)$$

$$\cdot S(\underline{\beta}) = \sum e_i^2 = \sum (y_i - \underline{x}_i^T \underline{\beta})^2 = (\underline{y} - \underline{X} \underline{\beta})^T (\underline{y} - \underline{X} \underline{\beta})$$

*Unbiased if  $y_i$  normal*

$$\frac{\partial S}{\partial \underline{\beta}} = 2 \underline{X}^T \underline{X} \underline{\beta} - 2 \underline{X}^T \underline{y} = 0 \Rightarrow \hat{\underline{\beta}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

$$\Rightarrow \hat{\underline{\beta}} \sim N(\underline{\beta}, \sigma^2 (\underline{X}^T \underline{X})^{-1})$$

$$\cdot \hat{\underline{y}} = \underline{X} \hat{\underline{\beta}} = \underline{H} \underline{y} \Rightarrow \underline{H} = \underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T \quad \left( \begin{array}{l} \text{tr}(\underline{H}) = p, \underline{H}^T \underline{H} = \underline{H} \\ \underline{H}^T = \underline{H} \end{array} \right)$$

$$\cdot \hat{\underline{e}} = \underline{y} - \hat{\underline{y}}, \quad \left( \begin{array}{l} \underline{X}^T \hat{\underline{e}} = 0 \quad (\underline{X}^T \underline{X} \underline{\beta} - \underline{X}^T \underline{y} = 0) \\ \underline{e}^T \underline{e} = 0 \quad (\underline{e}^T \underline{H} \underline{e} = 0) \end{array} \right)$$

*% Variance explained, low = don't trust model inferences*

$$\rightarrow \cdot RSS = \hat{\underline{e}}^T \hat{\underline{e}}, \quad E(\hat{\underline{e}}^T \hat{\underline{e}}) = E(\underline{y}^T (\underline{I} - \underline{H}) \underline{y}) = (N - p) \sigma^2$$

$$\Rightarrow \hat{\sigma}^2 = \frac{RSS}{N - p}, \quad \left( \frac{RSS}{\sigma^2} \sim \chi^2_{N-p} \right)$$

• WLS:  $e_i \sim N(0, V_{ii})$ ,  $V$  diagonal

$$S(\underline{\beta}) = \sum w_i (y_i - \hat{y}_i), \quad w_i = \frac{1}{\text{Var}(y_i)}$$

$$= (\underline{y} - \underline{X} \underline{\beta})^T \underline{V}^{-1} (\underline{y} - \underline{X} \underline{\beta})$$

$$\Rightarrow \hat{\underline{\beta}} = (\underline{X}^T \underline{V}^{-1} \underline{X})^{-1} \underline{X}^T \underline{V}^{-1} \underline{y}$$

$$\left( \begin{array}{l} E(\underline{A} \underline{x}) = \underline{A} E(\underline{x}) \\ V(\underline{A} \underline{x}) = \underline{A} V(\underline{x}) \underline{A}^T \end{array} \right)$$

$$\cdot R^2 = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{RSS}{CTSS}$$

$NSY_0 > 2$

$$\cdot r_i = \frac{\hat{e}_i}{\hat{\sigma} \sqrt{1 - h_{ii}}}$$

$$(\hat{e}_i \sim N(0, \sigma^2 (\underline{I} - \underline{H}))) \leftarrow (\Rightarrow \hat{e}_i \perp\!\!\!\perp \hat{e}_j, \text{ as } \sum e_i = 0)$$

$$\begin{aligned} V(\hat{\underline{e}}) &= V(\underline{y} - \hat{\underline{y}}) = V((\underline{I} - \underline{H}) \underline{y}) \\ &= (\underline{I} - \underline{H}) V(\underline{y}) (\underline{I} - \underline{H})^T \\ &= \sigma^2 (\underline{I} - \underline{H}) \end{aligned}$$

### • Regularisation

$$\cdot \hat{\underline{\beta}}_{RR} = (\underline{X}^T \underline{X} + \lambda \underline{I})^{-1} \underline{X}^T \underline{y}, \quad E(\hat{\underline{\beta}}_{RR}) = \underline{X}^T \underline{X} (\lambda \underline{I} + \underline{X}^T \underline{X})^{-1} \underline{\beta}$$

## Partial vs total regression coeff

• t-test:  $\hat{\beta}_j \sim N(\beta_j, \sigma^2 V_j)$ ,  $V_j = [(\underline{X}^T \underline{X})^{-1}]_{jj}$

$$\Rightarrow \frac{\hat{\beta}_j - \beta_j}{\text{se}(\hat{\beta}_j)} \sim t_{N-p} \quad (\text{se}(\hat{\beta}_j) = \sqrt{\hat{\sigma}^2 V_j^{-1}})$$

$$\hat{\beta}_j \pm t_{N-p, \frac{1}{2}\alpha} \text{ confidence interval } (1-\alpha)\%$$

$\{\psi = \rho\}$

• Robust Regression:

(+ve, symmetric  $\rho(0)=0$ )

• M-estimator:

$$S(\beta) = \sum \rho\left(\frac{e_i(\beta)}{\sigma}\right)$$

(bounded  $\psi = \text{bounded}$ )  
reject of large  $e_i$

analogous to WLS,  
except  $\underline{W}$  not  
constant

$$\hat{\beta}_M^{(k+1)} = (\underline{X}^T \underline{W}^{(k)} \underline{X})^{-1} \underline{X}^T \underline{W} \underline{y},$$

iterate to convergence

$$\begin{aligned} \underline{W}_i^{(k)} &= \frac{\psi(u_i^{(k)})}{u_i^{(k)}} \\ u_i^{(k)} &= (y_i - \hat{y}_i^{(k)}) / \hat{\sigma}^{(k)} \\ \hat{y}_i^{(k)} &= \underline{X}_i^T \hat{\beta}^{(k)} \\ \hat{\sigma}^{(k)} &\propto \text{MED}\{\hat{\beta}_i - \text{MED}\{\hat{\beta}_j\}\} \end{aligned}$$

• Variable Selection

• Backwards & forwards elimination

( $p=d+1$ )

$$AIC = -2\hat{\lambda}(\text{model}) + 2p = N \log \frac{RSS}{N} + 2p \quad \left( \hat{\lambda}(\text{model}) = \frac{N}{2} \log \hat{\sigma}^2 - \frac{RSS}{2\hat{\sigma}^2} \right)$$

$$\hat{\sigma}^2 = \frac{RSS}{N}$$

• LOO-CV

• Lasso (standardize  $x_{ij} = \frac{z_{ij} - \bar{z}_{ij}}{s_j}$ )

# Generalized Linear Models

$$\cdot \eta_i = \underline{x}_i^T \underline{\beta}$$

$$\text{link function} \rightarrow g(\mu_i) = \eta_i, \mu_i = E(y_i), y_i \sim \text{ExpGam}(\Theta, \phi)$$

$$\left. \begin{array}{l} y_i \sim N(\mu_i, \sigma^2) \Rightarrow g(\mu) = \mu \\ y_i \sim \text{Poisson}(\mu_i) \Rightarrow g(\mu) = \log(\mu) \\ \text{default for R} \rightarrow y_i \sim \text{Binomial}(n, p) \Rightarrow g(\mu) = \text{logit}(\mu) = \log \frac{\mu}{1-\mu} \end{array} \right\} \text{Canonical links, } g(\mu) = \Theta$$

$$L(\underline{\beta}) = \sum \log P(y_i)$$

$$\frac{\partial L}{\partial \beta_j} = \sum \left( \frac{y_i - \mu_i}{V(\mu_i)} \right) \frac{d\mu_i}{d\eta_i} x_{i,j} = 0$$

Estimat

Fisher

$$\underline{\beta} \rightarrow \text{iterations} \rightarrow \Rightarrow (\underline{\underline{X}}^T \underline{\underline{W}} \underline{\underline{X}})^{(s-1)} \hat{\underline{\beta}}^{(s)} = (\underline{\underline{X}}^T \underline{\underline{W}} \underline{\underline{Z}})^{(s-1)}$$

$$W_{ii} = \frac{1}{V(\mu_i^{(s)})} \left( \frac{d\mu_i}{d\eta_i} \right)^2$$

$$Z_i^{(s)} = \eta_i + (y_i - \mu_i^{(s)}) \left( \frac{d\eta_i}{d\mu_i^{(s)}} \right)$$

$$\text{Fisher Info Matrix} \rightarrow \underline{\underline{\Sigma}}_{\beta_j} = E \left( \frac{\partial L}{\partial \beta_j} \frac{\partial L}{\partial \beta_j} \right), \underline{\underline{\Sigma}} = \underline{\underline{X}}^T \underline{\underline{W}} \underline{\underline{X}} / \phi \quad (= \text{negative of Hessian})$$

$$\hat{\underline{\beta}} \sim N(\underline{\beta}, \underline{\underline{\Sigma}}^{-1}) \Rightarrow \frac{\hat{\beta}_j - \beta_j}{\text{se}(\hat{\beta}_j)} \sim N(0, 1)$$

Deviance:

$$D^* = 2 \left[ L(\hat{\beta}_{\text{sat}}) - L(\hat{\beta}) \right] \phi \quad \left( L(\hat{\beta}_{\text{sat}}) \text{ is likelihood when } \hat{\mu}_i = y_i \right)$$

$$D^* = D / \phi \sim \chi^2_{N-p}$$

$$\text{likelihood test of } H_0 \rightarrow D_o^* - D^* \sim \chi^2_{p-q}, \quad (D_o^* = \text{standardized deviance under } H_0: p-q \text{ all 0})$$

$$D_{M_1}^* - D_{M_2}^* \sim \chi^2_{p_1 - p_2}$$

$$\text{Estimat. } \phi \rightarrow \chi^2 = \sum \frac{(y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)}, \quad \chi^2 / \phi \sim \chi^2_{N-p} \Rightarrow \hat{\phi} = \frac{\hat{\chi}^2}{N-p}$$

D

## Exponential Family

$$g(y; \theta, \phi) = \exp \left[ \frac{y\theta - b(\theta)}{\phi} + c(y, \phi) \right]$$

	$\theta$	$\phi$	$b(\theta)$	$c(y, \phi)$
Poisson ( $\mu$ )	$\log \mu$	1	$e^\theta$	$-\log y!$
Bin ( $n, \pi$ )	$\logit \pi$	1	$n \log(1+e^\theta)$	$\log \binom{n}{y}$
$N(\mu, \sigma^2)$	$\mu$	$\sigma^2$	$\frac{1}{2}\theta^2$	$-\frac{1}{2}(y/\phi + \log(2\pi\phi))$

For GLM, Canonical link if  $g(\mu) = \theta$

$$E(Y) = b'(\theta), \quad \text{Var}(Y) = b''(\theta) \phi = V(\mu) \phi$$

## Generalized Additive Models

$$\eta = \mathbf{x}^T \theta + \sum_j g_j(x_j) \quad (\text{instead of } \eta = \mathbf{x}^T \beta \text{ for GLM})$$

$g_j$  are smooth functions, eg splines

Spline: weighted sum of basis functions,  $g_j(x_j) = \sum_k \beta_{j,k} b_{j,k}(x_j)$

- $\lambda$  = smoothing parameter; punish excessively wiggly fits

- estimate parameters w/ fixed  $\lambda$  by maximum likelihood

## Linear Models

$$Y_i = \vec{x}_i^T \vec{\beta} + e_i, \quad i=1 \dots N, \quad \vec{x}_i \in \mathbb{D}^M$$

- often assume  $e_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$  ( $y_i \stackrel{\text{iid}}{\sim} N(\vec{x}_i^T \vec{\beta}, \sigma^2)$ )  $D=R \Rightarrow$  multiple linear regression
- components of  $\vec{x}_i$  can be transforms of the explanatory variables, eg

$$\vec{x}_i = \begin{pmatrix} 1 \\ u_i^2 \\ \log(u_i) \end{pmatrix}$$

$D=\text{discrete} \Rightarrow$  analysis of variance (ANOVA)  
 $D=\text{mix} \Rightarrow$  general linear regression

useful if model assumptions (linearity, normality) violated for original variables

or multiplicative model  $Z_i = u_{i1} u_{i2} \dots u_{im} V_i$ , take logs  $\Rightarrow$   $\log Z_i = \log u_{i1} + \log u_{i2} + \dots + \log u_{im} + e_i$

or variable interactions,  $\vec{x}_i = (1, u_i, v_i, u_i v_i)^T$

or comparison of group means

$$Y_{ij} = \mu_i + e_{ij} \quad i=1 \dots I \quad \# \text{ groups}$$

$$= x_{ijk} \mu_1 + x_{ijk} \mu_2 + \dots + e_{ij} \quad j=1 \dots N_i$$

$$Y_{ij} = \vec{x}_{ij}^T \vec{\mu} + e_{ij}$$

$$\text{where } x_{ijk} = \begin{cases} 1 & \text{if } i=k \\ 0 & \text{otherwise} \end{cases}, \quad \vec{x}_{ij} \in \{0, 1\}^I$$

### Inference:

Parameter estimation:

Least squares:

$$\hat{\beta} = \underset{\beta}{\operatorname{arg\,min}} S(\vec{\beta}), \quad S(\vec{\beta}) = \sum_i e_i^2 = \sum_i (Y_i - \vec{x}_i^T \vec{\beta})^2$$

$$= (\vec{Y} - \vec{X} \vec{\beta})^T (\vec{Y} - \vec{X} \vec{\beta})$$

$$\Rightarrow \vec{X}^T \vec{X} \hat{\beta} = \vec{X}^T \vec{Y}$$

if  $p < N$ ,  $\vec{X}^T \vec{X}$  invertible  $\Rightarrow \hat{\beta} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{Y}$

$\cdot E(\hat{\beta}) = \vec{\beta}$ ,  $V(\hat{\beta}) = \sigma^2 (\vec{X}^T \vec{X})^{-1}$  (if  $e_i$  independent w/ var  $\sigma^2$ )

$$\hat{\beta} \sim N(\vec{\beta}, \sigma^2 (\vec{X}^T \vec{X})^{-1}) \text{ if } e_i \sim N(0, \sigma^2) \text{ iid}$$

## Maximum likelihood estimation

$$L(\vec{\beta}, \sigma^2) = \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^N \exp \left[ -\frac{S(\vec{\beta})}{2\sigma^2} \right]$$

max of  $L(\vec{\beta}, \sigma^2)$  for given  $\sigma$  equivalent to  $\min S(\vec{\beta})$

### Estimation of $\sigma^2$

#### Residuals

Residual  
sum of  
squares

$$\hat{e} = \vec{Y} - X\hat{\beta}, \quad \hat{\beta} = (X^T X)^{-1} X^T \vec{Y} \\ = (I_N - H) \vec{Y}, \quad H = X(X^T X)^{-1} X^T$$

$$\rightarrow \text{RSS} = \hat{e}^T \hat{e}$$

$$E(\text{RSS}) = E(\hat{e}^T \hat{e}) = E(\vec{Y}^T (I_N - H) \vec{Y}) = (N - p)\sigma^2$$

$$\Rightarrow \hat{\sigma}^2 = \frac{\text{RSS}}{N - p} \quad (\text{if } e_i \sim N(0, \sigma^2), \text{ IID} \Rightarrow \frac{\text{RSS}}{\sigma^2} \sim \chi^2_{N-p})$$

diagonal, i.e. independent errors

### Weighted Least Squares Estimation

use when error variance not constant

minimise  $S(\vec{\beta}) = \sum w_i (Y_i - \vec{x}_i^T \vec{\beta})^2$  ( $w_i \uparrow$  more reliable measurements,  $w_i \downarrow$  lower variance)

$$\text{if } w_i = \frac{1}{\text{Var}(x_i)}, \quad S(\vec{\beta}) = (\vec{Y} - X\vec{\beta})^T V^{-1} (\vec{Y} - X\vec{\beta})$$

$$e_i \sim N(0, V)$$

#### correlated errors:

$V$  no longer diagonal

normal equations:  $X^T V^{-1} X \hat{\beta} = X^T V^{-1} Y$

problem: need to estimate  $V$ ; see later

### Tests & Confidence intervals:

#### Regression params:

$$e_i \sim N(0, \sigma^2) \quad \text{element } (j,j) \text{ of } (X^T X)^{-1}$$

$$\hat{\beta}_j \sim N(\beta_j, \sigma^2 V_j)$$

$$\frac{\text{RSS}}{\sigma^2} \sim \chi^2_{N-p}$$

$$\hat{\sigma}^2 = \frac{\text{RSS}}{N-p}$$

$$\Rightarrow t = \frac{\hat{\beta}_j - \beta_j}{\sqrt{\hat{\sigma}^2 V_j}} \sim t_{N-p}$$

$$\text{SE}(\hat{\beta}_j) = \sqrt{\hat{\sigma}^2 V_j}$$

$$100(1-\alpha)\% \text{ CI} = \hat{\beta}_j \pm t_{N-p, 1-\alpha/2} \cdot \text{SE}(\hat{\beta}_j)$$

Test hypothesis of the form  $\beta_j = \beta^*$  for some  $\beta^*$ ,

$$\text{eg } H_0: \beta_j = 0$$

$$T = \frac{\hat{\beta}_j}{\text{se}(\hat{\beta}_j)} \sim F_{N-p} \text{ under } H_0$$

• Linear combination of regression params:

$$\Psi = \vec{C}^T \vec{\beta}, \quad \hat{\Psi} = \vec{C}^T \hat{\beta}$$

$$\hat{\Psi} \sim N(\Psi, \sigma^2 V) \quad \text{where } V = \vec{C}^T (X^T X)^{-1} \vec{C}$$

$$\frac{\hat{\Psi} - \Psi}{\sqrt{\sigma^2 V}} \sim t_{N-p}$$

leaving max  
q non-zero

• Tests about multiple params

• test subset  $p-q$  params are all 0,  
 $H_0$ : at least 1 of  $p-q$  params in  $H_0$  isn't 0

$\text{RSS}_0, \text{RSS}_1$  are RSS under  $H_0, H_1$

$$\frac{\text{RSS}_0 - \text{RSS}_1}{\sigma^2} \sim \chi^2_{p-q} \quad \left( \text{as } \frac{\text{RSS}_0}{\sigma^2} \sim \chi^2_{N-q}, \frac{\text{RSS}_1}{\sigma^2} \sim \chi^2_{N-p} \right)$$

$$\frac{(\text{RSS}_0 - \text{RSS}_1)/(p-q)}{\text{RSS}_1/(N-p)} \sim F_{p-q, N-p} \text{ under } H_0$$

## Linear Models - Multiple Linear Regression

$$Y_i = \vec{x}_i^T \vec{\beta} + \epsilon_i$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y, \quad \hat{\beta} \sim N(\vec{\beta}, \sigma^2 (X^T X)^{-1})$$

$$RSS = (\vec{Y} - X^T \vec{\beta})^T (\vec{Y} - X^T \vec{\beta}) = \vec{Y}^T \vec{Y} - \vec{\beta}^T X^T \vec{Y}$$

$$\sigma^2 = \frac{RSS}{N-m-1}$$

$$H_0: \beta_j = 0 \Rightarrow t = \frac{\hat{\beta}_j}{se(\hat{\beta}_j)} \sim t_{N-m-1}$$

$$H_0: \nu \text{ of parameters are all } 0$$

↓

$$F = \frac{(RSS_0 - RSS)/\nu}{RSS/(N-m-1)} \sim F_{\nu, N-m-1}$$

Interpretation of Regression params:

• Partial regression coefficient:

•  $\beta_j$  in multiple linear regression model measures  
rate of change in  $\vec{x}_j$  with all others fixed  
mean response when with

• Total regression coefficient:

• when only 1 response variable,  $\beta_j$  is total regression  
coeff  
• measure rate of change of mean response with  
 $\vec{x}_j$ , ignoring other explanatory variables

• Partial & total regression coeffs only equal if X-matrix  
columns are orthogonal

- Checking model adequacy:
  - Checking assumptions

X	Y
X <sub>1</sub>	
X <sub>2</sub>	

- Matrix plot

- predictors vs predictors:

check collinearity & leverage points  $\checkmark$  but problematic

- predictors vs response:

check linearity, outliers & homoscedasticity

- Residuals & Standardized results

$\hat{e}_i = (Y_i - \hat{X}_i^T \hat{\beta})$  estimators for  $e_i$

standardized residuals

$$\rightarrow r_i = \frac{\hat{e}_i}{\hat{\sigma} \sqrt{1 - h_{ii}}} \quad (H = X(X^T X)^{-1} X^T)$$

$r_i$  should have mean 0, std dev 1; assess their size ( $\sim 5\%$  larger than 2)

- Residual plots:

autocorrelation

- Predictor vs residuals: non-linearity, heteroscedasticity
- Fitted val vs residuals: same as above
- Observation order vs residuals: (if order is informative)
  - autocorrelation, heteroscedasticity

- Normal probability plots:

should look like a straight line

indicates deviation from normality & outliers

- Remedies for violated assumptions:

- Non-linearity:

transform predictors (eg log, exp...) or non-linear models

- Non-normality: (help skew)

transformation, robust regression (esp help outliers)

- Heteroscedasticity:

Weighted Least Squares, robust regression if only a few outliers

- Dependence (of errors):

possibly time series models

only really affects std devs & CIs

Assumptions:

• linearity

• Normality of  $e_i$  (outliers/skew bad)

• Homogeneity of variances

• variance independent of predictor variables

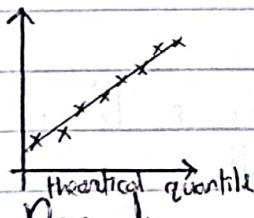
• independence of  $e_i$

not violations of assumptions



residual

sample quantile



straight line = good

- Remedies for violated assumptions:

- Non-linearity:

transform predictors (eg log, exp...) or non-linear models

- Non-normality: (help skew)

transformation, robust regression (esp help outliers)

- Heteroscedasticity:

Weighted Least Squares, robust regression if only a few outliers

- Dependence (of errors):

possibly time series models

only really affects std devs & CIs

- $R^2$ , coefficient of determination

$$R^2 = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{RSS}{CTSS} \quad \text{CTSS} = \sum (y_i - \bar{y})^2$$

- proportion of total variation explained by model
- small = assumption violated, or predictors missing or large error variance

### Outliers:

- Regression outliers: large residuals on plots

- Leverage points: potential of a single point to influence model  
• don't violate model assumptions

- good: is 'in-line' w/ other points, ie omitting this point won't change model much

- bad: is 'out of line'; removing point changes model a lot

standardized measure of leverage  $\rightarrow MD_i = \sqrt{(\vec{x}_i - \bar{\vec{x}})^T \hat{\Sigma}_x^{-1} (\vec{x}_i - \bar{\vec{x}})} ; MD_i^2 = (N-1) \left[ h_{ii} - \frac{1}{N} \right]$

- Cook's Statistic:  $(MD_i^2 \sim \chi^2_{p-1})$

$\hat{\vec{\beta}}_{(i)}$  is estimator of  $\vec{\beta}$ , omitting point  $i$

$$D_i = \frac{1}{p\hat{\sigma}^2} (\hat{\vec{\beta}} - \hat{\vec{\beta}}_{(i)})^T \vec{X}^T \vec{X} (\hat{\vec{\beta}} - \hat{\vec{\beta}}_{(i)})$$

$$\left( = \frac{1}{p} \left( \frac{h_{ii}}{1-h_{ii}} \right) r_i^2 \right) \leftarrow \text{efficient}$$

(residual)

- Problems:

- $\vec{X}^T \vec{X}$  non-invertible; Cook's statistic unstable as  $\vec{\beta}$  unstable

- Masking effect:

- leverage points 'hide' each other; need to use robust estimators, can then look at residual plots to help find outliers

# Linear Models - Robust Regression

- Outliers are common in real world data, often <sup>some</sup> explanatory variables are inevitably missed
  - may want to ignore them, or not (eg finance)
- In LS regression, <sup>Squared</sup> residuals make sensitive to outliers (ie try to minimize  $\hat{\sigma}^2$ )
  - Instead, try to minimize absolute or median residual (ie try to minimize  $\hat{\sigma}$ )

- Effect of Leverage points:
  - in simple linear regression,  $\hat{\beta}_1 = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2} = \sum v_i y_i$
  - largest weight to large  $(x_i - \bar{x})$  points (ie large leverage)
  - $R^2$  analogy  $\Rightarrow$  large leverage  $\uparrow R^2$   
small leverage  $\downarrow R^2$

- Efficiency:
  - ~~design~~ LS is optimal, robust methods designed to be less effected by outliers  
 $\Rightarrow$  asymptotic variance higher than that of LS estimator

- M-estimation:
  - minimizes sum of objective function,  $\rho$ , or residuals  $\sum \rho(r_i)$ , where  $\rho'(u) \geq 0$ ,  $\rho(0) = 0$ ,  $\rho(u) = \rho(-u)$ ,  $\rho(u) \geq \rho(u') \text{ if } |u| \geq |u'|$
  - $\hat{\beta}_M$  minimizes  $\sum \rho\left(\frac{e_i(\hat{\beta})}{\sigma}\right) = \sum \rho\left(\frac{Y_i - \vec{x}_i^T \hat{\beta}}{\sigma}\right)$ 
    - robustness weights
    - $w_i = \frac{\rho'(u_i)}{u_i}$
    - $u_i = (Y_i - \vec{x}_i^T \hat{\beta}_M) / \sigma$

$\sigma$  is unknown: estimate  $w_i$

$\sigma_{MAD} = \text{MED}(\hat{e}_i - \text{MED}(\hat{e}_i), 3/3)$   
(need to iterate this too)

if  $w_i$  fixed,  $\hat{\beta}_M = (X^T W X)^{-1} X^T W Y$  ( $W = \text{diag}(w_1, \dots, w_N)$ )  
as not, iterate

$$\hat{\beta}_M^{(k+1)} = (X^T W^{(k)} X)^{-1} X^T W^{(k)} Y \text{ until convergence}$$

- residuals only enter through  $\rho'$   $\Rightarrow$  bounded  $\rho' = \text{bounded effect of large residuals}$

• Distribution of M-estimators & efficiency:

• Consistent if  $E(\rho'(Z)) = 0$ , where  $Z \sim F$

• holds for all symmetric  $F$  if  $\rho'$  bounded

• Tests & CIs:

$$\sqrt{n}(\hat{\beta}_M - \beta) \xrightarrow{n \rightarrow \infty} N(0, V(\rho', F)L^{-1}), \quad V(\rho', F) \text{ complicated matrix}$$

$$L = \lim \frac{1}{n} X^T X$$

• Efficiency:

$$C = \sigma^2 (X^T X)^{-1} \quad \text{for LS}$$

$$V(\rho, F)L^{-1} = bC \quad \text{for M-estimation}$$

$1/b$  = efficiency, need  $b \propto$  as many observations to match LS estimator precision

• Objective function:

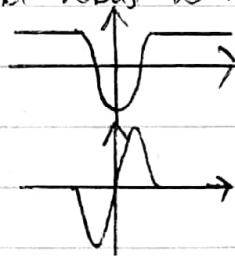
$$\rho(u) = u^2 \Rightarrow LS$$

$$\rho(u) = |u| \Rightarrow L_1\text{-Regression}$$

$\rho$ : not small for large residuals;  
unbounded  
not robust to  $\uparrow$  leverage points

Bisquare

$$\rho_B(u) = \begin{cases} \frac{1}{6}(1 - [1 - (\frac{u}{c})^2]^3) & \text{if } |u| \leq c \\ \frac{1}{6} & \text{if } |u| > c \end{cases}$$



0 for large residuals!

$c$  tunes efficiency vs robustness

• S-estimator:

Bisquare M-estimator relies on (unknown)  $\sigma$

$$\text{M-estimator of scale} \rightarrow M(\vec{z}) = \frac{1}{N} \sum_i \rho\left(\frac{z_i}{S_M(\vec{z})}\right) = E_{N(0,1)} \rho(Z) \quad \text{for observations } (z_1, \dots, z_N)$$

$$\cdot \rho(u) = u^2 \Rightarrow E_{N(0,1)} \rho(Z) = 1 \text{ and } S_M^2(\vec{z}) = \frac{1}{N} \sum_i z_i^2$$

• Linear regression:

$$\begin{aligned} \text{minimise } & S_M(e_1(\vec{\beta}), \dots, e_N(\vec{\beta})) = S_M \\ & = \min S_M\{e_1(\vec{\beta}), \dots, e_N(\vec{\beta})\} \end{aligned}$$

residuals

• Computation:

• find  $\vec{\beta}^{(k+1)}$  for given  $S_M^{(k)}(\vec{z}^{(k)})$  &  $\vec{\beta}^{(k)}$

• find  $S_M^{(k+1)}$  for given  $\vec{z}^{(k+1)} = \{e_1(\vec{\beta}^{(k+1)}), \dots, e_N(\vec{\beta}^{(k+1)})\}$

• repeat until convergence

• maybe local optima; repeat w/ diff  $\vec{\beta}^{(0)}$

- MM-estimator:

- combine S & M-estimation

- $S_M$  estimator of  $\sigma^2$ , then use this for M-estimator

- iterate

## Linear Models - Variable Selection

$$Y_i = \vec{x}_i^T \vec{\beta} + e_i$$

- aim to reduce #no. variables, ie find  $\beta_i = 0$  (or close to 0)
- motivation:
  - + p large, N small  $\Rightarrow X^T X$  nearly singular, unstable estimation
  - $N > 10p$  desired
  - + simpler, clear explanation
  - worse for prediction
  - unstable variable selection if some are correlated
  - bias  $\beta$  remaining to be high, so p-values biased low
- t-tests & F-tests:
  - F-test anova table: no variable affects response
  - F-test:  $\mathcal{Y}$  of  $\beta_i = 0$  (given other  $m-1$  in model)
  - t-test:  $\beta_j = 0$ , given other  $m-1$  in model
  - note: dropping 1 variable changes p-values for others!
- Best Subset Selection:
  - $m$  variables  $\Rightarrow 2^m$  possible models; which is best?
  - for fixed  $k$  no. variables, find model which:
    - minimizes  $RSS$ ,  $\hat{\sigma}^2$ , maximizes  $R^2$
    - cannot use this to compare models w/ different  $k$
- Criteria to compare models:
  - Akaike Information Criteria (AIC)
$$AIC = -2\hat{\lambda}(\text{model}) + 2p$$

$\hat{\lambda}(\text{model})$  is  $\max_{\text{log-likelihood function}}$   
no. regression params (=d+1)

$$\hat{\lambda}(\text{model}) = -\frac{N}{2} \log \hat{\sigma}^2 - \frac{RSS}{2N} + \text{const.}$$
$$\Rightarrow AIC = N \log \frac{RSS}{N} + 2p, \text{ using } \hat{\sigma}_{ML}^2 = \frac{RSS}{N}$$
  - Other functions of  $RSS$ :
    - Mallows  $C_p$  or adjusted  $R^2$ ; in general too 'soft'; doesn't reduce no. vars enough

Leave-one-out  
cross validation



- LOO-CV:

- train on  $N-1$  points, test on holdout point
- repeat  $N$  times & average losses
- computationally demanding;
- + not based on any model assumptions

- Stepwise methods:

- best subset is exponential in no. variables; too slow!

- Backward elimination:

- fit full model

- for  $k=m$  to 0:

- fit ~~mod~~ all models w/  $k-1$  vars ( $k$  diff models)

- find best

- get sequence of 'best' models for all  $k$ , use AIC or LOO-CV to pick best from sequence

- Forward selection

- Start from 0 variables, keep adding variable which produces largest decrease in RSS

- again select from sequence produced w/ AIC or LOO-CV

- can use A/F-tests here to decide when to stop

- backwards/forwards give diff models; often backwards better

- not guaranteed to find best subset model, & very unstable

- Lasso:

- standardise variables (to mean 0, stdv 1)

$$x_{ij} = \frac{z_{ij} - \bar{z}_i}{s_j}, \text{ minimise } S(\vec{\beta}) = \sum e_i^2$$

makes  $\beta_i$   
comparable size

subject to  $\sum_{i=1}^m |\beta_i| \leq t$

- set  $t$  by LOO-CV

- encourages sparse solutions, & equivalent to  $L_1$ -regularization

- + more stable & 'smooth' for large  $m$  & close-to-collinear situations

- biased, worse than stepwise in 'clear-cut' situations

- computation is complicated

# Analysis of Variance

## Linear Models: ANOVA

- $Y_i$  depends on categorical variables  $\in \{0, 1\}$
- ANOVA refers to breakdown of  $\text{CTSS}$  into sum of  $\text{RSS}$  &  $\text{SS}$
- Comparison of Groups - One Way Layout

$$Y_{ij} = \mu_i + e_{ij}, \quad i=1, \dots, I \leftarrow \text{groups}$$

$$j=1, \dots, n_i$$

$$\vec{Y} = \vec{X} \vec{\beta} + \vec{e}$$

$$\vec{Y} = \begin{pmatrix} Y_{11} \\ \vdots \\ Y_{In_I} \end{pmatrix}, \quad \vec{X} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix}, \quad \vec{\beta} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_I \end{pmatrix}, \quad \vec{e} = \begin{pmatrix} e_1 \\ \vdots \\ e_N \end{pmatrix}$$

- cols of  $X = I$  indicator variables, rows =  $n$  observations
- eg  $n_1=2, n_2=3, n_3=1 \Rightarrow X = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$
- this way,  $\vec{\beta} = \vec{\mu}$
- usual hypothesis:  $\mu_1 = \mu_2 = \dots = \mu_I : H_0, H_1: \text{at least 2 differ}$
- test with  $\text{RSS}_0 = \text{CTSS}$

- Two Way Layout:

• observations classified according to 2 factors

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + e_{ijk} \quad i=1, \dots, I, \quad j=1, \dots, J$$

$\uparrow$  overall effect of expected value  
 $\uparrow$  level  $i$  of factor 1  
 $\uparrow$  level  $j$  of factor 2  
 $\uparrow$  interaction effect

$$k=1, \dots, n_{ij}$$

• for  $I=J=n_{ij}=2$

$$Y = \begin{pmatrix} Y_{111} \\ Y_{112} \\ Y_{121} \\ Y_{122} \\ \vdots \\ Y_{222} \end{pmatrix}, \quad X = \begin{pmatrix} 1 & 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \end{pmatrix}, \quad \vec{\beta} = \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \beta_1 \\ \beta_2 \\ \gamma_{11} \\ \gamma_{12} \\ \gamma_{21} \\ \gamma_{22} \end{pmatrix}$$

• but  $X^T X$  not invertible on too many params

• need to introduce constraints, eg

$$\sum \alpha_i = 0, \quad \sum \beta_j = 0 \quad \text{and} \quad \sum \gamma_{ij} = 0, \quad \sum \gamma_{ij} = 0$$

(allow one of each to be determined by all others, effectively allowing us to ~~remove~~ shrink  $X$ )

• these lead to

$$X = \begin{pmatrix} 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & 1 & -1 \\ 1 & 1 & -1 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix}, \beta = \begin{pmatrix} \mu \\ \alpha_2 \\ \beta_2 \\ \gamma_{22} \end{pmatrix} \quad (\text{other params from constraints})$$

• Usual F-tests: ↗ (or AIC/LOO-CV)

no interactions  $\rightarrow$  all  $\gamma_{ij} = 0$

can remove a factor  $\begin{cases} \rightarrow \text{all } \alpha_i = 0 \\ \rightarrow \text{all } \beta_j = 0 \end{cases}$

done in order is ↘

for 'restricted back selection'

# Generalised Linear Models - General Theory

- Components:

- Random

$$Y_1, \dots, Y_N$$

- Systematic

$$\eta_i$$

- Link function

$$g(\mu_i) = \eta_i, \quad \mu_i = E(Y_i) \quad (\text{monotonic \& differentiable})$$

- Linearity:

$$\eta_i = \vec{x}_i^T \vec{\beta}$$

- Special cases:

- Linear Models:  $g(\mu) = \mu$  (identity link)

- Logistic model:

$$g(\mu) = \log \left( \frac{\mu}{1-\mu} \right) \leftarrow n=1 \Rightarrow \text{logit of } \mu$$

- used for Bernoulli trials ( $n=1$ ) data, or  $Y_i \sim \text{Binomial}(n, \pi_i)$

- Log-linear model:

$$Y_i \sim \text{Poisson}(\mu_i) \Rightarrow \log(\mu_i) = \vec{x}_i^T \vec{\beta} \Rightarrow g(\mu) = \log(\mu)$$

- Contingency tables:

- Exponential family:

$$f(y; \Theta, \phi) = \exp \left[ \frac{y\Theta - b(\Theta)}{\alpha(\phi)} + c(y, \phi) \right]$$

dispersion/scale parameter  $\alpha(\phi) = \phi/w$  for known  $w$

Distribution	$\Theta$	$\phi$	$b(\Theta)$	$c(y, \phi)$	
Poisson( $\mu$ )	$\log \mu$	1	$e^\Theta$	$-\log y!$	(for all, $w=1$ )
Bin( $n, \pi$ )	$\log \frac{\pi}{1-\pi}$	1	$n \log(1+e^\Theta)$	$\log(\frac{\pi}{y})$	
$N(\mu, \sigma^2)$	$\mu$	$\sigma^2$	$\frac{1}{2}\Theta^2$	$-\frac{1}{2} \left( \frac{y^2}{\phi} + \log(2\pi\phi) \right)$	

- Known  $\phi \Rightarrow$  exp family

- Unknown  $\phi \Rightarrow$  exp dispersion family

- Properties:

differentiate  
 $\int f(y) dy = 1$ ,  
use  $\alpha(\phi) \neq 0$

- $E(Y) = b'(\Theta)$

- $V_{\text{ar}}(Y) = b''(\Theta) \alpha(\phi)$

- $V_{\text{ar}}(Y) = V(\mu) \alpha(\phi) = V(\mu) \frac{\phi}{w}$

- Canonical links:

- $g$  such that  $g(\mu) = \Theta$

- log, logit & identity for Poisson, Bin & N

	$V(\mu)$
Poisson( $\mu$ )	$\mu$
Bin( $n, \pi$ )	$\mu(n-\mu)/n$
$N(\mu, \sigma^2)$	1

$$f(y_i; \theta_i, \phi) = \exp \left[ \frac{y_i \theta_i - b(\theta_i)}{a_i(\phi)} + c(y_i, \phi) \right]$$

- Some GLM theory:

- $Y_i \sim$  from distribution w/ pdf in exp family (& independent)

- let  $E(Y_i) = \mu_i$ ,  $\eta_i = \sum_j \beta_j x_{ij}$

- Estimation:

$$\lambda = \sum \lambda_i, \quad \lambda_i = \log f(y_i; \theta_i, \phi) [= \log P(y)]$$

$$\Rightarrow \frac{\partial \lambda}{\partial \beta_j} = \sum_i \left( \frac{y_i - \mu_i}{V(\mu_i)} \right) \frac{d\mu_i}{d\eta_i} x_{ij} = 0$$

- Iterative procedure:

- Estimation of  $\hat{\beta}$ :

- cannot solve likelihood eqns algebraically

- Fisher Scoring method:

- initialize  $\hat{\beta}^{(0)}$ ,

- iterate

$$\underbrace{(X^T W X)^{(s-1)}}_{\text{evaluated w/ } \hat{\beta}^{(s-1)}} \hat{\beta}^{(s)} = (X^T W \vec{z})^{s-1}$$

$$w_{ii} = \frac{1}{V(\mu_i)} \left( \frac{d\mu_i}{d\eta_i} \right)^2$$

$$z_i = \eta_i + (y_i - \mu_i) \left( \frac{d\eta_i}{d\mu_i} \right)$$

Pearson statistic

$$X^2 = \sum_i \frac{(y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)} \quad \text{O mean, unit variance RVs}$$

$$\frac{X^2}{\phi} \sim \chi^2_{N-p} \Rightarrow \hat{\phi} = \frac{\hat{X}^2}{N-p}$$

- Sampling distribution of  $\hat{\beta}$ :

Taylor expand  
around true  $\beta_0, \hat{\beta}$   
evaluate at  $\hat{\beta}$

$$\rightarrow \frac{\partial \lambda}{\partial \beta} \Big|_{\hat{\beta}} \approx \frac{\partial \lambda}{\partial \beta} \Big|_{\beta_0} + (\hat{\beta} - \beta_0) \frac{\partial^2 \lambda}{\partial \beta^2} \Big|_{\beta_0}$$

$$\Rightarrow (\hat{\beta} - \beta_0) \approx - \frac{\frac{\partial \lambda}{\partial \beta} \Big|_{\beta_0}}{\frac{\partial^2 \lambda}{\partial \beta^2} \Big|_{\beta_0}} \sim N(0, I^{-1}) \text{ as } n \rightarrow \infty \text{ (by CLT)}$$

$$\Rightarrow \hat{\beta} \sim N_p(\bar{\beta}, I^{-1})$$

$(I = X^T W X / \phi, \text{ for } W \text{ at convergence}$   
of Fisher Scoring method)

# GLM - General Theory

Covariance Matrix of  $\hat{\beta}$ :

Fisher information matrix  
( $j, k$ )<sup>th</sup> element

$$E\left(\frac{-\partial \lambda}{\partial \beta_j} \frac{\partial \lambda}{\partial \beta_k}\right) = \sum_i \sum_{i'} \left( \frac{E[(Y_i - \mu_i)(Y_{i'} - \mu_{i'})]}{\text{Var}(Y_i) \text{Var}(Y_{i'})} \right) \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \mu_{i'}}{\partial \eta_{i'}} x_{ij} x_{i'k}$$

(ie the -Hessian!)

$$I = \sum_i \left( \frac{x_{ij} x_{ik}}{\text{Var}(Y_i)} \right) \left( \frac{\partial \mu_i}{\partial \eta_i} \right)^2 \left( \begin{array}{l} \text{as} \\ E[(Y_i - \mu_i)(Y_{i'} - \mu_{i'})] \\ \text{Var}(Y_i) \text{ is } i=i' \\ \text{cov}(Y_i, Y_{i'})=0 \text{ is } i \neq i' \end{array} \right)$$

$$I = X^T W X / \phi$$

Covariance matrix of  $\hat{\beta}$  is  $I^{-1} = (X^T W X)^{-1} \phi$

Confidence intervals for model parameters:

assume known dispersion parameter,  $\phi$  (unknown  $\phi$  same, except must estimate  $\phi$  & use appropriate t-distribution)

Single parameters:

$$\hat{\beta} \sim \text{approx } N_p(\vec{\beta}, I^{-1})$$

$$\text{approx CI}_{(1-\alpha)Y_i} \rightarrow \hat{\beta}_j \pm z_{\alpha/2} \text{se}(\hat{\beta}_j), \text{ where } \text{se}(\hat{\beta}_j) = \sqrt{(I^{-1})_{jj}}$$

linear combination of parameters:

$$\psi = \vec{c}^T \hat{\beta}$$

$$\psi \sim \text{approx } N(\psi, \vec{c}^T I^{-1} \vec{c})$$

Hypothesis testing:

$$H_0: \beta_j = 0 \Rightarrow \hat{\beta}_j \sim \text{approx } N(0, 1) \text{ under } H_0$$

likelihood test

$H_0$ : subset of  $p-q$  params all 0

$$\Rightarrow 2[\lambda(\hat{\beta}) - \lambda(\hat{\beta}_0)] \sim \chi^2_{p-q} \text{ under } H_0 \quad (\text{applies when } \phi \text{ known})$$

Similar to RSS

Deviance:

$$D = 2[\lambda(\hat{\beta}_{\text{sat}}) - \lambda(\hat{\beta})] \phi \quad \left( \begin{array}{l} \lambda(\hat{\beta}_{\text{sat}}) \text{ is saturated likelihood,} \\ \text{achieved when } \mu_i = y_i \end{array} \right)$$

scaled deviance

$$D^* = D/\phi \sim \chi^2_{N-p}$$

(under  $H_0$  above)

can also reexpress likelihood test

$$D_o^* - D^* \sim \chi^2_{p-q}$$

• unknown  $\phi$ :

• under same  $H_0$ ,  $D_o^* - D^* \sim \chi^2_{p-q}$  &  $D^* \sim \chi^2_{N-p}$

• if  $(D_o^* - D^*) \perp\!\!\!\perp D^*$ , under  $H_0$  & in large sample limit

$$F = \frac{(D_o^* - D^*)/p-q}{D^*/(N-p)} = \frac{(D_o - D)/p-q}{D/(N-p)} \sim F_{p-q, N-p}$$

no need to know  $\phi$

# GLM - Binomial Data & Logistic Regression

$Y_1, \dots, Y_N$  independent,  $Y_i \sim \text{Bin}(n_i, \pi_i)$

$$g(\mu_i) = \vec{x}_i^T \vec{\beta}, \quad \mu_i = E(Y_i) = n_i \pi_i$$

logit link gives  
linear logistic model  
(CDF of logistic dist)  $\pi = \frac{1}{1+e^{-\eta}}$

$$\eta = \vec{x}_i^T \vec{\beta} = \log \left( \frac{\pi_i}{1-\pi_i} \right)$$

$\pi = \Phi(\eta)$

(alternative: probit (inverse standard normal CDF))

complementary log-log ( $\vec{x}_i^T \vec{\beta} = \log [-\log(1-\pi_i)]$ )

(CDF of extreme value dist)  $\pi = 1 - \exp(-e^{\eta})$

• these choices ensure  $0 \leq \pi \leq 1 \quad \& \quad -\infty < g(\mu) < \infty$   
 $(\Rightarrow \text{no constraints on } \vec{\beta})$

• Interpretation of logistic regression parameters:

- logit( $\pi$ ) =  $\log \left( \frac{\pi}{1-\pi} \right)$  = log-odds
- $\beta_j$  measures rate of change of log odds w/  $x_j$   
 $(\& \text{other } x_i \text{ held constant})$
- e.g.  $x_i \in \{0, 1\}, \Rightarrow$  odds on success w/  $x_i = 1$  is  
 odds of success w/  $x_i = 0 \times \exp[\beta_i]$

• Some basic results:

• likelihood:  $\frac{\partial \lambda}{\partial \beta_j} = \sum_i \frac{y_i - \mu_i}{V(\mu_i)} \cdot \frac{\partial \mu_i}{\partial \eta_i} x_{ij} = \sum_i (y_i - \hat{\mu}_i) x_{ij} = 0$

• MLE:

Fisher iterations

$$(X^T W X)^{(s-1)} \hat{\beta}^{(s)} = (X^T W Z)^{(s-1)}$$

• for logistic model,  $W_{ii} = \frac{1}{V_i} \left( \frac{\partial \mu_i}{\partial \eta_i} \right)^2$   $\left( V_i = \text{Var}(Y_i) = V(\mu_i) = n_i \pi_i (1-\pi_i) \right)$

$$W_{ii} = n_i \pi_i (1-\pi_i) \quad \left( \frac{\partial \mu_i}{\partial \eta_i} = n_i \pi_i (1-\pi_i) \right)$$

$$Z_i = \eta_i + (y_i - \mu_i) \left( \frac{d\eta_i}{d\mu_i} \right), \quad \left( \frac{d\eta_i}{d\mu_i} = \frac{1}{n_i \pi_i (1-\pi_i)} \right)$$

• Sampling distribution of  $\hat{\beta}$ :

$$\hat{\beta} \stackrel{\text{approx}}{\sim} N_p(\vec{\beta}, (X^T W X)^{-1}) \quad \text{for large } N$$

$$\text{Deviance: } D = 2 \sum_i \left[ y_i \log \left( \frac{y_i}{\hat{\mu}_i} \right) + (n_i - y_i) \log \left( \frac{n_i - y_i}{n_i - \hat{\mu}_i} \right) \right]$$

$D \sim \chi^2_{N-p}$  if model is true

$H_0: \beta_j = 0$

$$\Rightarrow \frac{\hat{\beta}_j}{\text{se}(\hat{\beta}_j)} \sim N(0, 1)$$

$H_0: \nu \text{ of } \beta_1, \dots, \beta_m \text{ are } 0$

$$\Rightarrow D_0 - D \sim \chi^2_{\nu}$$

alternating  
goodness of  
fit test

Pearson chi-squared statistic:

$$X^2 = \sum_i \frac{(Y_i - n_i \hat{\pi}_i)^2}{n_i \hat{\pi}_i (1 - \hat{\pi}_i)} \sim \chi^2_{N-p} \text{ under model}$$

(can be shown  $D \approx X^2$ )

$$(X^2 = \sum \frac{(O - E)^2}{E})$$

Binary data:

if  $n_i = 1 \forall i$ , results still hold but  $D$  no longer good for assess quality of fit

Can test by grouping observations into  $g \approx 10$  groups according to  $\hat{\pi}_{ij}$ , then calculate  $X^2 \sim \chi^2_{g-2}$

Checking model adequacy:

Plot residuals:

Raw

$$\hat{e}_i = Y_i - n_i \hat{\pi}_i$$

$$\left\{ \begin{array}{l} h_{ii} = H_{ii} \\ H = W^{1/2} X (X^T W X)^{-1} X^T W^{1/2} \end{array} \right.$$

Pearson

$$X_i = \hat{e}_i / \sqrt{n_i \hat{\pi}_i (1 - \hat{\pi}_i)}$$

Standardized Pearson

$$r_{p,i} = x_i / \sqrt{1 - h_{ii}}$$

Deviance

$$\sum d_i^2 = D$$

Standardized Deviance

$$r_{D,i} = d_i / \sqrt{1 - h_{ii}}$$

Cook's Statistic:

$$D_i = \frac{1}{p} \left( \frac{h_{ii}}{1 - h_{ii}} \right) r_{p,i}^2$$

## GLM - Binomial Data & Logistic Regression

- Model selection:

- Stepwise / best subset selection:

- can then use AIC, which for binomial data

$$AIC = D + 2p$$

- Compare model deviances:

$$D_{M_1} - D_{M_2} \sim \chi^2_{p_1 - p_2} \quad (\text{equivalent to likelihood test})$$

- Analysis of deviance:

## GLM - Contingency tables: log-linear modelling

A	B	$B_1 \dots B_J$
$A_1$		$n_{11} \dots n_{1J}$
$A_2$		$n_{21} \dots n_{2J}$

$\pi_{ij} = \text{prob observation in cell } (i,j)$

$M_{ij} = N \pi_{ij} = \text{expected frequency}$

expected row frequency

expected column frequency

$$\text{if } A \perp\!\!\!\perp B, \pi_{ij} = \pi_{i+} \cdot \pi_{+j}, M_{ij} = \frac{\pi_{i+} \cdot \pi_{+j}}{N}$$

$$\cdot H_0: A \perp\!\!\!\perp B \Rightarrow \log M_{ij} = \log \pi_{i+} + \log \pi_{+j} - \log N$$

$$\cdot H_1: A \nperp\!\!\!\nperp B \Rightarrow \log M_{ij} = \alpha_i + \beta_j + \lambda + \phi_{ij} \quad \begin{matrix} \text{log-linear} \\ \text{models} \\ \text{(log of} \\ \text{expected} \\ \text{response)} \end{matrix}$$

• this is over-parameterised; need constraints on parameters to solve

$$\cdot \text{eg } \alpha_i = \beta_j = \phi_{ij} = \phi_{i+} = 0$$

$$\cdot \text{or } \sum \alpha_i = \sum \beta_j = \sum \phi_{ij} = 0$$

•  $H_1$  is saturated; as many effective params as table cells,  
 $\Rightarrow$  can fit perfectly,  $N \pi_{ij} = n_{ij}$

• no clear interpretation of params; mainly want to test if zero

• Three-way contingency tables:

• now variables A, B, C w/ I, J, K categories

• Saturated model:

$$\log M_{ijk} = A + B + C + A:B + A:C + B:C + A:B:C = A * B * C$$

• again, over-parameterised, need similar constraints to reduce to IJK effective parameters

• ModSubmodels:

$$1 \cdot A + B + C$$

$$2 \cdot A + B + C + A:B$$

$$3 \cdot A + B + C + A:B + A:C$$

$$4 \cdot A + B + C + A:B + A:C + B:C$$

Independences

$A \perp\!\!\!\perp B \perp\!\!\!\perp C$

$A, B \perp\!\!\!\perp C$

$B \perp\!\!\!\perp C \mid A$

$\pi_{ij+} \cdot \pi_{i+k} / \pi_{i++}$

$\pi_{ijR}$

$\pi_{i++} \cdot \pi_{+j+} \cdot \pi_{++k}$

$\pi_{ij+} \cdot \pi_{i+k}$

$\pi_{i++}$

$P(B, C \mid A) P(C \mid A) P(A)$

"

$P(B, C \mid A) P(A)$

- Fitting the models:

- maximal likelihood (joint dist of observation is multinomial)
- equivalent to joint distribution of independent Poisson RVs, conditioned on their sum  
 $\Rightarrow$  fit as if observed frequencies are independent Poisson RVs

1

- Goodness of fit & model checking:

- done as for Poisson data, using  $G^2$   

$$G^2 = 2 \sum o \log \left( \frac{o}{e} \right) \quad (o = \text{observed}, e = \text{expected})$$

- can compare nested models using  $G^2$

- Residuals:

- Raw

$$\hat{e}_{ij} = n_{ij} - \hat{\mu}_{ij}$$

- Pearson

$$X_{ij} = \hat{e}_{ij} / \sqrt{\hat{\mu}_{ij}}$$

- Standardised Pearson

$$r_{P_{ij}} = X_{ij} / \sqrt{1 - h_{ij}}$$

- Deviance

$$\sum_i d_{ij}^2 = D$$

- Standardised Deviance

$$r_{D_{ij}} = d_{ij} / \sqrt{1 - h_{ij}}$$

## Generalized Additive Models: An Overview

$$g[E(Y_i)] = \eta_i = \vec{X}_i^* \vec{\Theta} + f_1(x_{i1}) + f_2(x_{i2}) + f_3(x_{i3}, x_{i4}) + \dots$$

$X_i^*$  =  $i^{\text{th}}$  row of  $X^*$

$X^*$  = model matrix for any parametric model components

$\vec{\Theta}$  = parameter vector

$f_j$  = smooth function for covariates  $x_{ij}$   
note  $\sum_j f_j(x_{ij}) = 0$

- Difference w/ GLM:

• GLM models effect of  $x_{ij}$  as  $\Theta_0 + \sum_j \Theta_j x_{ij}$  linear!

• GAM models

as  $\sum_j f_j(x_{ij})$

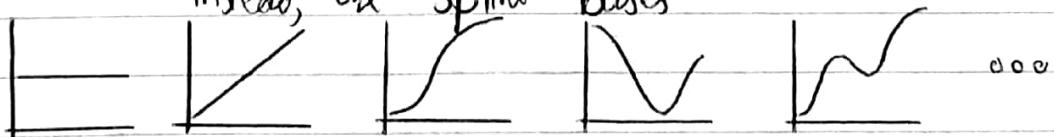
- Smooth functions:

represented using regression splines

$f_j(x_j) = \sum_k \beta_{jk} b_{jk}(x_j)$ , where  $b_{jk}$  are basis functions &  $\beta_{jk}$  regression parameters

• basis functions could be polynomial ( $x, x^2, x^3, \dots$ ), but these become highly collinear

• instead, use spline bases



- Parameter estimation:

• estimated using MLE, adding smoothing parameters  $\lambda$  to reduce overfitting (choose  $\lambda$  by, e.g., LOO-CV)

$$\text{Maximize } \lambda(\vec{\beta}) - \frac{1}{2} \sum_j \lambda_j \cdot \int [f_j^{(d)}(x_j)]^2 dx_j$$

$$= \lambda(\vec{\beta}) - \frac{1}{2} \sum_j \lambda_j \vec{\beta}^T S_j \vec{\beta}$$

↑  
don't estimate  $\lambda$ ; would just choose very low regularization!

derivative of  $f_j$   
(d, usually = 0)

encourage  
less 'rough'  
functions

known coefficients matrix  $\rightarrow S = \sum_j \lambda_j S_j$

$$\Rightarrow \hat{\beta} = (X^T W X + \lambda I)^{-1} X^T W \vec{z} \quad (\text{note biased estimator, due to penalty})$$

- Inference:

- Inference not standard due to fixed  $\lambda_j$   
 $\hat{\beta} | \vec{y} \sim N(\hat{\beta}, (I + \lambda I)^{-1})$  allows approximate inference (eg CIs) to be done

- Variable selection & model comparison also not standard