## DEVIANCE

- Output from GLM in R
- Extra criterion, for how good prediction is.
- Deviance is a log likelihood ratio statistic that compares the statistical model with the proposed GLM model

- Max likelihood estimate of the Saturated model is computed by:
$$(\hat{\theta}_1, \ldots, \hat{\theta}_n) = argmax \{ Log \lambda (\theta_1, \ldots, \theta_N) \}$$
and the value $log \lambda (\hat{\theta}_1, \ldots, \hat{\theta}_N)$ is therefore the maximum value of log likelihood function of the saturated models.

- When considering a GLM, a link function $g$ is used to constrain the parameters such that $\theta_i \& g^{-1}(x_i^T \beta)$   $i = 1, \ldots, n$
  - In this case likelihood is written $\lambda(\beta)$ and log likelihood $log [\lambda(\beta)]$
  - Parameter $\beta$ often a lower dimension than $\theta$s $(dim \beta \leq N)$ and the maximum
  - likelihood is computed such that:
$$\hat{\beta} = argmax \{ log \lambda(\beta) \}$$
- Maximum likelihood value for a GLM is then $log \lambda(\hat{\beta})$

$$Log \lambda (\beta) = \sum_{i=1}^{N} a(y_i) b(g^{-1}(x_i^T \beta)) + c(g^{-1}(x_i^T \beta)) + d(y_i)$$

Deviance also called log likelihood ratio statistic compares the saturated model with the proposed GLM
$$D = 2[ Log \lambda(\hat{\theta}_1, \ldots, \hat{\theta}_N) - Log \lambda(\hat{\beta})] = 2 Log \left[ \frac{\lambda(\hat{\theta}_1, \ldots, \hat{\theta}_N)}{\lambda(\hat{\beta})} \right] \quad \begin{array}{l} \text{Saturated} \\ \text{-GLM} \end{array}$$

- $Log(\hat{\theta}_1, \ldots, \hat{\theta}_N)$ max value of log likelihood function for saturated model
- $Log(\hat{\beta})$  value of log likelihood function when fitting model $g(E[y]) = x^T \beta$

Approximation of log likelihood $f^n$ when near its maximum
$$\theta = (\theta_1, \ldots, \theta_N)^T \qquad \theta_i = (\hat{\theta}_1, \ldots, \hat{\theta}_N)$$
$$Log \lambda(\theta) \quad near \ (\hat{\theta})$$

- Using Taylor expansion, Log likelihood can be approximated near the maximum likelihood extreme
- For saturated model with notation $(\theta_1, ..., \theta_N)^T = 0$, when $\theta$ close to $\hat{\theta}$ :

$$Log \lambda (\theta) \simeq log \lambda (\hat{\theta}) + (\theta - \hat{\theta})^T \nabla_{\hat{\theta}} + \frac{1}{2} (\theta - \hat{\theta})^T H_{\hat{\theta}} (\theta - \hat{\theta})$$

$\nabla_{\hat{\theta}}$ : Gradient of $log \lambda$ function computed at $\hat{\theta}$

$$\nabla \hat{\theta} = \begin{bmatrix} \frac{d \, log(\lambda)}{d\theta_1}|_{\theta} \\ \vdots \\ \frac{d (log \lambda)}{d\theta_N}|_{\theta} \end{bmatrix} \quad \text{vector of dim (N)}$$

At max, derivative = 0 or equal 0 bc max $\hat{\theta}$ is max likelihood solution so term is going away

$H_{\hat{\theta}}$ : Hessian Matrix of $log \lambda$ $f^n$ computed at $\hat{\theta}$

Matrix will return a negative value (bc it is a second derivative)

Hessian Matrix is symmetric and $\mathbb{R}$ ⇒ Can compute eigenvalues/vectors

Similarly for GLM model, when $\beta$ is close to $\hat{\beta}$ :

$$Log \lambda (\beta) \simeq log \lambda (\hat{\beta}) + (\beta - \hat{\beta})^T \nabla_{\hat{\beta}} + \frac{1}{2} (\beta - \hat{\beta})^T H_{\hat{\beta}} (\beta - \hat{\beta})$$

- In both case $\nabla_{\hat{\beta}}$, $\nabla_{\hat{\theta}}$ are zero vectors since $\hat{\beta}$ and $\hat{\theta}$ are maxima of $log \lambda$

- Re write Deviance:  approximation of deviance near $\hat{\theta}$ and $\hat{\beta}$

$$D \simeq 2 \{ log \lambda (\hat{\theta}) - Log \lambda (\hat{\beta}) \}$$

$$\simeq 2 \{ log \lambda(\theta) - \frac{1}{2} (\theta - \hat{\theta})^T H_{\hat{\theta}} (\theta - \hat{\theta}) - log \lambda(\beta) - \frac{1}{2} (\beta - \hat{\beta})^T H_{\hat{\beta}} (\beta - \hat{\beta}) \}$$

$$\simeq 2 Log \left[ \frac{\lambda \hat{\theta}}{\lambda \hat{\beta}} \right] - (\theta - \hat{\theta})^T H_{\hat{\theta}} (\theta - \hat{\theta}) + (\beta - \hat{\beta}) H_{\hat{\beta}} (\beta - \hat{\beta})$$

- The term $v$ is positive and will be near 0 if GLM model fits data almost as well as the saturated model does).
- Hessian matrix computed at max $H_{\hat{\beta}}$ $H_{\hat{\theta}}$ will be negative

09/02/16   ALSM 2

## Sampling Distribution of Deviance

- The likelihood function for $\beta$ can be approximated by a normal distribution near the estimate $\hat{\beta}$ such that

$$\lambda(\beta) = \frac{1}{\sqrt{2\pi|\Sigma|}} \exp\left[ -\frac{1}{2}(\beta-\hat{\beta})^T \Sigma^{-1}(\beta-\hat{\beta}) \right]$$
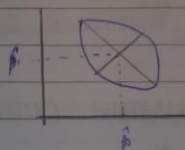
$$\log \lambda(\beta) = \text{constant} - \frac{1}{2}(\beta-\hat{\beta})^T \Sigma^{-1}(\beta-\hat{\beta})$$

$$\Sigma^{-1} = -H_\beta \quad \text{(Hessian Matrix)}$$

Can find a set of eigen vectors

$$\Sigma = L \Lambda L^T$$

↗ eigen values (diagonal matrix)

$\hat{\beta}$ - - - -

$\Sigma$ controlling shape

$\hat{\beta}$

Any covariance Matrix $\Sigma$ is symmetric and real

→ $\Sigma\, LL^T \Lambda\, LL$

Diagonal matrix which have eigenvals on diagonal and zeros otherwise

eigen vectors (orthonormal basis of $\Sigma$)    $LL^T L = I$

$$(\beta-\hat{\beta})\underbrace{\Sigma^{-1}}_{-H_\beta}(\beta-\hat{\beta}) = \left[LL(\beta-\hat{\beta})\right]^T \Lambda^{-1}\left[LL(\beta-\hat{\beta})\right]$$

$$\Lambda = \gamma^2 \qquad = \left[\gamma^2 LL(\beta-\hat{\beta})\right]^T I \left[\gamma^{-1} LL(\beta-\hat{\beta})\right]$$

Change of variable $z = \gamma^{-1} LL(\beta-\hat{\beta})$

$$Z \sim N(0, I) \qquad \text{Identity Matrix}$$

$-(\beta-\hat{\beta})H_\beta(\beta-\hat{\beta})$ is a sum of $m$ variables with distribution $N(0,1)$

identity matrix in $\mathbb{R}^m$ with $m = \text{Dim}(\beta)$

This term $\sum_{i=1}^{m} z_i^2$ has $\chi^2$ dist with $m$ d.f.

$$(\theta-\hat{\theta})^T H_\theta(\theta-\hat{\theta}) \text{ follows } \chi^2(\dim(\theta)=N)$$
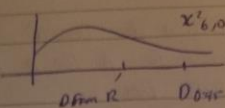
$$D \simeq 2\log\left[\frac{\lambda\hat{\theta}}{\lambda\hat{\varphi}}\right] \quad \frac{-(\theta-\hat{\theta})^T H_\theta(\theta-\hat{\theta}) + (\beta-\hat{\beta})^T H_\beta(\beta-\hat{\beta})}{\chi^2(\dim(\theta)-\dim(\beta))}$$

↗

offset should be near zero if good model

$D \sim \chi^2_{(n-m, N \to \infty)}$     Back D.E



$x^2_{6, p}$

D from R     $D_{0.45}$

Is $D$ value below $D_{0.45}$? if it is $\Rightarrow$ good model
(value from table below that of calculated)

- work backwards in Scott table book
- Value from table 12·5
- Value from R: 11·2

## AKAIKE INFORMATION CRITERION

- $\{(y_i, x_i)\}_{i=1, \dots N}$

likelihood: $L(\theta_1, \dots, \theta_N) = \prod_{i=1}^{N} P_{y|\theta}(y_i | \theta_i)$

with $P_{y|\theta}$ from an exponential family probability distribution

$L(\theta_1, \dots, \theta_N) = \prod_{i=1}^{N} \text{Exp}[a(y_i) b(\theta_i) + c(\theta_i) + d(y_i)]$

- log transform usually computed instead:

$\log L(\theta_1, \dots, \theta_N) = \sum_{i=1}^{N} [a(y_i) b(\theta_i) + c(\theta_i) + d(y_i)]$

- MLE for saturated model then

$(\theta_1, \dots, \theta_N) = \text{argmax} \log L(\theta_1, \dots, \theta_N)$

- When a GLM is used, a link function $g$ is used to constrain the parameters such that $\theta_i \propto g^{-1}(x_i^T; \beta) \quad \forall i=1, \dots, N$

- The parameter $\beta$ has often a lower dimension than $\theta_i$ i.e. $\dim(\beta) \leq N$ and the maximum likelihood estimate computed: $\beta = \text{argmax} \log L(\beta)$

- AIC is a measure of goodness of fit defined by: $AIC = -2\log L(\hat{\beta}) + 2p$

where: • $p = \dim(\beta)$ # of parameters to be estimated by the model

• $\hat{\beta}$ are the estimated parameters that maximise the likelihood or log likelihood

• $\log L(\hat{\beta})$ is the maximum value of the log likelihood.

- Best model is a trade off between one that maximises the likelihood with also having the minimum number of parameters

- Select model with lowest AIC

## CHI-Square Distribution

Assuming that $z \sim N(0,1)$ $i=1,...,N$ show that $x = z_i^2$ has a $\chi^2$ distribution
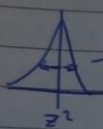
$$P_x(x) = \int P_{xz}(x,z)\,dz$$

$$= \int P_{x|z}(x|z)\,P_z(z)\,dz \qquad \text{conditional distribution}$$

$$z \sim N(0,1)$$

No uncertainty in $x$ when given $z$ ($x$ is $z^2$ remebr)

Dirac Distribution has variance of $0$: no uncertainty

$$P_{x|z}(x|z) = \delta(x - z^2)$$



width of $0$    $\sigma^2 = 0$

$z^2$    $\int_{-\infty}^{\infty} \delta(x-z)^2\,dx = 1$ (Property of dirac function)

$$= \int \delta(x-z)^2 \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-z^2}{2}\right) dz$$

Use Dirac property to rewrite dirac put.

$$\int \frac{1}{2|\sqrt{x}|} \left[ \delta(\sqrt{x}+z) + \delta(\sqrt{x}-z) \right]$$

$$= \int \frac{1}{2|\sqrt{x}|} \left[ \delta(\sqrt{x}+z) + \delta(\sqrt{x}-z) \right] \frac{1}{2\pi} \exp\left[\frac{-z^2}{2}\right] dz$$

two

$$= \frac{1}{2|\sqrt{x}|} \int \delta(\sqrt{x}-z) \exp\left[\frac{-z^2}{2}\right] + \int \delta(\sqrt{x}+z) \frac{1}{\sqrt{2\pi}} \exp\left[\frac{-z^2}{2}\right] dz$$

Next Property: $\int_{-\infty}^{\infty} f(t)\,\delta(t-T)\,dt = F(T)$

$$= \frac{1}{2\sqrt{x}} \left[ \frac{1}{\sqrt{2\pi}} \exp\left[\frac{-x}{2}\right] + \frac{1}{\sqrt{2\pi}} \exp\left[\frac{-x}{2}\right] \right]$$

Chi-Square $f(x,k) = \dfrac{x^{(\frac{k}{2}-1)} e^{\frac{-x}{2}}}{2^{k/2}\, \Gamma(\frac{k}{2})}$    $x \geq 0$

gamma

$$\mathbb{E}[x] = k$$

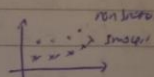- look for derivare new value of $k$ as an indicate

16/02/16   ALSM2

Smoothing Dataset
- Split into rows according to yes/no outcome
  $\hat{\theta}$ for Saturated Model $= y_i / n_i$ for each row

- Nature of explanatory variables
  • nominal   red, green, black
  • ordinal   natural ordering - week days
  • continuous   e.g. weight, age, temperature

For the $\hat{\theta}$'s $\left(\frac{y_i}{n_i}\right)$ draw a graph of all values to see difference
- GLM tries to fit a curve between these points
- Need to determine $\beta^T x$
  → could use $age^2$
  → Smoking is a binary variable, $Smoke^2$ has no effect
  → Multiply Smoke by age → 1,2,3,4,5, 0,0,0,0   ← called "mixed effect" when these variables are used
  → Multiply $age^2$ by Smoke → 1,4,9,1625, 0,0,0,0
  → i.e. Can make many explanatory variables out of 2 variables

$\beta_0 + \beta_1 x^{age} + \beta_2 x^{Smoke} + \beta_3 x^{age} x^{Smoke} + \beta_4 (x^{age})^2 + \beta_5 (x^{age})^2 x^{Smoke}$
  → largest model for what we have defined in our table

- When Smokes $x^{Smoke} = 1$ :  $\beta_0 + \beta_2 + (\beta_1 + \beta_3) x^{age} + (\beta_4 + \beta_5)(x^{age})^2$  "mixed model"
  $x^{Smoke} = 0$ :  $\beta_0 + \beta_1 x^{age} + \beta_4 (x^{age})^2$
- 6 parameters here
- Binomial or Poisson is done here

- Try all possible models and chose one with lowest AIC
- Fitted Model: $\beta_0 + \beta_1 x^A + \beta_2 x^s + \beta_3 x^A x^s + \beta_4 (x^a)^2$  has lowest AIC
- AIC:  Poisson + log link : 66.7
       Binomial + logit : 66.63
       Binomial + probit : 66.33
- Tried with $age^2$ removed - made AIC worse

- When mixed effect was removed, AIC increased
- Write the GLM $g(\mu) = \beta^T X$
- Poisson + log: $\log(\mu) = \beta_0^i + \beta_1^i x^a + \beta_2^i x^s + \beta_3^i x^a x^s + \beta_4^i x^{o^2}$
- Binomial logit: $\log[\frac{\pi}{1-\pi}] = \beta^i$ $\beta$'s will be different because of different likelihood calculation

Poisson: $P(y|\lambda) = \frac{\lambda^y}{y!} \exp(-\lambda)$     10 groups, $N=10$

$L(\lambda_1, \ldots, \lambda_{10}) = \prod_{i=1}^{10} \frac{\lambda_i^{y_i}}{y_i!} \exp(-\lambda_i)$

$\log(L) = \sum_{i=1}^{10} y_i \log(\lambda_i) - \lambda_i - \log(y_i!)$

$\hat{\lambda}_i = y_i$     Saturated model — Max likelihood shown

Deviance/AIC: $\log L(\lambda_1, \ldots, \lambda_{10})$

$\quad\quad -n\lambda_i + \sum y_i \log(\lambda) - \log(\pi y_i)$

$\frac{d\log(L)}{d\mu} = -n + \sum y_i x$

$\hat{\lambda} = \frac{\sum y}{n}$

Binomial: $P(y|\theta) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$

$L(\theta_1, \ldots, \theta_{10}) = \prod_{i=1}^{10} \binom{n_i}{y_i} \theta_i^x (1-\theta_i)^{n-x}$     (Applied to each group)

$L(\theta_1, \ldots, \theta_{10}) = \prod_{i=1}^{10} \left( \frac{n!}{x_i!(n-x_i)!} \right) p^{\sum x_i} (1-p)^{n-\sum x_i}$

$\log L = \sum y_i \log(p) + (n - \sum y_i) \log(1-p)$

$\frac{d\log L}{dp} = \frac{1}{p_i} \sum y_i + \frac{1}{1-p}(n - \sum y_i) = 0$

$\hat{\theta}_i = \frac{\sum y_i}{n_i}$

Saturated model replace $\theta_i$ with $\hat{\theta}_i = \frac{y_i}{n_i}$

Difference between $a/c$? look at coefficients $\beta$'s to determine.
- look at standard errors associated with the $\beta$'s
- "couldn't remove any explanatory variable from model → variable is important"