

In case of linear models the true distⁿ of f_i is of the form $g(y; \theta)$ for some θ if $f \in \mathcal{L}(x)$.

$$\Leftrightarrow \lambda = 0$$

When this holds $\frac{n\hat{\sigma}^2}{\sigma_0^2} \sim \chi_{n-p}^2$

for a x_p^n r.v. X , $E\left(\frac{1}{X}\right) = \frac{1}{\gamma-2}$

$$\text{So, } E\left[\frac{\sigma_0^2}{\hat{\sigma}^2}\right] = \frac{n}{n-p-2}$$

$$E[2\Delta] = n E[\log(2\hat{\sigma})^2] + \frac{n(n+p)}{n-p-2}$$

- Both 2Δ and AIC have expectationally of order $O(n)$.

- The bias in estimating $E[2\Delta]$ by AIC is

$$E[AIC - 2\Delta] = n + 2(p+1) - \frac{n(n+p)}{n-p-2}$$

$$= 2(p+1) - \frac{2n(p+1)}{n-p-2} = O(1)$$

- Although the bias is of smaller order than $E[2\Delta]$, we can get an exactly unbiased of $E[2\Delta]$ by using a modified criterion.

$$AIC_c = n \log 2\hat{\sigma}^2 + \frac{n(n+p)}{n-p-2}$$

In many cases, AIC_c is a much better estimate of 2Δ than AIC.

If \hat{y}_i is not in $\mathcal{E}(x)$, then

$$n E[\hat{\sigma}^2] = \sigma_0^2(n-p+2) \quad E[RSS_p] = E(ME) + \sigma^2(n-2p)$$

$$E\left[\frac{\sigma_0^2}{\hat{\sigma}^2}\right] \approx \frac{\sigma_0^2}{E[\hat{\sigma}^2]} = \frac{n}{2+n-p} = 2\sigma^2 + \sigma^2 p + \sigma^2(n-2p)$$

$$\text{So, } E[AIC - 2\Delta] \approx n + 2(p+1) - \frac{(n+p+2)n}{2+n-p} = 2\left[p+1 - \frac{np}{2+n-p}\right] \approx \frac{np}{2+n-p} \approx p$$

for large n .

- A slightly different criterion is if we assume σ^2 is unknown.

- Say the no. of unknown is p and

$$-2\log(g(\hat{y}_j | \hat{\theta}(\hat{y})) = n \log(2\pi\sigma^2) + \frac{(\underline{y} - \underline{x}\hat{\beta})'(\underline{y} - \underline{x}\hat{\beta})}{\sigma^2}$$

$$ATC = \frac{RSS_p}{\sigma^2} + 2p$$

\rightarrow very similar to Cp.

- An obvious generalization of this version of AIC is to consider the criterion of the form

$$ATC = \frac{RSS_p}{\sigma^2} + a_n p$$

- a_n is allowed to depend on n

- $a_n = \log n \rightarrow$ Bayes information criterion (BIC).

Polynomial Regression:

Polynomial in one variable

Polynomial in one variable:

- A major problem when fitting a high degree polynomial.

Set $x_{ij} = x_i^{j+q}$, $k = p-1$ ($\leq n-1$) in general multiple regression set up.

$$y_i = \beta_0 + \beta_1 x_{it} + \beta_2 x_{it}^2 + \dots + \beta_k x_{it}^k, \quad i=1, 2, \dots, n$$

When $k > 6$, X become

Assume that $x_i \sim U[0,1]$, then for large n

$$(x'x)_{rs} = n^{-1} \sum_{i=1}^n x_i^r x_i^s = n^{-1} \left[\sum_{i=0}^n x_i^r x_i^s \right]$$

$$\approx n \int_0^1 x^r x^s dx$$

$$= n \int_0^1 x^{r+s} dx = \frac{n}{r+s+1}$$

So, $(x'x) \rightarrow (n+1) \times (n+1)$ principle minor of the Hilbert matrix.

$$H = \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \dots & \dots \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \dots & \dots \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \dots & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots & \ddots \end{pmatrix}$$

When $K=9$, the inverse of H_{10} (10×10 principle minor of H) has elements of the magnitude 3×10^{10} .

If a small error of 10^{-10} in one element of $\underline{x}'\underline{y}$

→ will lead to an error of observation about 3 in an element of $\hat{\underline{\beta}} = (\underline{x}'\underline{x})^{-1}\underline{x}'\underline{y}$

- Two things can be done.

① Normalize x_i so that they run from -1 to 1.

$$x_i' = \frac{2x_i - \max(x_i) - \min(x_i)}{\max(x_i) - \min(x_i)}$$

② Use orthogonal polynomial.

Consider the model.

$$\gamma_i = \gamma_0 \phi_0(x_i) + \gamma_1 \phi_1(x_i) + \cdots + \gamma_k \phi_k(x_i)$$

When $\phi_r(x_i) \rightarrow r^{\text{th}}$ degree polynomial in x_i , $r=0, 1, \dots, k$.

$$\sum_{i=1}^n \phi_r(x_i) \phi_s(x_i) = 0, \forall r, s, r \neq s.$$

Then $\underline{Y} = \underline{X}\underline{\gamma} + \underline{\epsilon}$

$$\underline{X} = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \cdots & \phi_k(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \cdots & \phi_k(x_2) \\ \vdots & & & \\ \phi_0(x_n) & \cdots & & \phi_k(x_n) \end{pmatrix}$$

$$x'x = \left[\begin{array}{c} \sum \phi_0^2(x_i) \\ \vdots \\ \sum \phi_k^2(x_i) \end{array} \right] \quad \text{and} \quad y = \left[\begin{array}{c} y_0 \\ \vdots \\ y_K \end{array} \right]$$

$$\hat{\gamma} = (x'x)^{-1} x'y \Rightarrow \hat{\gamma}_r = \frac{\sum \phi_r^2(x_i) y_i}{\sum \phi_r^2(x_i)}$$

for $r=0, 1, \dots, K$

Generating orthogonal polynomials.

Generating Orthogonal Polynomials.

$$\phi_{r+1}(x) = 2(x - a_{r+1})\phi_r(x) - b_r\phi_{r-1}(x)$$

$$\phi_0(x) = 1, \quad \phi_1(x) = 2(x - a_1)$$

x is normalized so that $-1 \leq x \leq 1$

a_{r+1} & b_r are chosen as

$$a_{r+1} = \frac{\sum x_i \phi_r^2(x_i)}{\sum_{i=1}^n \phi_r^2(x_i)} ; \quad b_r = \frac{\sum \phi_r^2(x_i)}{\sum_{i=1}^n \phi_{r-1}^2(x_i)}$$

$$r=0, 1, \dots, K-1$$

$$\hat{\gamma}_0 = \bar{y}$$

HW: find the F-test for testing $\hat{\gamma}_K = 0$.

Robust Regression :

- ISE's are the most efficient unbiased estimates of the regression coefficient when the errors are normally distributed.

- However, they are not very efficient when the dist' of the errors is long-tailed.

→ If we expect outlier in the data.

[under normality assumption].

- When fitted a regression, we minimize some average measure of the size of the residuals.

LS: Least mean of squares.

- fits a regression by minimizing the mean of the squared residuals.

⇒ the sum of the squared residuals.

$$\min_b \frac{1}{n} \sum_{i=1}^n e_i^2(b)$$

Here average is interpreted as the mean and size as the square.

- The sensitivity of LS to outliers is due to two factors.

① If we measure size using the squared residuals, any residual with large magnitude will have a very large size relative to the others.

② By using a measure of location such as a mean that is not robust, any large square will have a very strong impact on the criterion.

Two remedies:

① Measure size in some other way

→ by replacing e^2 by some $f(e)$ which reflects the type of the residual in a less extreme way.

To be a sensible measure of size, the $f(e)$ should be — symmetric $[f(e) = f(-e)]$

— non-ve $[f(e) \geq 0]$

— monotone $[f(|e_1|) \geq f(|e_2|)$

if $|e_1| > |e_2|$

→ leads to M-estimation

② Replace the sum (or the mean) by a more robust measure of location such as the median or a trimmed mean.

↳ least median square (LMS)

→ least trimmed squares (LTS)

M-estimates:

Suppose the observed responses y_i , independent and have the density f^n .

$$f_i(y_i | \beta, \sigma) = \frac{1}{\sigma} f\left(\frac{y_i - x_i' \beta}{\sigma}\right) \quad (1)$$

where σ is a scale parameter.

Then the log likelihood corresponding to the density f^n is

$$\ell(\beta, \sigma) = -n \log \sigma + \sum_{i=1}^n \log \left[f\left(\frac{y_i - x_i' \beta}{\sigma}\right) \right]$$

Write $\rho = -\log f$

$$\ell(\beta, \sigma) = -\left\{ n \log \sigma + \sum_{i=1}^n \rho\left[\frac{y_i - x_i' \beta}{\sigma}\right] \right\}$$

So, to estimate β & σ using ML method, we minimize

$$n \log \sigma + \sum_{i=1}^n \rho\left[\frac{e_i(\beta)}{\sigma}\right]$$

Differentiating leads to the estimating eq.

$$\sum_{i=1}^n \psi\left(e_i(\beta)/\sigma\right) x_i = 0 \quad \text{or } \psi = \rho' \quad (2)$$

$$\sum_{i=1}^n \psi\left(e_i(\beta)/\sigma\right) e_i(\beta) = ns$$

Example: if $\rho(x) = \frac{1}{2}x^2$, $\psi(x) = x$.

(2) reduces to the normal equation of LSE, $\hat{\beta}$

$$\hat{\beta} = (x'x)^{-1} x'y$$

$$(3) \rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n e_i^2(\hat{\beta})$$

(2) $\rho(x) = |x|$, the corresponding estimates are values of β & σ that minimize,

$$n \log s + \frac{1}{s} \sum_{i=1}^n |\epsilon_i(b)|$$

- is called the L_1 estimate.

- The L_1 estimate is the MLE if f is the double exponential density proportional to $\exp(-|x|)$.
- An alt. name is LAD (least absolute deviation) estimate.
- If we have no particular density in mind, we can choose φ to get the estimate by choosing a ρ for which $\Psi = \rho'$ is bdd. (e.g. $\rho = (\alpha, \beta)$).
- We can generalize ③ & ④

$$\sum_{i=1}^n \Psi \left[\frac{\epsilon_i(b)}{s} \right] x_i = 0$$

$$\sum_{i=1}^n X \left[\frac{\epsilon_i(b)}{s} \right] = 0$$

where X is also chosen to make the scales estimate robust.

The resulting estimates [solⁿ of ③] is called the M-estimates. Since their definition is motivated by the ML estimates eq's.

- There is no requirement that Ψ & X be related to the density f_h of the responses.

Example:

$$\text{Let, } \Psi(x) = \begin{cases} -K & x < -K \\ x & -K \leq x \leq K \\ K & x > K \end{cases}$$

When K is a const to be chosen.

- The value of K usually chosen is 1.5, which gives a reasonable compromise between LS & L₁ estimate.
- It can be shown that a necessary condⁿ for considering ~~other where~~ when the parameter are estimated using ⑤ is

$$E[\Psi(z)] = 0$$

$$E[X(z)] = 0$$

when Z has density $f(z)$ s.t. $\int f(z) dz = 0$

For example, ⑥ will be satisfied if f is symmetric about zero and of Ψ is $\Rightarrow \Psi(z) = \Psi(-z)$.

Example: Suppose $\Psi(x)$ is same as in Example ⑤.

Then ⑥ is satisfied.

Then the scale parameter is $\hat{\theta}$ estimated by taking $\Psi(\hat{x}) = \Psi^2(\hat{x}) - c$ for some constant c .

which is chosen to make the estimate consi. When f is the normal density.

From ⑦, we require that $c = E[\Psi(z)^2]$, where z is standard normal.

Example: Another choice of X is $\chi(z) = \text{sign}(|z| - 1/2)$ for some const c .

corresponding

Then the estimating equation

$$\sum_{i=1}^n \text{sign}(|x_i(z)| - 1/2) = 0$$

$$\sum_{i=1}^n x_i (e_i(b)/s) = 0$$

- has sofn.

$$s = c \text{ median}_i |e_i(b)|$$

- mean absolute deviation

(MAD estimate)

To make it consider, well write set was full.

then requires, $c^* = \Phi^{-1}(3/4) = 0.6749$

$$c = 1.4326$$

- Regression coeff estimated using M-estimation

- method are almost as efficient as the LS method if the errors are normal but are much more robust if the error dist is long tailed.

- $e_i(b) = (y_i - \hat{y}_i)^2$ is less sensitive

- M estimates of regression coefficients are as vulnerable as LS to outliers in the expl variables.

- Estimate Based on Robust Location & Scale measures

LMS: least median Sq estimates.

Replacing the mean by a robust measure of location, but retain the squared residuals,

↳ LMS which minimizes

$$\text{median}_i \otimes e_i^2(b).$$

An alternative to LMS estimator is to use the trimmed mean rather than the median.

\hookrightarrow LTS (least trimmed mean) estimates.

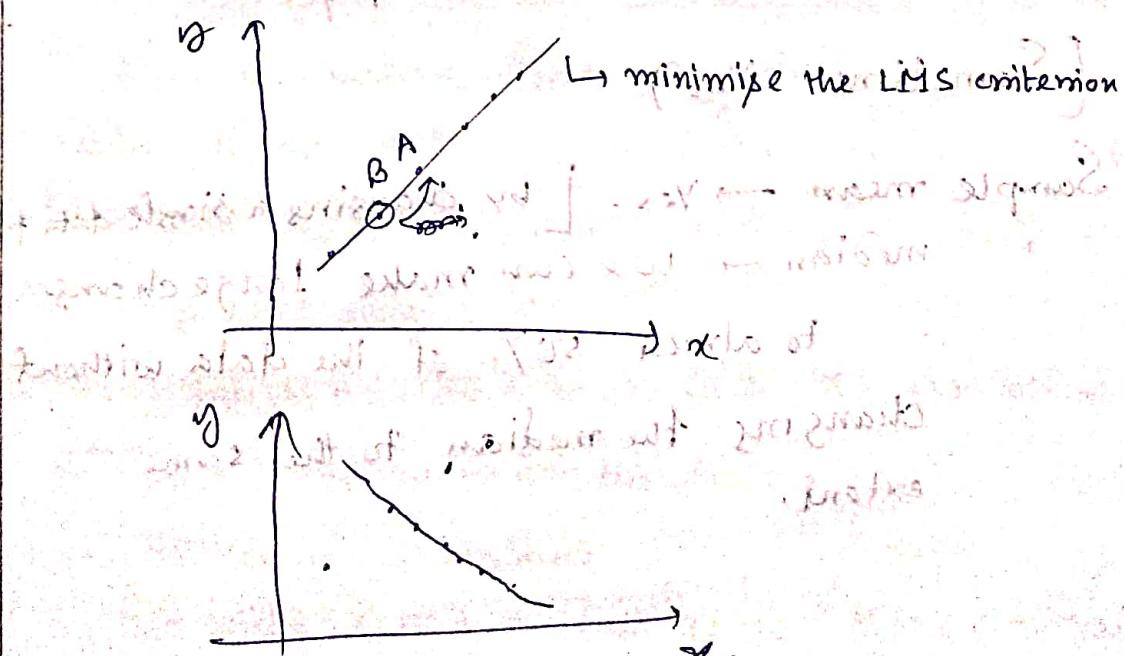
This minimizes $\sum_{i=1}^n e_i^2(b)$

where b is chosen to achieve a robust estimate of the estimator and $e_{(1)}^2(b) \leq e_{(2)}^2(b) \leq \dots \leq e_{(n)}^2(b)$ are ordered squared residuals.

The amount of trimming has to be quite severe to make it robust.

$n = [n/2] + 1$ is a popular choice \rightarrow which amounts to trimming of 50% of residuals.

LMS & LTS are very robust to outliers in both errors and explanatory variables but can be very unstable in a different way.



- In addition, these estimates are very inefficient compared to LSE of the data if the data are actually normally distributed.

- Asymptotic relative efficiency of LMS, relative to LSE is zero.

$$\left[\frac{V(LSE)}{V(LMS)} \rightarrow 0 \text{ as sample size } \rightarrow \infty \right]$$

$$\text{and also } \frac{Var(LSE)}{Var(LTS)} = 0.08$$

Measuring Robustness:

Breakdown Point (BD): measure how well an estimate can resist gross corruption of a fraction of the data.

Suppose that we select a fraction of the data. Can we cause an arbitrarily large change in the estimate by making a suitable large change in the selected data pts?

[Sometimes Yes]

Sample mean \rightarrow Yes. [by changing a single data pt]
median \rightarrow We can make large changes to almost 50% of the data without changing the median to the same extent.

Defⁿ: The finite sample breakdown point of an estimate is the fraction of the data that can be given arbitrary values without making the estimator arbitrarily ~~bad~~.

- The sample mean has BD pt $\frac{1}{n}$ and the sample median a BD pt of almost $\frac{1}{2}$.
- A BD point of $\frac{1}{2}$ is the best possible. If more than 50% of the sample is contaminated, it is impossible to distinguish between "good" or "bad" obs. Since the outliers are now typical obs of the sample.
- LSE of a regression coefficient is a linear comb^b of the responses. BD pt of LSE is $\frac{1}{n}$.
- $\sum |y_i - \theta|$ w.r.t. θ Least absolute deviation estimator $\sum_{i=1}^n |y_i - x_i' \beta|$
- Since the median has a very high BD point and median of the data y_1, \dots, y_n minimizes the least absolute deviation $\sum |y_i - \theta|$ as a fn coefficients would also have a high BD point.
- This is not the case, in fact, the BD pt of L₁ estimator is the same as the LSE.
- It can be shown that when X is of full rank there is a value of $\underline{\beta}$ minimizing $\sum_{i=1}^n |e_i(\underline{\beta})|$ for which at least p

residuals are zero.

- Further if one data point is arbitrarily far from the others, this data pt must have a zero residual.
- It follows that by moving the data pt can arbitrary amount, we must also be moving the fitted plane by an arbitrary amount, since the fitted planes passes through the extreme data pt.
- The same is true for M-estimator.
- The LMS & LTS estimators were inefficient compared to M-estimator.

But they have BD points almost $\frac{1}{2}$, the best possible one.

Influence Curve

F: K dim distⁿ.

θ : poplⁿ parameter that depends on F.

We write $\theta = T(F)$

we call T as a statistical functional.

Example: Simplest example of statistical functional is the mean $E_F(x)$ of a r.v. x.

$$T(F) = E_F(x) = \int x dF(x)$$

Example: \underline{z} is a random vector with distⁿ f, then $E(\underline{z}\underline{z}') = \int \underline{z}\underline{z}' dF(\underline{z})$

Defⁿ: If x_1, x_2, \dots, x_n are iid RVS each with distⁿ f , then the empirical distⁿ (edf) $F_n(\hat{F}_n)$

is the distⁿ which places $\frac{1}{n}$ at each of the pts x_i ,

($i=1, 2, \dots, n$). In other words, if t is a number then

$$F_n(t) = \frac{\text{no of elements in the sample } \leq t}{n}.$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{I}(x_i \leq t)$$

Integration w.r.t. edf is just the averaging. If h is a function then

$$\int h(x) dF_n(x) = \frac{1}{n} \sum_{i=1}^n h(x_i)$$

Many statistic used to estimate parameters $T(F)$ are the plug-in-estimator of the $T(F_n)$, where F_n is the edf based on a sample of size n .

Example: z_1, z_2, \dots, z_n random sample from a multivariate distⁿ F . The plug-in-estimator of

$$T(F) = \int z dF(z)$$

$$\text{as } T(F_n) = \int z dF_n(z) = \frac{1}{n} \sum z_i \rightarrow \text{the sample mean.}$$

the plug-in estimator of

$$T(F) = \int z^2 dF(z)$$

$$T(F_n) = \int z^2 dF_n(z) = \frac{1}{n} \sum_{i=1}^n z_i^2$$

- Consider a regression with a response variable y and exp variables x_1, \dots, x_{p-1} .
- Regard the regression data, $(\underline{x}_i, \underline{y}_i)$; $i=1, 2, \dots, n$ as n iid random $(p+1)$ vectors, distributed as $(\underline{x}, \underline{y})$, having a joint distⁿ f .

- We take the vector \underline{x} is random and having initial element as 1; if the regression model has a constant term.
- We also assume that the conditional distⁿ $y|\underline{x}$ has density $f^n g\left(\frac{y - \underline{x}' \beta}{\sigma}\right)$; where g is known density.

(LS): Consider the functional

$$T(f) = [E_f(\underline{x}' \underline{x}')]^{-1} E_f[\underline{x}' \underline{y}]$$

The plug-in-estimator of this T is

$$T(f_n) = \left[\bar{n} \sum \underline{x}_i \underline{x}_i' \right]^{-1} \left[\bar{n} \sum \underline{x}_i \underline{y}_i \right]$$

$$= (\underline{\underline{x}}' \underline{\underline{x}})^{-1} \underline{\underline{x}}' \underline{\underline{y}}$$

- Suppose that F is a distⁿ pm.
- We can model a small change in F at a fixed (non-random) value $\underline{z} = (\underline{x}_0', \underline{y}_0')$. by considering a mixture of distⁿ

$$F_t = (1-t)F + t\delta_{\underline{z}_0}$$

Where $\delta_{\underline{z}_0}$ is the distⁿ pm of the constant \underline{z}_0 and t_0 is close to zero.

The sensitivity of T can be measured by the rate at which $T(f_t)$ changes for small values of t .

Defⁿ: The influence curve (IC):

of a statistical functional T is the derivative of $T(f_t)$ w.r.t. t evaluated at $t=0$, and is a measure of the rate at which T responds to a small change amount of contamination at x_0 .

IC depends on both F & x_0 .

$$IC(F, x_0) = \frac{d T(f_t)}{dt} \Big|_{t=0}$$

Example:

$$\text{Let, } T(F) = \int x dF(x).$$

$$T(f_t) = \int x dF_t(x) = (1-t) \int x dF(x) + t \int x d\delta_{x_0}(x)$$

$$= (1-t)T(F) + t x_0$$

$$\text{So, } \frac{T(f_t) - T(F)}{t} = x_0 - T(F)$$

$$IC(F, x_0) = x_0 - T(F) + t x_0 (1-t) = x_0$$

This suggests that a small amount of contamination can cause an arbitrarily large change.

→ The term is highly non-robust.

Example: (IC for LSE)

Let, T be the LSE functional.

$$[E_F(\underline{x} \underline{x}')]^{-1} E_F(\underline{x} \underline{y})$$

Write $\Sigma_{F_t} = E_{F_t}(\underline{x} \underline{x}')$ & $\gamma_{F_t} = E_F(\underline{x}, \underline{y})$

$$\text{Then } T(F_t) = (\Sigma_{F_t})^{-1} \gamma_{F_t}$$

$$\Sigma_{F_t} = E_{F_t}(\underline{x} \underline{x}') = (1-t) E_F(\underline{x} \underline{x}') + t \underline{x}_0 \underline{x}'_0$$

$$= (1-t) \Sigma_F + t \underline{x}_0 \underline{x}'_0$$

$$= (1-t) \left\{ \Sigma_F + t' \underline{x}_0 \underline{x}'_0 \right\}$$

$$t' = \frac{t}{1-t}$$

$$\Sigma_{F_t}^{-1} = (1-t)^{-1} \left[\Sigma_F^{-1} - \frac{t' \Sigma_F^{-1} \underline{x}_0 \underline{x}'_0 \Sigma_F^{-1}}{1+t' \underline{x}_0 \Sigma_F^{-1} \underline{x}'_0} \right]$$

$$= (1-t)^{-1} \left[\Sigma_F^{-1} - \frac{t' \Sigma_F^{-1} \underline{x}_0 \underline{x}'_0 \Sigma_F^{-1}}{(1-t) \Sigma_F^{-1} + t \underline{x}_0 \Sigma_F^{-1} \underline{x}'_0} \right]$$

$$\gamma_{F_t} = (1-t) \gamma_F + t \underline{x}_0 \gamma_0$$

$$T(F_t) = T(F) + t' \Sigma_F^{-1} \underline{x}_0 \gamma_0 - t' \Sigma_F^{-1} \underline{x}_0 \underline{x}'_0 \Sigma_F^{-1} + o(t)$$

$$\underline{x}_t = o(t) \Rightarrow \frac{1}{t} \underline{x}_t \rightarrow 0 \text{ as } t \rightarrow 0$$

$$\Rightarrow \frac{T(F_t) - T(F)}{t} = \Sigma_F^{-1} \underline{x}_0 \gamma_0 - \Sigma_F^{-1} \underline{x}_0 \underline{x}'_0 \Sigma_F^{-1} + o(1)$$

Take $t \rightarrow 0$;

$$IC(F, \underline{x}_0) = \Sigma_F^{-1} \underline{x}_0 [\gamma_0 - \underline{x}'_0 T(F)]$$

→ is ~~a~~ unbounded in both \underline{x}_0 & \underline{y}_0 which indicate
that LSE is not robust.

Stepwise Regression:

Suppose that we have the model

$$Y = X\beta + \epsilon$$

X is $n \times (k+1)$ want to identify the significant variables having non-zero regression coefficients.

Suppose, we divide the k variable into two groups.

1. The 1st one consists of $p-1$ variables that we consider as important.
2. The second are which contains $k-p+1$ variables consists of variables whose coefficient we suspect are zero.

We can test if the second set contains no significant variable using

$$F = \frac{\frac{RSS_p - RSS_{K+1}}{n-K-1}}{\frac{RSS_{K+1}}{K-p+1}}$$

- test which discriminate between two models having $p-1$ and K variables.

If $K=p$; then

$$F^* = \frac{RSS_p - RSS_{p+1}}{RSS_{p+1}} (n-p-1)$$

- If addition of a specified extra variable is necessary.

Forward Selection:

1. Start with the model with constant term by only compute F^* with $p=1$ for all the available exp variables and pick that variable for which F^* is the largest.
2. Repeat the procedure for $p=2, 3, \dots$ selecting at each step the variable currently not included that gives the maximum value of F^* .
3. Stop if the maximum value of F^* at any stage does not exceed some threshold value of F_{IN} .
4. This procedure is called forward selection.

A slight variation on FS is to pick up at each stage the variable not currently included in the model having the largest possible partial correlation with the response given the variable currently included.

Backward elimination:

As an alternative to FS, we can start the full model using all K variables (provided $K < n$) and compute F with $p=K$ for each of the K variables.

- Eliminate the variable having the smallest F-statistic from the model provided F_{out} is less than some threshold value.
- The procedure is continued until all variable are eliminated or the smallest F fails to be less than F_{out} .

- This procedure is called the backward elimination.

Stepwise Regression

A method that combines FS and BE.

- This is just a FS followed by a BE step at each stage.

- Start with a constant term alone then perform a FS step, adding a single variable.

- This is followed by a BE step, removing a variable if the corresponding F is less than F_{out} .

- This combination of an FS step followed by a BS step is replaced until no further variable is added at the FS step. Provided $F_{out} \leq F_{IN}$

the stepwise regression algorithm must eventually terminates.

- Suppose that at the beginning of an FS step there are $p-1$ variables in the current model, which has a residual RSS_p . Then either, no variable can be added, and the algorithm terminates or an additional variable is added, resulting in the new RSS_{p+1} which must satisfy:

$$\frac{\text{RSS}_p - \text{RSS}_{p+1}}{\text{RSS}_{p+1}} > F_{IN}$$

$$\text{or, } RSS_p > RSS_{p+1} \left(1 + \frac{F_{IN}}{n-p-1}\right) > RSS_{p+1}$$

Now, do a BE step. Either no variable is deleted or we get a new model with p_1 variables and a new RSS equal to RSS_{p+1}^*

$$RSS_{p+1}^* \leq RSS_{p+1} \left(1 + \frac{F_{out}}{n-p-1}\right) \leq RSS_{p+1} \left(1 + \frac{F_{IN}}{n-p-1}\right)$$

$$\text{So, } RSS_{p+1}^* \leq RSS_{p+1}$$

If $F_{out} \leq F_{IN}$,

Testing for Heteroscedasticity: Breusch-Pagan Test

The linear regression equation with heteroscedastic error

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{p-1} x_{ip-1} + \epsilon_i ; i=1,2,\dots,n.$$

$$E(\epsilon_i) = 0$$

$$V(\epsilon_i) = \sigma_i^2$$

Assume that error variables is of the linear term of independent variables,

$$\sigma_i^2 = b_0 + b_1 x_{i1} + \dots + b_k x_{ik-1} ; k: \text{some exp variables}$$

Choose the set of ~~exp~~ variables with which we can assume that σ_i^2 vary.

$$H_0: \text{Var}(\epsilon_i | x_i) = \sigma_i^2 \quad H_1: \text{not } H_0$$

$$\Leftrightarrow H_0: \sigma_1^2 = \dots = \sigma_{k-1}^2 = 0$$

$H_1: \text{not all } \sigma_i^2 \text{ are zero}$

• Since we never know the actual error, we use then estimates $\hat{\epsilon}_i$, which are OLS residuals.

$$\text{consider } \hat{\epsilon}_i^2 = \delta_0 + \delta_1 x_{i1} + \dots + \delta_{k-1} x_{ik-1} + \hat{\epsilon}_i$$

We can use a F-test.

- F statistic for testing this hypothesis is the goodness of fit test.

Let, $R_{\hat{\epsilon}}^2$ be the multiple correlation coeff of this regression.

$$F = \frac{R_{\hat{\epsilon}}^2}{1 - R_{\hat{\epsilon}}^2} \frac{n-k}{n-1} \sim F_{n-1, n-k} \text{ under } H_0.$$

Reject H_0 if $F > F_{n-1, n-k}(\alpha)$

A form of Breusch-Pagan test for heteroscedasticity is $n \times R_{\hat{\epsilon}}^2 \sim \chi^2_{k-1}$ [$k-1 \rightarrow \text{no of variables}$]

(Lagrange's Multiplier Statistic)

$$LM = \left(\frac{\partial L}{\partial \theta} \right)^T \left[-E \left[\frac{\partial^2 L}{\partial \theta \partial \theta'} \right] \right] \left(\frac{\partial L}{\partial \theta} \right)$$

BP test has been shown to be sensitive to any violation of normality assumption.

Auxiliary equation for error variance can be

$$\sigma_i^2 = \delta_0 + \delta_1 x_{i1} + \dots + \delta_{k-1} x_{ik-1}$$

$$\epsilon_i = \delta_0 + \delta_1 x_{i1} + \dots + \delta_{k-1} x_{ik-1}$$

$$\ln(\sigma_i^2) = \delta_0 + \delta_1 x_{i1} + \dots + \delta_{k-1} x_{ik-1}$$

$$\text{or } \sigma_i^2 = \exp[\delta_0 + \delta_1 x_{i1} + \dots + \delta_{k-1} x_{ik-1}]$$

We do not know σ_i^2 , estimate $\hat{\epsilon}_i^2$ and use $\hat{\epsilon}_i^2$ for σ_i^2 . $|\hat{\epsilon}_i|$ for σ_i and $\ln(\hat{\epsilon}_i^2)$ for $\ln(\sigma_i^2)$

White test; Goldfeldt - Quandt Test; Breusch Godfrey test

White test: This is explicitly intended to test for form of heteroscedasticity.

- The relation of ϵ^2 with all the independent variables (x_i), the squares (x_i^2) and all cross-product ($x_i x_j$, $i \neq j$)

- Just like BP test, we consider $\hat{\epsilon}_i^2$ on all this above variables and compute R^2 . Let this total no of $(x_i, x_j^2, x_i x_j)$ be k .

- $n R^2$ is distributed as χ_k^2 for large sample.

The abundance of independent variables is a weakness in this original ~~as~~ of white test.

The idea is, derived by observing the original regression model.

- Consider in the 2nd stage.

$$\hat{\epsilon}^2 = \delta_0 + \delta_1 \hat{y} + \delta_2 \hat{y}^2 + v$$

Where \hat{y} is the predicted value of y .

Use F or LM test $H_0: \delta_1 = \delta_2 = 0$

[very bad for small sample]

→ a sp. case of White test

→ perform very badly for small samples.

The holdfield & Guandt Test:

- White's test and BP test both focus on smoothly changing variances for the error term.

- The GS test for heteroscedasticity divides n observations into groups and tests whether two specific groups have disturbances with the same variances, because the errors are unknown we compare the MSEs for these groups.

Test:

1. Divide n observations into k groups.

Let, n_1, n_2, \dots, n_k denote the no. of observations in groups 1, 2, ..., k.

2. Choose two groups 1, 2, say, for which to test $H_0: \sigma_1^2 = \sigma_2^2$ v/s. $H_1: \sigma_1^2 \neq \sigma_2^2$

3. Estimate K explanatory model of interest by OLS using n_i obs in group 1 let RSS be CSR ,

4. Do it for group-2, let it be SSR_2

5. Name the larger one of $\frac{SSR_1}{n_1-k}$ and $\frac{SSR_2}{n_2-k}$

as $\frac{SSR_L}{n_L-k}$ and the smaller one as $\frac{SSR_S}{n_S-k}$

6. Form the statistic

$$G = \frac{\frac{SSR_L}{(n_L-k)}}{\frac{SSR_S}{(n_S-k)}}$$

7. Reject H_0 if $G > F_{n_L-k, n_S-k}$

8. The GLS test is valid if the underlying errors are normally distributed.

9. It is approximately valid if both n_L and n_S-k are large and errors are not normally distributed.

10. The most frequent approaches of GLS test involves dividing the sample in two groups.

Consider $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$
 $V(\epsilon_i) = \sigma^2$, $Cov(\epsilon_i, \epsilon_j) = 0, i \neq j$

how can you estimate $Var(\hat{\beta}_1 | x)$

$$Var(\hat{\beta}_1 | x) = \frac{\sum (x_i - \bar{x})^2 \sigma^2}{\left[\sum (x_i - \bar{x})^2 \right]}$$

$$\neq \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \\ &= \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2} \end{aligned}$$

Breusch - Godfrey Test for AR(1) :

$$Y_t = \beta_0 + \beta_1 X_t + \epsilon_t$$

$$\epsilon_t = p_1 \epsilon_{t-1} + p_2 \epsilon_{t-2} + \dots + p_q \epsilon_{t-q} + \eta_t$$

$$\epsilon_t \sim \text{iid } N(0, \sigma^2)$$

or uncorrelated.

$$H_0: p_1 = p_2 = \dots = p_q = 0$$

* Estimation β_0 and β_1 using OLS and obtain $\hat{\epsilon}_t$ (residuals)

Regress the estimated residuals on the lagged values upto order 1 ag q and all original variables.

$$\hat{\epsilon}_t = p_1 \hat{\epsilon}_{t-1} + \dots + p_q \hat{\epsilon}_{t-q} + \gamma_1 X_t + \gamma_0 + \eta_t$$

→ Treat it as an regression model and obtain R^2

$$NR^2 \sim \chi^2_q \quad [\text{asymptotic approximation}]$$

↳ no of data pts available for the seemed stage regression

$$(n-q) R^2 \sim \chi^2_{q-1}$$

Ridge Regression

In general case, we need a reliable method of estimating the optimal value of k if ridge estimate is to be superior to LSE.

- Popular method is based on CV or GCV.

The estimate of the prediction error of the ridge predictor

$$\frac{1}{n} \sum_{i=1}^n [y_i - x_i' \hat{\beta}_i(k)]^2$$

$\hat{\beta}_i(k)$: ridge estimate of β calculated with ridge coeff k but leaving out the i th data point.

Simplify it $\frac{1}{n} \sum_{i=1}^n [y_i - x_i' \hat{\beta}(k)]^2$

$$= \text{diag}[I - a_{ii}(k)]^2$$

Where, $a_{ii}(k)$: i th diagonal element of $A(k)$

$$= \text{diag}[I - X(X'X + kI_p)^{-1}X']^2$$

The value of k is then chosen to minimise this criterion.

- The ridge estimate can be regarded as the soln of

a constrained LS problem.

Consider minimising the SS

$$\|y - Xb\|^2$$

subject to the constraint $\sum b_j^2 \leq s$, where s is some specified constant.

- If $\sum \hat{\beta}_j^2 \leq s$, the LSE solves this constrained probm

If $\sum \hat{\beta}_j^2 > s$, then the solⁿ is different

and given by the Langrange's Multipliers.

$$x'x\hat{\beta} - x'y + \lambda \hat{\beta} = 0 ; \sum \hat{\beta}_j^2 = s.$$

which has the solⁿ of the form

$$\hat{\beta}(\lambda) = (x'x + \lambda I_p)^{-1} x'y$$

Where λ satisfies $\sum \hat{\beta}_j^2(\lambda) = s$.

- If $x'x = I_p$, the ridge estimate is

$$\hat{\beta}(\lambda) = s_j \frac{\hat{\beta}}{1+\lambda} \rightarrow \text{therefore, each coefficient}$$

is shrunk by a constant factor.

Garrote & Lasso Estimate

- One drawback of the ridge estimate is

unlike the f-subset selection it retains all the variation in the model. No possibility of a simpler model with fewer variables.

An alternative to ridge to get a simple model

Garrote.

In Garrote, the individual LS Coefficients

$\hat{\beta}_j$ are shrunk by non-ve quantity c_j

- leads to the Garrote estimate

$$\tilde{\beta}_j = c_j \hat{\beta}_j \rightarrow \text{LSE}$$

↑

• Shrinkage factors.

The c_j 's are chosen to minimise the LS estimators

$$\|\underline{y} - \underline{x}\hat{\beta}\|^2 = \sum_{i=1}^n \left(y_i - \sum_{j=0}^{p-1} x_{ij} c_j \hat{\beta}_j \right)^2$$

subject to the constraints $c_j \geq 0$, $j=0, 1, 2, \dots, p-1$

$$\& \sum_{j=0}^{p-1} c_j \leq s$$

Where s is some specified positive constant.

for $s > p$, the choice $c_j = 1, j = 0 \dots p-1$

gives $\tilde{\beta}_j = \hat{\beta}_j$, unconstrained minimum

As s is reduced, the shrinkage coefficients get smaller and some are even forced to zero

Thus, garrote is regarded as a compromise between the ridge and the subset selection

- The problem of finding the shrinkage coefficients $\tilde{\beta}_j$ can be reduced to a standard LS problem (Constrained)

- write $\|y - \hat{x}\beta\|^2 = \|y - \hat{x}\beta + \hat{x}\beta - \hat{x}\beta\|^2$

$$= \|y - \hat{x}\beta\|^2 + \|\hat{x}\beta - \hat{x}\beta\|^2 + 2(y - \hat{x}\beta)'(\hat{x}\beta - \hat{x}\beta)$$

$$= RSS + \|\hat{x}\beta - \hat{x}\beta\|^2 + 0$$

$$\text{Let } a_{ij} = x_{ij} \hat{\beta}_j, A = ((a_{ij})), \underline{c} = (c_0, \dots, c_p)'$$

$$\hat{x}\beta = Ac, \underline{d} = \hat{x}\beta$$

$$\therefore \|y - \hat{x}\beta\|^2 = RSS + \|\underline{d} - Ac\|^2 \quad (*)_2$$

the vector \underline{c} which minimizes $(*)_2$, subject to the constraint $c_j \geq 0, j = 0 \dots p-1$ & $\sum c_j \leq s$ is the same as the vector \underline{c} that minimizes $(*)_2$, so

$$\|Ac - \underline{d}\|^2 \text{ subject to the same constraints}$$

A version of garrote which drops the requirement that the shrinkage coefficients be non-negative satisfy $\sum \beta_j^2 \leq s$ similar to the ridge

Breiman : non-garrote $\rightarrow \sum c_j \leq s$
garrote $\rightarrow \sum \beta_j^2 \leq s$

LASSO

Date: / /

consider a different estimator which minimizes

$$\sum_i (\gamma_i - \mathbf{x}_i' \beta)^2 + \lambda \sum_{j=0}^p |\beta_j|, \quad \lambda \rightarrow \text{tuning parameter}$$

The only difference from RR being that absolute values instead of squares are used in the penalty ()

- The change in penalty is small, but the impact on the estimator is huge
- Like RR, penalizing the absolute values of the coefficients introduces shrinkage.
- However, unlike RR, some of the coefficients are shrunk all the way to zero
 - such solutions with multiple values that are identically equal to zero, are said to be sparse
- The penalty thereby performs a sort of variable selection
The resulting estimator was thus named LASSO for "Least Absolute Shrinkage and Selection Operator"

Note: centered and scaled model: $\mathbf{Y} = \mathbf{X}_S \boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \mathbf{X}_S = (\mathbf{1}, \mathbf{X}^*)$

if $\boldsymbol{\beta}$ has a $N_p(\mathbf{m}, \sigma^2 \mathbf{V})$ prior distribution $\boldsymbol{\beta} = (\beta_0, \mathbf{v})'$

then the posterior mean of $\boldsymbol{\beta}$ is $(\mathbf{X}_S' \mathbf{X}_S + \mathbf{V}^{-1}) \mathbf{V}^{-1} \mathbf{m} + \mathbf{X}_S' \boldsymbol{\epsilon}$

if $\mathbf{m} = \mathbf{0}$ and $\mathbf{V} = \begin{pmatrix} \mathbf{C}^{-1} & \mathbf{0} \\ \mathbf{0} & k^{-1} I_p \end{pmatrix}$

the posterior mean (the Bayes estimate) of $\boldsymbol{\beta}$ is exactly the ridge estimate

Generalized Linear Model (GLM)

GLMs extend ordinary regression models to encompass non-normal response distributions and modelling functions of the mean

They have three components:

1. A random component, identifies the response variable \mathbf{Y} and its probability distribution
2. A systematic component, specifies explanatory variables used in a linear predictor (\mathbf{X})

3. A link (), specifies the function $E(Y)$ that the model relates to the linear predictor

- The random component of a GLM consists of a response variable y with independent observations (y_1, y_2, \dots, y_n) from a distribution in the natural exponential family:

$$f(y_i | \theta_i) = a(\theta_i) b(y_i) \exp\{v_i \varphi(\theta_i)\}$$

- the value of the parameter θ_i varies for $i=1(1)n$ as a function of the values of the explanatory variables.

- the parameter $\varphi(\theta)$ is called the natural parameter

- the systematic component of a GLM relates a vector $(\eta_1, \eta_2, \dots, \eta_n)$ to the explanatory variable through a linear model

- let x_{ij} denote the value of explanatory variable j ($j=0(1)p-1$) for subject i

$\eta_i = \sum_{j=0}^{p-1} x_{ij} \beta_j ; i=1(1)n$ is called the linear predictor

- usually $x_{0i}=1$ + i , representing the coefficient of an intercept term β_0 in the model

- the third component of a GLM is a link function that connects the Random component and the systematic component

- let $\mu_i = E(Y_i) ; i=1(1)n$

- the model links μ_i to η_i by $\eta_i = g(\mu_i)$ where the link function g is a monotonic, differentiable function

$$g(\mu_i) = \sum_{j=0}^{p-1} x_{ij} \beta_j ; i=1(1)n$$

- $g(\mu) = \mu$ + identity function $\eta_i = \mu_i$

→ specifies linear model for the mean itself

→ link for ordinary regression with normally distributed y

- the link function that transforms the mean to the natural parameter is called the canonical link

- Date: / /
- We present the "success" and "failure" outcomes by 1 and 0.
 - A Bernoulli trial has possibilities: $P(Y=1) = \pi$, $P(Y=0) = 1-\pi$
 $\therefore E(Y) = \pi$

The probability mass function $f(y; \pi) = \pi^y (1-\pi)^{1-y}$
 $= (1-\pi) \left(\frac{\pi}{1-\pi}\right)^y$
 $\Rightarrow f(y; \pi) = (1-\pi) \exp\{y \log\left(\frac{\pi}{1-\pi}\right)\}, \text{ for } y=0, 1$
 $\theta = \pi$, $a(\pi) = 1-\pi$, $\varphi(\pi) = \log\left(\frac{\pi}{1-\pi}\right)$; $b(y) = 1$

The natural parameter $\log\left(\frac{\pi}{1-\pi}\right)$ is the log odds of response outcome 1, the logit of π

This is the canonical link function

↳ GLM with logit link are referred to as logistic regression models or logit models

Poisson loglinear model for count data

$$f(y; \mu) = \frac{e^{-\mu} \mu^y}{y!} = \frac{1}{y!} \exp\{-\mu + y \log(\mu)\}$$

$$\therefore \theta = \mu, a(\theta) = e^{-\mu}; b(y) = \frac{1}{y!}, \varphi(\theta) = \log(\mu)$$

So the natural parameter is $\log \mu$ and canonical link function is the log link, $\eta = \log(\mu)$

$$\log(\mu_i) = \sum_{j=0}^{p-1} \alpha_{ij} \beta_j; i = 1, n$$

↳ Poisson log linear model

Deviance of a GLM / Deviance function:

- for a particular GLM with observations (y_1, \dots, y_n) , let $L(\lambda; \mathbf{y})$ denote the log-likelihood () expressed in terms of the mean
- let $\hat{L}(\hat{\lambda}; \mathbf{y})$ denote the maximum of the log-likelihood for the model
- consider for all possible models, the maximum achievable log-likelihood is $L(\lambda; \mathbf{y})$
 \hookrightarrow this occurs for the most general model having a separate parameter for each observation and

↳ called the "saturated model"
(contains as many parameters as there are data points)
Date:

- this is not ^a useful model as it doesn't provide any reduction
 - however, it serves as a baseline for comparison with other model fits
- The deviance () is defined/based on likelihood ratio as:

$$D = -2 \log_e \left[\frac{\text{likelihood of the fitted model}}{\text{likelihood of the saturated model}} \right]$$

$$\Rightarrow D = -2 \log_e [L(\hat{\theta}, \bar{y}) - L(\theta^*, \bar{y})]$$

This is the LR statistic for testing the null hypothesis that the model under consideration holds against the general alternative (i.e., the saturated model). Deviance is used for model checking and for inferential comparison of models.

GLM for Binary Data

y : binary response variable (success, failure)

Each obs has one of the two outcomes, denoted by 1 or 0. → we call it trial for a single Bernoulli trial

$E(y) = P(Y=1) = \pi(\underline{x})$, say, reflecting its value of $\underline{x} = (x_1, \dots, x_n)$ of expt variable.

$$\text{Var}(Y) = \pi(\underline{x})(1-\pi(\underline{x}))$$

Linear Probability Model

$$\pi(\underline{x}) = \alpha + \beta_1 x_1 + \dots + \beta_n x_n.$$

$$E(Y|\underline{x}) \quad \hookrightarrow \text{a linear prob model}$$

↪ with independent obs, it is a GLM with Binomial PC and identity link.

The model has a structural problem

$\pi(\underline{x}) \in [0, 1]$, whereas the linear predictor takes values in the entire real line.

$\pi(\underline{x}) < 0$ and/or $\pi(\underline{x}) > 1$ for some values of \underline{x} .

so the model is valid on a very restricted region of \underline{x} values.

β_j is the change in $\pi(\underline{x})$ for one unit change in x_j

↪ constant variance

so the OLS's will not be optimal. So ML estimate is more efficient than OLS's.

→ Mostly unsatisfactory $\hat{\pi}(x)$ as an estimator.

Also, y , being binary, is very far from normality, so the usual t test, F test sampling distn for LSE's do not apply.

Logistic Regression Model.

- Usually binary data result from a unknown relationship between $\pi(x)$ and x .
 - A fixed change in x often has less impact when $\pi(x)$ is near 0 or 1, than when $\pi(x)$ is near 0.5.
- In practice, non-linear relationship between $\pi(x)$ and x are often monotone with $\pi(x)$ increasing continuously or decreasing continuously as x decreases.
 - The S-shaped curve are typically observed



The important curve with this shape is the model

$$\pi(x) = \frac{e^{ax+b}}{1+e^{ax+b}}$$

↪ Logistic regression model with one exp variable.

As $x \uparrow$, $\pi(x) \uparrow$ when $a > 0$ and \downarrow when $a < 0$.

$$\text{General Form: } Y_i | x_i \sim \text{Bernoulli}(\pi_i), \quad \pi_i = \pi(x_i) = P(Y_i=1 | x_i)$$

$$\text{logit } (\pi_i) = \log \frac{\pi_i}{1-\pi_i} = \beta_0 + \beta_1 x_{i1} + \dots + \beta_n x_{in}$$

$$= \beta_0 + \beta' x_i'$$

$$\beta = (\beta_0, \dots, \beta_n)'$$

$\beta_0 \rightarrow$ intercept $\beta \rightarrow$ vector of parameters
corresponding to the independent variables.

$$\gamma_i = \beta_0 + \beta' x_i$$

In general all odds are defined as the ratio of probability of $Y=1$ to that of $Y=0$.

Note: To compute $\beta' x_i$, it requires x_i to be a vector of numeric values.

- It does not directly apply to categorical variables.

- For a binary covariate, we can represent the differential effect asserted by the two levels using a binary indicator or dummy variable, taking values 0 and 1.

- For categorical variable with k levels ($k \geq 2$), we may designate one level as a reference and use $k-1$ binary indicators to represent the individual difference for each of the remaining $k-1$ levels.

the name of the logistic regression model.

$$F(x) = \frac{\exp(x)}{1 + \exp(x)} \quad -\infty < x < \infty$$

\hookrightarrow coeff of the standard logistic r.v.

- If y_i is the no of subjects with the event and m_i is the size of the stratum for the i th unique value of x_i , then $y_i \sim \text{Binomial}(m_i, \pi_i)$, π_i is the prob of success.

Binomial GLM for 2x2 contingency table

- For a \Rightarrow binary response, we have a single exp variable that is also categorical. Values are 0 and 1.

For a given link fn.

$$\text{link}[\pi(x)] = \alpha + \beta x.$$

as the effect of x is described as

$$\beta = \text{link}[\pi(1)] - \text{link}[\pi(0)]$$

y	0	1
0	n_{00}	n_{01}
1	n_{10}	n_{11}

- For the identity link, $\beta = \pi(1) - \pi(0)$ is the difference between proportions.

For log link, $\beta = \log\left(\frac{\pi(1)}{\pi(0)}\right)$ is the log of relative risk.

For logit link, $\beta = \text{logit}(\pi(1)) - \text{logit}(\pi(0))$

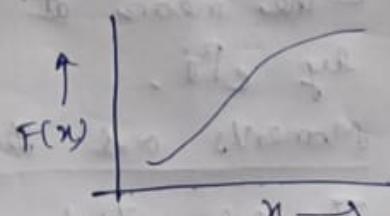
$$= \log \frac{\pi(1)}{1 - \pi(1)} - \log \frac{\pi(0)}{1 - \pi(0)}$$

$$= \log \left[\frac{\pi(1) / (1 - \pi(1))}{\pi(0) / (1 - \pi(0))} \right] \rightarrow \text{is the } \underline{\text{log odds ratio.}}$$

Biject and inverse of cdf link function

- This suggests a model for binary response having form

$$\pi(x) = F(x) \text{ for some cdf}$$



$$\pi(x) = \Phi(\alpha + \beta x)$$

when Φ is strictly increasing over the entire real line,
its inverse exists and

$$\Phi^{-1}(\pi(x)) = \alpha + \beta x.$$

So the link fn for the GLM is Φ^{-1} .

- This curve has the shape of normal cdf when Φ is $N(0,1)$ → called probit model.
- when $\beta > 0$, the logistic regression curve is a cdf for the logistic distri.
- when $\beta < 0$, the curve for $1 - \pi(x)$ has that appearance.
The cdf of the logistic distri with mean μ and dispersion parameter $\tau^2 > 0$, then

$$F(x) = \frac{\exp[(x-\mu)/\tau]}{1 + \exp[(x-\mu)/\tau]} \quad -\infty < x < \infty.$$

Poisson Log-linear model

$$\log \mu(x) = (\beta_0 + \beta_1 x_1 + \dots + \beta_h x_h)$$

$$\rightarrow \mu(x) = \exp(\beta_0 + \beta_1 x_1 + \dots + \beta_h x_h)$$

$$= e^{\beta_0} (e^{\beta_1})^{x_1} (e^{\beta_2})^{x_2} \dots (e^{\beta_h})^{x_h}$$

- A 1-unit increase in x_j has a multiplicative impact of e^{β_j} .
- The mean at x_{ij+1} equals the mean at x_{ij} multiplied by e^{β_j} .

Moments and likelihood for GLM

The exponential dispersion family

y_1, \dots, y_n are independent with pdf or pmf for y_i :

$$\frac{\partial \ell_i}{\partial \theta_i} = L_i = \left[y_i \theta_i - b(\theta_i) \right] / a(\phi) + c(y_i; \phi)$$

$$\frac{\partial L_i}{\partial \theta_i} = \frac{y_i - b'(\theta_i)}{a(\phi)}, \quad \frac{\partial^2 L_i}{\partial \theta_i^2} = -\frac{b''(\theta_i)}{a(\phi)}.$$

$b'(\theta), b''(\theta) >$ 1st and 2nd derivative of $b(\theta)$.

• Applying general likelihood results

$$E\left[\frac{\partial L}{\partial \theta_i}\right] = 0, \quad -E\left[\frac{\partial^2 L}{\partial \theta^2}\right] = E\left[\left(\frac{\partial L}{\partial \theta}\right)^2\right]$$

This holds under regularity conditions satisfied by the exponential family

$$E\left[\frac{y_i - b'(\theta_i)}{a(\phi)}\right] = 0 \Rightarrow \mu_i = E(y_i) = b'(\theta_i)$$

$$\frac{b''(\theta_i)}{a(\phi)} = E\left[\left(\frac{y_i - b'(\theta_i)}{a(\phi)}\right)^2\right] = \frac{\text{var}(y_i)}{\{a(\phi)\}^2}$$

$$\Rightarrow \text{var}(y_i) = b''(\theta_i) \cdot a(\phi)$$

The first term determines the moments of y_i

when y_i is Poisson

$$f(y_i; \mu_i) = \frac{e^{-\mu_i} \mu_i^{y_i}}{y_i!} = \exp[y_i \log \mu_i - \mu_i - \log y_i!] \\ = \exp[y_i \theta_i - \exp(\theta_i) - \log y_i!]$$

When $\theta_i = \log \mu_i$, $b(\theta_i) = \exp(\theta_i)$, $a(\phi) = 1$

$$E(y_i) = b'(\theta_i) = \exp(\theta_i) = \mu_i \quad \therefore c(y_i; \phi) = -\log y_i!$$

$$\text{var}(y_i) = b''(\theta_i) = \exp(\theta_i) = \mu_i$$

• Suppose $y_{ij} \sim \text{Bin}(n_{ij}, p_{ij})$, p_{ij} is the sample proportion of successes.

$$l(\beta) = \sum_{i=1}^n l_i = \sum_{i=1}^n \ln a(\phi) - \frac{1}{a(\phi)} + \sum_{i=1}^n c(y_i, \phi).$$

$l(\beta)$ depends on β through θ_i

Likelihood eq. $\frac{\partial l(\beta)}{\partial \beta_j} = \sum_{i=1}^n \frac{\partial l_i}{\partial \beta_j} = 0 + j.$

$$\frac{\partial l_i}{\partial \beta_j} = \frac{\partial l_i}{\partial \theta_i} \cdot \frac{\partial \theta_i}{\partial \mu_i} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_j}$$

since $\frac{\partial l_i}{\partial \theta_i} = [y_i - l^*(\theta_i)]/a(\phi)$, $\mu_i = l^*(\theta_i)$

$$\frac{\partial l_i}{\partial \theta_i} = \frac{y_i - \mu_i}{a(\phi)}, \quad \frac{\partial \mu_i}{\partial \theta_i} = l''(\theta_i) = \text{var}(y_i)/a(\phi)$$

$$\eta_i = \sum \beta_j x_{ij}, \quad \frac{\partial \eta_i}{\partial \beta_j} = x_{ij}$$

and finally $\eta_i = g(\mu_i)$, $\frac{\partial \mu_i}{\partial \eta_i}$ depends on the link function

$$\begin{aligned} \frac{\partial l_i}{\partial \beta_j} &= \frac{y_i - \mu_i}{a(\phi)} \cdot \frac{a(\phi)}{\text{var}(y_i)} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot x_{ij} \\ &= \frac{(y_i - \mu_i)x_{ij}}{\text{var}(y_i)} \cdot \frac{\partial \mu_i}{\partial \eta_i} \end{aligned}$$

Summing over all obs.

$$\sum \frac{(y_i - \mu_i)x_{ij}}{\text{var}(y_i)} \cdot \frac{\partial \mu_i}{\partial \eta_i} = 0, \quad j = 0, 1, 2, \dots$$

$$\mu_i = g^{-1}(\sum \beta_j x_{ij})$$

- The LR equations depend on the dist of y_i only through μ_i and $\text{var}(y_i)$
- The variance depends on the mean through a particular form $\text{var}(y_i) = V(\mu_i)$ for some variance V for $V(\cdot)$

For example $V(\mu_i) = \mu_i$, for Poisson.

$$V(\mu_i) = \mu_i(1-\mu_i)/n_i \text{ for Binomial}$$

$$V(\mu_i) = \sigma^2, \text{ a constant for normal}$$

Likelihood eq for Binomial GLMs

Suppose $\pi_i | y_i \sim \text{Bin}(n_i, \pi_i)$ (y_i is the sample proportion of successes for n_i trials)

$$E(y_i) = \pi_i$$

The Binomial GLM's for single predictor, $[\pi(x) = \Phi(a + bx)]$ extends with several predictors to

$$\pi_i = \Phi\left(\sum_j \beta_j x_{ij}\right), \Phi - \text{standard cdf of some cont-dist}$$

$$\frac{\sum (y_{ij} - \pi_{ij}) x_{ij}}{\pi_i(1-\pi_i)/n_i}, \quad \pi_i = \mu_i = \Phi(\eta_i)$$

$$\frac{\partial \mu_i}{\partial \eta_i} = \phi(\eta_i) = \phi(\sum_j \beta_j x_{ij})$$

$$\phi(u) = \Theta \frac{d\Phi(u)}{du}$$

Since $\text{var}(y_i) = \pi_i(1-\pi_i)/n_i$, the likelihood eq simplifies to

$$\frac{\sum n_i (y_{ij} - \pi_{ij}) x_{ij}}{\pi_i(1-\pi_i)} \phi\left(\sum_j \beta_j x_{ij}\right) = 0, \quad \pi_i = \Phi\left(\sum_j \beta_j x_{ij}\right)$$

$j=0, 1, \dots$

For the logit link,

$$\Rightarrow \sum_{j=0}^n (y_i - \mu_i) x_{ij} = 0, \quad j=0, 1, \dots, n.$$

Likelihood eq. for Poisson log link model

$$Y_i \sim \text{Poi}(\mu_i)$$

$$E(Y_i) = \mu_i$$

$$\text{Var}(Y_i) = \mu_i$$

$$\sum_j \frac{(y_i - \mu_i) x_{ij}}{\text{Var}(y_i)} \cancel{\frac{\partial \mu_i}{\partial \eta_i}} = 0$$

$$\Rightarrow \sum_j \frac{(y_i - \mu_i) x_{ij}}{\mu_i} \cancel{\frac{\partial \mu_i}{\partial \eta_i}} = 0 \Rightarrow \sum_j (y_i - \mu_i) x_{ij} = 0$$

$$\eta_i = \log \mu_i = \sum_j \beta_j x_{ij}$$

$$\Rightarrow \frac{\partial \eta_i}{\partial \mu_i} = \frac{1}{\mu_i} \Rightarrow \frac{\partial \mu_i}{\partial \eta_i} = \mu_i$$

Asymptotic Covariance Matrix of Model Parameter Estimates

- The likelihood fn for the GLM also determines the asymptotic covariance matrix of the ML estimator $\hat{\beta}$.

$$\text{Information matrix } I = \left(E \left[\frac{\partial^2 L}{\partial \beta_i \partial \beta_j} \right] \right)$$

L_i : contribution from the i th sample fit to the log likelihood.

- For the exponential family, we have

$$E \left[\frac{\partial^2 L_i}{\partial \beta_i \partial \beta_j} \right] = -E \left(\frac{\partial L_i}{\partial \beta_i} \right) \left(\frac{\partial L_i}{\partial \beta_j} \right)$$

$$E \left[\frac{\partial^2 l(\beta)}{\partial \beta_k \partial \beta_j} \right] = \sum_{i=1}^n \frac{x_{ik} x_{ij}}{\text{var}(y_i)} \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2$$

- Let W be a diagonal matrix with i th diagonal entry

$$w_i = \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2 / \text{var}(y_i)$$

Then $I = X^T W X$

↳ the form of w and hence I depend on the link function.

- The asymptotic variance-covariance matrix of $\hat{\beta}$ is estimated by: $\hat{\text{cov}}(\hat{\beta}) = I^{-1} = (X^T \hat{W} X)^{-1}$

where \hat{W} is W evaluated at $\hat{\beta}$.

In case of Poisson log linear model, the estimated asymptotic covariance matrix will be

$$w_i = \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2 / \text{var}(y_i) = \mu_i$$

$$\log \mu_i = \sum_j \beta_j x_{ij}$$

Deviance & Goodness of fit

- The standard GLM has a separate parameter for each observation. It gives a perfect fit.
- Sounds good, but not a helpful model.
- Does not smooth data and hence no advantage of having a simpler model.
- It serves as a baseline for other models, for example, for checking model fit.

Let, $\hat{\theta}$ denote the estimate of θ , for the saturated model, corresponding is estimated means $\tilde{\mu}_i = y_i \forall i$.

- For a particular unsaturated model, denote the corresponding ML estimators by $\hat{\theta}$ and $\hat{\mu}_i$.

- $2(\hat{\mu}; y)$: maximized log likelihood for the unsaturated model.
- $2(\hat{\mu}_i; Y)$: maximized log likelihood for the saturated model.

- $-2 [2(\hat{\mu}; y) - 2(\hat{\mu}_i; Y)]$ describes the lack of fit

- LR statistic for testing the null hypothesis that the model holds against the alt fact a more general model holds.

$$\begin{aligned} & -2 \left[l(\hat{\mu}; \mathbf{y}) - l(\boldsymbol{\gamma}; \mathbf{y}) \right] \\ &= 2 \sum \left[y_i \tilde{\phi}_i - b(\tilde{\phi}_i) \right] / \alpha(\phi) - 2 \sum \left[y_i \hat{\phi}_i - b(\hat{\phi}_i) \right] / \alpha(\phi) \end{aligned}$$

When $\alpha(\phi) = \phi/\omega_i$ this is equal to

$$\begin{aligned} &= 2 \sum_{i=1}^n \omega_i \left[y_i (\tilde{\phi}_i - \hat{\phi}_i) - b(\tilde{\phi}_i) + b(\hat{\phi}_i) \right] / \phi \\ &= \frac{D(\boldsymbol{\gamma}; \hat{\mu})}{\phi} \end{aligned}$$

is called the scaled deviance & $D(\boldsymbol{\gamma}; \hat{\mu})$ is the deviance.

Deviance for Poisson GLMs

$$\hat{\phi}_i = \log \hat{\mu}_i$$

$$b(\hat{\phi}_i) = \exp(\hat{\phi}_i) = \hat{\mu}_i$$

$$\tilde{\phi}_i = \log y_i$$

$$b(\tilde{\phi}_i) = y_i$$

$\therefore \alpha(\phi) = 1$, so the deviance and scaled deviance are equal and

$$D(\boldsymbol{\gamma}; \hat{\mu}) = 2 \sum_{i=1}^n \left[y_i \log \left(\frac{y_i}{\hat{\mu}_i} \right) - y_i + \hat{\mu}_i \right]$$

- When a model with log link contain an intercept term

The likelihood eq of Poisson GLM with log link

$$\sum (y_{ij} - \mu_{ij}) x_{ij} = 0$$

$$\sum (y_{ij} - \mu_{ij}) = 0 \quad \text{corresponding to intercept.}$$

$$\sum y_{ij} = \sum \hat{\mu}_{ij}$$

$$D(\underline{y}; \hat{\mu}) = -2 \sum [y_{ij} \log(y_{ij}/\hat{\mu}_{ij})]$$

$$a(\phi) = y_{ni} \Rightarrow \phi = 1 \quad w_i = n_i \quad [\text{Binomial GLM}]$$

Likelihood Ratio Model Comparison using deviances

for a poisson / binomial model denoted by μ , $\phi = 1$, so

$$D(\underline{y}; \hat{\mu}) = -2 [l(\hat{\mu}; \underline{y}) - l(\underline{y}; \underline{\mu})]$$

Consider two models M_0 with fitted values $\hat{\mu}_0$, & M_1 with fitted values $\hat{\mu}_1$, with M_0 a special case of μ

Since M_0 is simpler than μ , a smaller set of parameter values satisfies M_0 than satisfy μ ,

- maximizing the log-likelihood over a smaller space cannot yield a larger maxima

$$l(\hat{\mu}_0; \underline{y}) \leq l(\hat{\mu}; \underline{y})$$

- with the same model for each μ_0 &

$$\hat{\mu}_1 [l(\hat{\mu}; \underline{y})]$$

$$D(\underline{y}; \hat{\mu}_1) \leq D(\underline{y}; \hat{\mu}_0)$$

→ simpler models have larger deviances.

- Assuming that model μ_0 holds, the LR test of hypothesis that μ_0 holds uses the test statistic.

$$= -2 [l(\hat{\mu}_0; \underline{y}) - l(\hat{\mu}; \underline{y})]$$

$$= -2 [l(\hat{\mu}_0; \underline{y}) - l(\underline{y}; \underline{y})] - \{-2 [l(\hat{\mu}_1; \underline{y}) - l(\underline{y}; \underline{y})]\}$$

$$= D(\underline{y}; \hat{\mu}_0) - D(\underline{y}; \hat{\mu}_1)$$

- The LR statistic comparing two models is simply the difference between the deviance.

- This is large when μ_0 fits poorly compared to μ_1 .

$$D(\underline{y}; \hat{\mu}_0) - D(\underline{y}; \hat{\mu}_1) = \hat{\mu}_0 \rightarrow \hat{\theta}_{0i} \\ \hat{\mu}_1 \rightarrow \hat{\theta}_{1i}$$

$$= 2 \sum_i w_i [y_i (\hat{\theta}_{1i} - \hat{\theta}_{0i}) - b(\hat{\theta}_{1i}) + b(\hat{\theta}_{0i})]$$

approx χ^2_{df} df: difference between the no. of

parameters in the two models.

under certain

regularity conditions.

for Poisson likelihood log linear model with

$$D(\underline{y}; \hat{\mu}) = 2 \sum_i y_i \log \left(\frac{y_i}{\hat{\mu}_i} \right)$$

$$D(\underline{y}; \hat{\mu}_0) - D(\underline{y}; \hat{\mu}_1) = 2 \sum_i y_i \log \left(\frac{\hat{\mu}_1}{\hat{\mu}_0} \right)$$

NR: $\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} = H^{(t)}^{-1} \underline{u}^{(t)}$

assuming that the Hessian matrix $H^{(t)}$ is ns.

Fisher scoring $\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} + (I^{(t)})^{-1} \underline{u}^{(t)}$

$I^{(t)}$: information matrix, for $\hat{\beta}^{(t)}$

($\hat{\beta}$ is measured replaced by $\hat{\beta}^{(t)}$)

$\underline{u}^{(t)}$: sets of first order derivative evaluated at $\hat{\beta}^{(t)}$

$I = x'w x$, w is diagonal with w_i

$$w_i = \left(\frac{\partial \mu_i}{\partial \eta_i} \right)^2 / v(y_i)$$

$I^{(t)} = x'w^{(t)}x$, $w^{(t)}$ evaluated at $\hat{\beta}^{(t)}$.

for both Fisher scoring & MR algorithm, the score \underline{u} has the form

$$\underline{u}_j = \frac{\partial l(\beta)}{\partial \beta_j} = \sum_{i=1}^n \frac{(y_i - \mu_i) x_{ij}}{v(y_i)} \frac{\partial \mu_i}{\partial \eta_i}$$

Define $D = \text{diag}\{d_1, \dots, d_n\}$, $d_i = \frac{\partial \mu_i}{\partial \eta_i}$

The GLMs likelihood eqn can be written as

$$\underline{u} = x'w D^{-1} (\underline{y} - \underline{\mu}) = 0$$

Consider maximizes the loglikelihood based on obsev y from $\text{Bin}(n, \pi)$. distn.

$$l(\pi) = y \log \pi + (n-y) \log(1-\pi) \quad [\text{Take } \pi \text{ as the unknown parameter}]$$

$$u = \frac{(y-n\pi)}{\pi(1-\pi)} ; \rightarrow \infty$$

$$H = - \left[\frac{n}{\pi^2} + \frac{n-y}{(1-\pi)^2} \right]$$

Each NR step has the form

$$\pi^{(t+1)} = \pi^{(t)} + \left[\frac{y}{(\pi^{(t)})^2} + \frac{n-y}{(1-\pi^{(t)})^2} \right]^{-1} \frac{y-n\pi^{(t)}}{\pi^{(t)}(1-\pi^{(t)})}$$

- (i) in positive position when

$$y/n > \pi^{(t)}$$

$\pi^{(t+1)}$ increases if $y/n > \pi^{(t)}$ & decreases

when $y/n < \pi^{(t)}$.

$$\text{When } \pi^{(0)} = 1/2; \quad \pi^{(1)} = \frac{y}{n}$$

(with $\pi^{(t)} = y/n$; no adjustment occurs &

$$\pi^{(t+1)} = y/n \text{ which is the correct answer.}$$

for starting values other than 0.5; adequate

(ii) convergence usually happens in just a few more

(iii) steps.

The information in this case.

so the steps for Fisher scoring stages

$$\pi^{(t+1)} = \pi^{(t)} + \left[\frac{n}{\pi^{(t)} + (1-\pi^{(t)})} \right]^{-1} \frac{y-n\pi^{(t)}}{\pi^{(t)}(1-\pi^{(t)})}$$

$$= \pi^{(t)} + \frac{y-n\pi^{(t)}}{n} = y/n$$

This gives the correct answer for $\hat{\pi}$ after a single iteration and stays at that value for successive iterations.

ML & as Iterative Reweighted Least Square

$$\underline{z} = \underline{x}\beta + \underline{\epsilon}; v(\underline{\epsilon}) = V$$

The WLS of β is $(\underline{x}'V^{-1}\underline{x})^{-1}\underline{x}'V^{-1}\underline{z}$

$$T = \underline{x}'W\underline{x}$$

from Fisher Scoring

$$T^{(t)} \hat{\beta}^{(t+1)} = T^{(t)} \hat{\beta}^{(t)} + \underline{\mu}^{(t)}$$

$$= (\underline{x}'W^{(t)}\underline{x}) \hat{\beta}^{(t)} + \underline{x}'W^{(t)}(\underline{\mu}^{(t)})^{-1}(\underline{y} - \underline{\mu}^{(t)})$$

$$= \underline{x}'W^{(t)} \left[\hat{\beta}^{(t)} + D^{(t)}^{-1}(\underline{y} - \underline{\mu}^{(t)}) \right]$$

$$= \underline{x}'W^{(t)} \underline{z}^{(t)}$$

When $\underline{z}^{(t)}$ has element

$$z_i^{(t)} = \sum_j x_{ij} \beta_j^{(t)} + (y_i - \mu_i^{(t)}) \frac{\partial \mu_i^{(t)}}{\partial \mu_i^{(t)}}$$

$$= n_i^{(t)} + (y_i - \mu_i^{(t)}) \frac{\partial \mu_i^{(t)}}{\partial \mu_i^{(t)}}$$

$$T^{(t)} \hat{\beta}^{(t+1)} = T^{(t)} \hat{\beta}^{(t)} + \underline{\mu}^{(t)}$$

$$\Rightarrow (\underline{x}'W^{(t)}\underline{x}) \hat{\beta}^{(t+1)} = \underline{x}'W^{(t)}\underline{z}^{(t)}$$

$$\hat{\beta}^{(t+1)} = (\underline{x}'W^{(t)}\underline{x})^{-1} \underline{x}'W^{(t)}\underline{z}^{(t)}$$

The vector $\underline{z}^{(t)}$ is an estimated linearized form of the link f^h , evaluated at y_i .

$$g(y_i) \approx g(\mu_i^{(t)}) + (y_i - \mu_i^{(t)}) g'(\mu_i^{(t)})$$

$$= n_i^{(t)} + (y_i - \mu_i^{(t)}) \frac{\partial \mu_i^{(t)}}{\partial \mu_i^{(t)}} = z_i^{(t)}$$

• At cycle t : regresses $\underline{z}^{(t)}$ on x with weight $w^{(t)}$

- obtain new estimate $\hat{\beta}^{(t+1)}$

- new linear predictor $\underline{\eta}^{(t+1)} = x \hat{\beta}^{(t+1)}$

- new adjusted $\underline{z}^{(t+1)}$

• The ML estimator results from iterative use of WLS in which the weight matrix changes at each cycle.

→ ② is called iterative reweighted LS.

$$u_i^{(0)} = y_i - \text{Bernoulli response} \quad i=1, \dots, n$$

logit

Poisson response, $M_i = Y_i > 0$

$y_i \geq 0$, we start $y_i = 1$

$y_i \geq 0, 1, \dots, Y_i$

Residuals for GLM

• Pearson residual • for obj i

$$\epsilon_i = \frac{y_i - \hat{\mu}_i}{\sqrt{v(u_i)}} ; v(y_i) = v(\mu_i)$$

Deviance residual:

$$D(y; \hat{\mu}) = 2 \sum_i w_i [y_i (\tilde{\theta}_i - \hat{\theta}_i) - b(\tilde{\theta}_i) + b(\hat{\theta}_i)]$$

$$= \sum_i d_i$$

Deviance residual is $\sqrt{d_i} \operatorname{sign}(y_i - \hat{\mu}_i)$

↳ sum of d_i is the deviance.

$\tilde{\theta}_i \rightarrow$ saturated
 $\hat{\theta}_i \rightarrow$

Logistic Regression:

Interpreting Parameters in Logistic Regression

- for binary response variable y and an expl variable x , let, $\pi(x) = P(y=1|x=x) = 1 - P(y=0|x=x)$

The logistic regression model

$$\pi(x) = \frac{\exp(\alpha + \beta x)}{1 + \exp(\alpha + \beta x)} \quad \textcircled{1}$$

or logit has the linear relationship.

$$\text{logit} [\pi(x)] = \log \frac{\pi(x)}{1 - \pi(x)} = \alpha + \beta x \quad \textcircled{2}$$

Q: How can we interpret β .

- Its sign determines whether $\pi(x) \uparrow$ or \downarrow as x increases.
- The rate of climb or descent increases on $|\beta|$ increases as $\beta \rightarrow 0$, the curve flattens to a horizontal st. line.
- When $\beta = 0$; y is independent of x
- for quantitative x with $\beta > 0$, the curve for $\pi(x)$ has the slope of the CDF of the logistic distn.
- Some logistic density is symmetric, $\pi(x)$ approaches 1 at the same rate as it approaches to zero.

$$\frac{\pi(x)}{1-\pi(x)} = e^{\alpha + \beta x}$$

The odds are exponential for $\alpha + \beta x$

The odds multiply by e^β for unit increase in x .

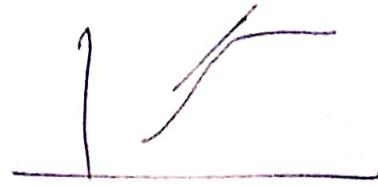
In other words, e^β is an odds ratio.

\rightarrow odds at x : $\beta = x+1$ / odds at $x=0$

Since it has a curved rather than a linear appearance, the logistic regression fn ①

\Rightarrow the rate of change in $\pi(x)$ per unit

change in x varies



$$\frac{d\pi(x)}{dx} = \frac{\beta e^{\alpha + \beta x}}{(1 + e^{\alpha + \beta x})^2}$$

$$= \beta \pi(x) (1 - \pi(x))$$

The line tangent to the curve at x for which $\pi(x)=1/2$ has slope. $\beta/9$

Slope $\rightarrow 0$ as $\pi(x) \rightarrow 0$ or 1.0

When $\pi(x) = 0.9$ or 0.1 ; slope $= 0.09\beta$

Slope $\rightarrow 0$ as $\pi(x) \rightarrow 1.0$ or 0

The steepest slope occurs at x for which $\pi(x)=1/2$

$$\frac{\pi(x)}{1-\pi(x)} = e^{\beta} (\alpha + \beta x) \Rightarrow x = -\alpha/\beta.$$

This x value ($-\alpha/\beta$) is sometimes called median effective level.

In toxicology it is called lethal dose.

[The dose with 50% chance of a lethal result].

[END]

Categorial

data analysis

→ Book (for GIM)

(D) \rightarrow categorical data with a correspondence analysis (CAT) approach to also test the hypothesis.

→ χ^2 test for independence of variables

→ χ^2 test for independence of variables