[quick background on Logit]　　Logit Model

Consider a dichotomous and response variables. $y$ with two measurement levels.

Let $\pi(x) = P(y=1 \mid X=x) = (1-P(y=0 \mid X=x)$

odds ratio $\dfrac{\pi(x)}{1-\pi(x)}$

$\log(\text{odds ratio}) = \text{Logit}[\pi(x)] = \log\left(\dfrac{\pi(x)}{1-\pi(x)}\right)$

$$= \alpha + \beta x$$

Thus the odds $= \exp(\alpha + \beta x)$

↓　↓

•) Rate of increase or decrease of the S shaped curve of $\pi(x)$

') The sign of $\beta'$ indicates whether curve ascends $(\beta > 0)$ or descends $(\beta < 0)$

**Multiple Logit Model.**

Let $k'$ denote nos of predictors for a binary response. $y'$ by $x_1, x_2 \ldots x_k$.

Then we have. $\text{Logit}[P(y=1)] = \alpha + \beta_1 x_1 + \cdots \beta_k x_k$.

or $\dfrac{\pi(x)}{1-\pi(x)} = \alpha + \beta_1 x_1 + \cdots \beta_k x_k$

$$\pi(x) = \dfrac{\exp(\alpha + \beta_1 x_1 + \cdots \beta_k x_k)}{1 + \exp(\alpha + \beta_1 x_1 + \cdots \beta_k x_k)}$$

# Multinomial Logit

Consider the following.

- > 'n' independent observations
- > 'p' explanatory variables.
- > 'k' Categories.

Take any category as base. -> Let us take category 'j' as base.

> $\pi_{ij}$ -> multinomial probability of an 'ith' observation falling in the 'j'th category

Then

$$\eta_{ij} = \log \frac{\pi_{ij}}{\pi_{iJ}} = \alpha_j + x_i' \beta_j.$$

we assume that the log of odds (wrt to base category) follows a *linear* model.

## Linear predictor:

Consider the linear predictor function. $f(k,i)$ to predict the probability that observation 'i' has outcome 'k'

$$f(k,i) = \beta_{0,k} + \beta_{1,k} x_{1,i} + \beta_{2,k} x_{2,i} + \cdots \beta_{M,k} x_{M,i}$$

$\beta_{m,k}$ -> Regression coefficient for 'm'th explanatory variable. and 'k'th outcome

writing more compactly we have.

$$f(k,i) \quad f(k,i) = \beta_k \cdot X_i$$

As a set of independent binary regressions:

For 'k' possible outcomes, run k-1 independent binary Logistics regression models.

one outcome 'say the last' is chosen as pivot.

Thus we have 'k-1' equations as follows.

$$\ln \frac{Pr(y_i = 1)}{Pr(y_i = k)} = \beta_1 \cdot X_i \quad \left] \Rightarrow Pr(y_i = 1) = Pr(y_i = k) e^{\beta_1} \right.$$

$$\vdots$$

$$\ln \frac{Pr(y_i = k-1)}{Pr(y_i = k)} = \beta_{k-1} \cdot X_i$$

Sum of probabilities must = 1: Thus we have.

$$Pr(y_i = k) = 1 - \sum_{k=1}^{k-1} Pr(y_i = k)$$

$$= 1 - \sum_{k=1}^{k-1} Pr(y_i = k) e^{\beta_k \cdot X_i}$$

$$\Rightarrow Pr(y_i = k) = \frac{1}{1 + \sum_{k=1}^{k-1} e^{\beta_k \cdot X_i}} \quad \left] - Ⓐ \right.$$

We can Now use Ⓐ to find other Probabilities.

$$Pr(y_i = 1) = \frac{e^{\beta_1 X_i}}{1 + \sum_{k=1}^{k-1} e^{\beta_k \cdot X_i}}$$

Background : Recap Logistics Regression

$\Pi_i$ = Probability of success For any given obs.

odds ratio = $\dfrac{\Pi_i}{1 - \Pi_i}$

$Logit = \ln\left(\dfrac{\Pi_i}{1-\Pi_i}\right) = \sum_{k=0}^{K} x_{ik} \beta_k$ , $i = 1, 2, \dots N.$
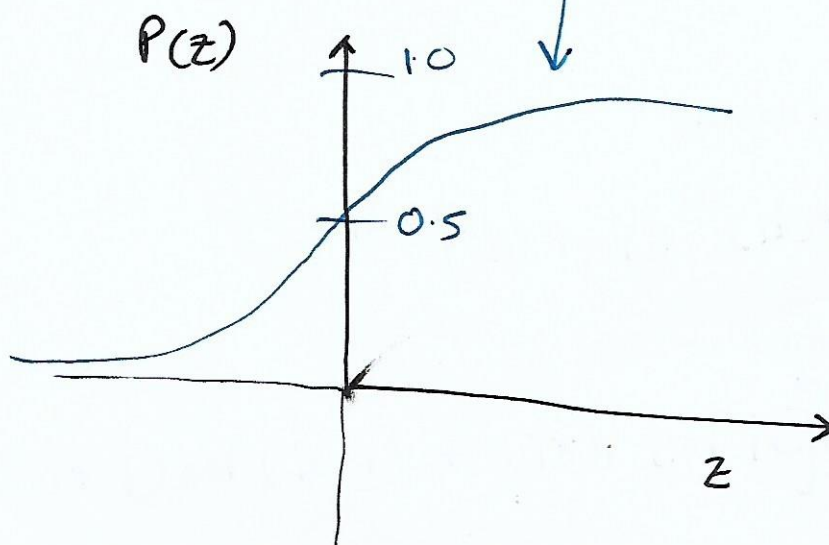
Solving For $\Pi_j = P(x)$

we have $\dfrac{P(x)}{1 - P(x)} = \exp\left(\sum_{j=0}^{k} x_j \cdot \beta_j\right) = \exp(z)$ say.

$\therefore \dfrac{P(x)}{1 - P(x)} = \exp(z)$

$\therefore P(x) = \dfrac{\exp(z)}{1 + \exp(z)}$

or $P(z) = \dfrac{\exp(z)}{1 + \exp(z)} = \dfrac{1}{1 + \exp(-z)}$

Sigmoid Function.

Maps a real line to $(0,1)$



$P(z)$

**A-srd. Proof:** $P'(z) = P(z)(1-P(z))$

we have $P(z) = \dfrac{1}{1-e^{-z}} = (1-e^{-z})^{-1}$

$\therefore P'(z) = -1(1-e^{-z})^{-2}(-1)(-e^{-z})$

$\therefore P'(z) = \dfrac{-e^{-z}}{(1-e^{-z})(1-e^{-z})}$

Now $P(z) = \dfrac{1}{(1-e^{-z})}$ and $1 - P(z)$

$$= \dfrac{1 - \dfrac{1}{1-e^{-z}}}{(1-e^{-z})} = \dfrac{(1-e^{-z})-1}{(1-e^{-z})} = \dfrac{-e^{-z}}{(1-e^{-z})}$$

$\therefore \boxed{P'(z) = P(z)(1-P(z))}$

**Determining the coefficients:** $\Big]$ Maximizing the Log Likelihood.

## Likelihood Function for Logistics Regression. →

· For each training data point · we have a vector of features $x_i$
and an observed class $y_i$

Prob of that class $= p$ if $y_i = 1$, or $1-p$ if $y_i = 0$

← · The likelihood Function is →

: Product of predicted probabilities of the N individual observations

The likelihood is written as

$$L(\beta_0, \beta) = \Pi\, P(x_i)^{y_i}(1-P(x_i))^{1-y_i}$$

The log likelihood is written as

$$ll(\beta_0, \beta) = \sum_{i=1}^{N}\Big[y_i \log P(x_i) + (1-y_i)\log(1-P(x_i))\Big]$$

Maximize the log likelihood.

We take derivative wrt $\beta'$

$$\nabla_b(\ell L) = \sum_{i=1}^{N} y_i \frac{P_i'}{P_i} x_i + \sum_{i=1}^{N} (1-y_i) \frac{P_i' x_i}{1-P_i}$$

Now $P_i' = P_i(1-P_i)$

$$\therefore \nabla_b(\ell L) = \sum_{i=1}^{N} \left[ y_i(1-P_i) - (1-y_i)P_i \right] x_i$$

$\Rightarrow$ Setting $\sum_{i=1}^{N} y_i x_i - P_i x_i = 0$

How do we solve for coefficients.

Say we have a vector valued function. $y = f(b)$.
We want
guess. $(b_0)$ $f(b_{opt}) = 0$. Assume we start with initial

$$f(b_0 + \Delta) \simeq f(b_0) + \Delta f'(b_0) = 0$$

Recap:
Linear Regression

$$\therefore \Delta_0 = -\frac{f'(b_0)}{f(b_0)}$$

$$y = x^T b$$

upgrade rule

$$b_1 = b_0 + \Delta_0$$

$$xy = xx^T b$$
$$b = (xx^T)^{-1} xy$$

Now we have $f = \nabla_b(\ell L) = \sum y_i x_i - P_i x_i = 0$

also $H = \frac{\partial}{\partial b}(\nabla_b \ell\ell) = -\sum x_i (\nabla_b(P_i))$

$$= -\sum x_i P_i(1-P_i) x_i^T$$

Now in matrix form.

$$= XWX^T$$

$$\nabla_b(\ell L) = X(y - P_R)$$

$$\text{or } \left[ \Delta_R = (XW_R X^T)^{-1} X (y - P_R) \right]$$

# Log likelihood For Multinomial case:

Likelihood $L = \prod\limits_{i=1}^{n} \prod\limits_{h=0}^{q} P_{ih}^{y_{ih}}$

$y_{ih} \rightarrow$ observed values

$P_{ih} \rightarrow$ Theoretical Values.

∴ Log likelihood $LL = \sum\limits_{i=1}^{n} \sum\limits_{h=0}^{q} y_{ih} \ln P_{ih}$ — Ⓐ

Let $B_h = \{b_j\} b_{hj}$ be the $(k+1) \times 1$ Col vector.

of binary logistics reg. coefficient of the outcome 'h' compared to '0'

Let $B$ be the $q(k+1) \times 1$ Col vector. consisting of

$B_0 \cdots B_r$ arranged in a column.

Let $X$ be the design matrix $n \times (k+1)$

For outcomes 'h' and 'l' Let $V_{hl}$ be the $n \times n$ diag matrix whose main diag. contains elements of the form.

$V_{ii} = \begin{cases} P_{ih}(1-P_{ih}) & h = l \\ -P_{ih}P_{il} & h \neq l \end{cases}$

Let $C_{hl} = X^T V_{hl} X$   Now define the $n \times n$, $q \times n q$ matrices.

$$C = \begin{bmatrix} C_{11} & \cdots & C_{1q} \\ C_{q1} & & C_{qq} \end{bmatrix}$$

Then $S = C^{-1}$ is the Covariance matrix for B.

For max Ⓐ log likelihood.

We have.

$\sum\limits^{n} (y_{ih} - P_{ih}) = 0$ and $\sum\limits_{i=1}^{n} x_{ij} (y_{ih} - P_{ih}) = 0$

Thus we get the following matrix eq.

$$x^T (y - p) = 0$$

Let $B^0$ be initial guess for $B$

For Each $m$th iter. we have.

$$B^{m+1} = B^{(m)} + S^{-1} x^T (y - p^{(m)})$$

For sufficiently Large $m$.

$$B^{(m+1)} \simeq B^{(m)}$$ is a good approx.

For $\hat{B}$