

Generalized Linear Models

response variable
and one or more
predictor variables.

→ Used to model relationship b/w

GLM analysis comes into

play when the error distribution is not normal

and / or

When a vector of non linear

functions of the responses

$$\eta(y) = (\eta_1(y_1), \eta_2(y_2), \dots, \eta_n(y_n)),$$

is equal to $x\beta$ and no y itself;

has expectation $x\beta$.

Multiple Linear Regression:

$$y = x\beta + \epsilon$$

$$E(y) = x\beta, \quad E(\epsilon) = 0$$

here

$$\boxed{E(\eta(y)) = x\beta.}$$

In GLM, the response variable distribution must be a member of the exponential Family.

The exponential Family of distribution

A random variable u that belongs to the exponential family with a single parameter θ has a probability density function (pdf)

$$f(u, \theta) = s(u) t(\theta) e^{a(u)b(\theta)}$$

Where:

s, t, a, b are all known functions

Rewrite:

$$f(u, \theta) = \exp \{ a(u)b(\theta) + d(u) + c(\theta) \}$$

Where

$$d(u) = \ln(s(u)) \quad c(\theta) = \ln(t(\theta))$$

Where $a(u) = u$, the distribution is said to be in Canonical form &

$b(\theta)$ is called Natural parameter.

Parameters other than the parameter of interest θ are called nuisance parameters.

Some members of the exponential family will be discussed in part.

Family will be discussed in part.

1. Normal distribution $N(\mu, \sigma^2)$

pdf of normal distribution

$$f(u; \mu) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(u-\mu)^2}{\sigma^2}}, -\infty < u < \infty$$

$$= \exp \left\{ u \cdot \frac{\mu}{\sigma^2} + \left[\frac{-\mu^2}{2\sigma^2} - \frac{1}{2} \ln 2\pi\sigma^2 \right] \right\}$$

$$a(u) = u, b(\theta) = \frac{\mu}{\sigma^2} = \mu$$

$$c(\theta) = \frac{-\mu^2}{2\sigma^2} - \frac{1}{2} \ln 2\pi\sigma^2$$

$$d(u) = \frac{u^2}{2\sigma^2}$$

2. Binomial Distribution $\text{Bin}(n, p)$

P - parameter of interest

n - nuisance parameter

$$f(u, p) = \binom{n}{u} p^u (1-p)^{n-u}, \quad u=0, 1, \dots, n$$

$$= \binom{n}{u} \left(\frac{p}{1-p} \right)^u (1-p)^n$$

$$= \exp \left\{ u \ln \left(\frac{p}{1-p} \right) + n \ln(1-p) + \ln(n) \right\}$$

$$a(\theta) = u, \quad b(\theta) = \ln \left(\frac{p}{1-p} \right) = \text{natural Parameter}$$

$$c(\theta) = n \ln(1-p)$$

$$d(\theta) = \ln \left(\frac{n}{u} \right)$$

3. Poisson Distribution

$$f(u, \lambda) = \frac{e^{-\lambda} \lambda^u}{u!}, \quad u=0,1,\dots$$

Probability of getting u successes in λ trials

$$= \exp \{u \ln \lambda - \lambda - \ln u!\}$$

$$a(u) = u, \quad b(\theta) = \ln \lambda$$

↳ Natural

$$c(\theta) = -\lambda$$

↳ Parameter

$$d(u) = -\ln u!$$

4 Gamma distribution

With parameter θ of interest &

α as nuisance parameter

$$f(u, \theta) = \frac{\theta^\alpha u^{\alpha-1} e^{-\theta u}}{\Gamma \alpha}, \quad \alpha, \theta > 0, u \geq 0$$

$$= \exp \left\{ -\theta u + \alpha \ln \theta - \ln \Gamma \alpha + (\alpha-1) \ln u \right\}$$

$$a(u) = u, \quad b(\theta) = -\theta$$

$$c(\theta) = \alpha \ln \theta - \ln \Gamma \alpha$$

$$d(u) = (\alpha-1) \ln u$$

5. Exponential distribution

Probability density function $f(u|\theta) = \theta e^{-\theta u}$, $u > 0, \theta > 0$

$$= \exp \{ -u\theta + \ln \theta \}.$$

$$b(\theta) = \theta \rightarrow \text{natural parameter}$$

$$a(u) = u$$

Parameter

6. Negative Binomial distribution

The variable U is the no. of failures observed to attain r successes in binomial trials with probability of success θ .

Successes θ

Probability mass function can be written as

$$f(u|\theta) = \binom{r+u-1}{r-1} \theta^r (1-\theta)^u,$$

$$u = 0, 1, 2, \dots$$

$$= \exp \left\{ u \ln(1-\theta) + r \ln \theta + \ln \binom{r+u-1}{r-1} \right\}$$

$$a(u) = u, b(\theta) = \ln(1-\theta) \rightarrow \text{natural parameter}$$

$$d(\theta) = r \ln \theta \quad d(u) = \ln \binom{r+u-1}{r-1}$$

Expected Value & Variance of $a(u)$

$$E(a(u)) = \frac{c'(\theta)}{b'(\theta)} \quad \text{&} \quad V(a(u)) = \frac{b''(\theta)c'(\theta) - c''(\theta)b'(\theta)}{[b'(\theta)]^3}$$

$$V(a(u)) = \frac{b''(\theta)c'(\theta) - c''(\theta)b'(\theta)}{[b'(\theta)]^3}$$

Example:

Binomial.

$$a(u) = u$$

$$b(\theta) = \ln\left(\frac{P}{1-P}\right) \quad d(\theta) = n \ln(1-P)$$

$$b'(\theta) = \frac{1}{P(1-P)} \quad c'(\theta) = \frac{-n}{1-P}$$

$$b''(\theta) = \frac{2P-1}{[P(1-P)]^2} \quad c''(\theta) = \frac{-n}{(1-P)^2}$$

$$E(a(u)) = E(u) = \frac{c'(\theta)}{b'(\theta)} = \frac{n}{1-P} * P(1-P)$$

$$= np$$

$$V(a(u)) = V(u) = np(1-P)$$

Suppose we have a set of independent

Observations variable.

$$(y_i, \tilde{x}_i^t), i=1(1)n, \tilde{x}_i^t = (x_{i1}, x_{i2}, \dots, x_{ip})$$

from some exponential type of distribution of canonical form

[ie $a(y) = y$]. Then the joint probability density function is

$$\begin{aligned} f(y_1, y_2, \dots, y_n, \theta, \phi) &= \exp \left\{ \sum_{i=1}^n y_i b(\theta) + \sum_{i=1}^n c(\theta_i) + \right. \\ &= \exp \left\{ \sum_{i=1}^n y_i b(\theta) + \sum_{i=1}^n d(y_i) \right\} \\ &= \prod_{i=1}^n \exp \left\{ y_i b(\theta_i) + c(\theta_i) + d(y_i) \right\} \end{aligned}$$

where ϕ is a vector of nuisance parameters that occur with $b(\cdot)$, $c(\cdot)$ and $d(\cdot)$.

$\theta = (\theta_1, \theta_2, \dots, \theta_n)'$ vector of parameters of interest

The variance in response variable
can be explained in terms of
 \hat{y}_i values

$$\hat{x}_i = (\hat{x}_{i1}, \hat{x}_{i2}, \dots, \hat{x}_{ip})'$$

Consider the set of parameters

$$\beta = (\beta_1, \beta_2, \dots, \beta_p)'$$

We find some suitable link function

$g(\cdot)$ such that

$$\left. \begin{array}{l} y_i \text{ follows binomial} \\ g(\mu_i) = \hat{x}_i' \beta \\ \ln\left(\frac{P}{1-P}\right) = \hat{x}_i' \beta \quad y_i \text{ is Normal} \\ \mu_i \text{ is } E(y_i) = \hat{x}_i' \beta \end{array} \right\} \begin{array}{l} y_i = \hat{x}_i' \beta + \epsilon \\ E(y_i) = \hat{x}_i' \beta \end{array}$$

A link function that is often regarded as sensible one is natural parameter.

GLM analysis comes into play when the error distribution is not normal but the error distribution, the response variable distribution must be a member of Exponential family.

Given (Y_i, \tilde{x}_i')

We would hope that variance in Y_i or $E(Y_i)$ values could be explained by terms of \tilde{x}_i' value.

We would hope that we could find suitable link function

$g(\theta_i)$ such that the model

$$g(\theta_i) = \tilde{x}_i' \beta \text{ held.}$$

Where

$\beta = (\beta_1, \beta_2, \dots, \beta_p)'$ is vector of regression coefficient

link function is often the natural parameter.

If y_i from normal distribution

$$y_i = \tilde{x}_i' \beta + \epsilon$$

$$\theta_i = E(y_i) = \tilde{x}_i' \beta$$

$$\theta_i = \tilde{x}_i' \beta$$

Natural parameter for normal distribution is

$$\theta$$

$$g(\theta) = \theta$$

We go for model

$$g(\theta_i) = \tilde{x}_i' \beta$$

$$\ln\left(\frac{\theta_i}{1-\theta_i}\right) = \tilde{x}_i' \beta$$

$$\theta_i = \frac{\exp(\tilde{x}_i' \beta)}{1 + \exp(\tilde{x}_i' \beta)}$$

$$E(y_i) = \frac{\exp(\tilde{x}_i' \beta)}{1 + \exp(\tilde{x}_i' \beta)}$$

Model - for
Binomial
distribution.

In detail about Binomial Distribution

Suppose we have data (y_i, \tilde{x}_i') from a binomial distribution $\text{Bin}(n_i, p_i)$

The single observation y_i is of the form $\frac{r_i}{n_i}$, where r_i is the no. of successes in n_i trials, each having probability p_i of success. and \tilde{x}_i'

$\tilde{x}_i' = (x_{i1}, x_{i2}, \dots, x_{ip})$ is a

set of observations of P regressions associated with y_i .

Binomial distribution is a member of the exponential family

$$\text{Joint Pdt} = f(y_1, y_2, \dots, y_n)$$

$$= \prod_{i=1}^n \binom{n_i}{y_i} (p_i)^{y_i} (1-p_i)^{n_i-y_i}$$

$$= \prod_{i=1}^n \exp \left\{ y_i \ln \left(\frac{p_i}{1-p_i} \right) + n_i \ln(1-p_i) + \ln \binom{n_i}{y_i} \right\}$$

$$= \exp \left\{ \sum_{i=1}^n y_i \ln \left(\frac{p_i}{1-p_i} \right) + \sum_{i=1}^n n_i \ln(1-p_i) + \sum_{i=1}^n \ln \binom{n_i}{y_i} \right\}$$

Given $y_i \sim x_i$ try to explain the variability of x_i in y_i

We would hope that the variation in the y_i | $E(y_i) = p_i$ $y_i = \frac{r_i}{n_i}$

could be explained in terms of the x_i values,

ie, we would hope that we could find a suitable link function $g(\cdot)$

such that

$$g(P_i) = \tilde{x}_i' \beta$$

Binomial distribution

Natural parameter = $\ln\left(\frac{P_i}{1-P_i}\right)$

we fit the model

$$\ln \frac{P_i}{1-P_i} = \tilde{x}_i' \beta = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

$$E(Y_i) = P_i = \frac{\exp(\tilde{x}_i' \beta)}{1 + \exp(\tilde{x}_i' \beta)}$$

This is the model instead

of fitting $y_i = \tilde{x}_i' \beta + \epsilon$

finally,

$$P_i = \frac{\exp(\tilde{x}_i' \beta)}{1 + \exp(\tilde{x}_i' \beta)} \quad \text{--- } \star$$

When $\tilde{x}_i' \beta = \beta_1 + \beta_2 x_{i2}$, is this is
 a situation \star
 called Logistic function.

Estimation via Maximum Likelihood

To estimate β , we use the
 method of maximum likelihoods.
 Compute likelihood function L .

$$L = \exp \left\{ \sum_{i=1}^n y_i \ln \left(\frac{p_i}{1-p_i} \right) + \sum_{i=1}^n n_i \ln (1-p_i) + \sum_{i=1}^n \ln \left(\frac{n_i}{y_i} \right) \right\}$$

$$\ln L = \sum_{i=1}^n y_i \ln \frac{p_i}{1-p_i} + \sum_{i=1}^n n_i \ln (1-p_i) + \sum_{i=1}^n \ln \left(\frac{n_i}{y_i} \right)$$

fit the model $p_i = \frac{\exp(\tilde{x}_i' \beta)}{1 + \exp(\tilde{x}_i' \beta)}$

Write loglikelihood in terms of β

$$= \sum_{i=1}^n y_i \tilde{x}_i' \beta + \sum_{i=1}^n n_i \ln(1 + \exp(\tilde{x}_i' \beta)) \\ + \sum_{i=1}^n \ln\left(\frac{n_i}{y_i}\right)$$

Maximize log likelihood $\ln L$ with respect to β

Differentiate $\ln L$ in terms of β

It is not easy

Numerical Search Method / Iteratively reweighted least square (IRLS) could be used to compute MLE's of β .

Example Pneumoconiosis Data

| Number of year of Exposure | No. of Severe cases | Total no of miners | Proportion of severe cases, y |
|-------------------------------|------------------------|-----------------------|------------------------------------|
| 5-8 | 0 | 98 | 0 |
| 15-0 | 1 | 54 | 0.0185 |
| 21-5 | 3 | 43 | 0.0698 |
| 27-5 | 8 | 48 | 0.1667 |
| 33-5 | 9 | 51 | 0.1765 |
| 39-5 | 8 | 38 | 0.2105 |
| 46-0 | 10 | 28 | 0.3571 |
| 51-5 | 5 | 11 | 0.4545 |

y : proportion of miners who have severe symptoms.

x_i : No. of years of exposure

Probability distribution for the number of severe cases is binomial

We will fit logistic regression model to the data:

$$P_i = E(Y_i) = \frac{\exp(\tilde{x}_i' \beta)}{1 + \exp(\tilde{x}_i' \beta)}$$

$$\tilde{x}_i' \beta = \underbrace{\beta_1 + \beta_2 x_i}$$

finally fitted Model

$$y_i = \frac{\exp(4.79 - 0.0935x_i)}{1 + \exp(4.79 - 0.0935x_i)}$$

Poisson distribution

Data (y_i, \tilde{x}_i') from Poisson $P(\mu_i)$,

$$E(y_i) = \mu_i$$

$$f(y, \mu) = \exp \{ y \ln \mu - \mu - \ln y! \}$$

$\ln \mu$ is the natural parameter

The variation in y_i could be explained in terms of the \tilde{x}_i' values.
we fit the model

$$g(\mu_i) = \tilde{x}_i' \beta$$

$$\ln \mu_i = \tilde{x}_i' \beta = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

$$y_i = e^{\tilde{x}_i' \beta} + \epsilon$$

Choice of Link function

| Distribution | Link function | Name |
|--|----------------------------|-----------------|
| Normal (continuous data) | $g(\mu) = \mu$ | identity link |
| Binomial (Binary or Proportion data) | $g(p) = \ln \frac{p}{1-p}$ | logistic link |
| Poisson (for count data) | $g(\mu) = \ln \mu$ | log link |
| Exponential | $g(\mu) = \frac{1}{\mu}$ | reciprocal link |
| Gamma (positive Continuous data) | $g(\mu) = \frac{1}{\mu}$ | reciprocal link |

Link Function

→ It is the transformation
that connects predicted values of the
dependent variables to the observed
values