

GLMs (Generalized Linear Models) Introduction

1. Assumption

Distribution of the dependent variable (for example $Y \sim \text{Bin}(n, p)$, $Y \sim \text{Gamma}(\alpha, \beta)$)

2. Specify a link function $g(\cdot)$

("Linearize Y ") such that $g(E(Y)) = X\beta$

If $E[y_i] = \mu_i$, $g(\mu_i) = \eta_i$, $\eta_i = X_i\beta$

Fitting a GLM.

Suppose $Y_i \sim \text{Bin}(n_i, p_i)$.

and we wish to predict Y_i/n_i .

Then $E(Y_i/n_i) = p_i$

A common link function is logit: $\log\left(\frac{p}{1-p}\right)$

$$g(\mu_i) = \log\left(\frac{\mu_i}{n_i - \mu_i}\right), \mu_i = n_i p_i$$

$$P(y_i = y_i | p_i, n_i) = \binom{n_i}{y_i} p_i^{y_i} (1-p)^{n_i - y_i}$$

Suppose now that you know $y_i \stackrel{iid}{\sim} \text{Bin}(n, p)$

What is the likelihood function?

$$L = \prod_{i=1}^K \binom{n}{y_i} p^{y_i} (1-p)^{n - y_i}$$

$$\log L = \sum_{i=1}^K y_i \log p + (n - y_i) \log(1-p) + \log \binom{n}{y_i}$$

What value of p maximizes the likelihood?

$$p = \arg \max_{p \in [0,1]} \log L$$

Great! How do we connect this with $X\beta$ - ?

Remember the logit function

$$X_i \beta = \log \left(\frac{\mu_i}{1 - \mu_i} \right) = \log \left(\frac{np}{1 - np} \right)$$

$$\Rightarrow p = \text{inv logit}(X_i \beta) = \frac{e^{X_i \beta}}{1 + e^{X_i \beta}}$$

$$\Rightarrow p = \frac{e^{X_i \beta}}{1 + e^{X_i \beta}}$$

↳ substitute this into the log-likelihood

$$\log L = \sum_{i=1}^n \left[\log \binom{n}{y_i} + y_i \cdot \log \left[\frac{1}{n} \cdot \frac{e^{X_i \beta}}{1 + e^{X_i \beta}} \right] + (1 - y_i) \cdot \log \left[\frac{1}{n} \cdot \frac{1}{1 + e^{X_i \beta}} \right] \right]$$

Now you can solve

$$\hat{\beta} = \underset{\beta}{\text{argmax}} \log L$$

The same logic is applied to the deviance

$$D = -2(l(\beta_{PM}) - l(\beta_{SM}))$$

D-deviance

$l(\beta_{PM})$: log-likelihood of the
Proposed model (PM)

$l(\beta_{SM})$: log-likelihood of the
Saturated model (SM)

(one parameter per observation - fits the
data perfectly) [n parameters to estimate]

$$l(\beta_{SM}) = \underline{1}$$

Residual deviance : $-2l(\beta_{PM})$

Null deviance : $-2l(\beta_{NM})$

NM - no model (empty model)

Exponential family of distributions
(normal, binomial, poisson, gamma, exponential ...)

Density:

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

ϕ - dispersion parameter

θ - canonical parameter

It can be shown that

$$E[Y] = b'(\theta) = \mu$$

(' denotes a derivative)

$$\text{and } \text{Var}[Y] = \phi b''(\theta) = \phi V(\mu)$$

Link function

$$g(\mu_i) = X_i \beta$$

$$g = (b')^{-1} \Rightarrow g(\mu_i) = \theta_i$$

\uparrow
canonical link fct.

The log-likelihood:

$$l = \sum_{i=1}^n \left[\frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i, \phi_i) \right]$$

The max. likelihood estimates are obtained by solving score equations

$$s(\beta_j) = \frac{\partial l}{\partial \beta_j} = 0, \text{ for } \beta_j$$