

MAT-468: Sesión 7, Estimación en modelos lineales generalizados

Felipe Osorio

<http://fosorios.mat.utfsm.cl>

Departamento de Matemática, UTFSM



Definición 1 (familia exponencial)

Sea Y una variable aleatoria con densidad

$$f(y; \theta, \phi) = \exp[\phi\{y\theta - b(\theta)\} + c(y, \phi)],$$

donde

$$E(Y) = \mu = b'(\theta), \quad \text{var}(Y) = \phi^{-1}V(\mu),$$

con $V(\mu) = d\mu/d\theta = b''(\theta)$ la **función de varianza** y $\phi^{-1} > 0$ el parámetro de dispersión. En este caso anotamos $Y \sim \text{FE}(\theta, \phi)$.



Definición 2 (modelo lineal generalizado)

Considere Y_1, \dots, Y_n variables aleatorias independientes. Un **modelo lineal generalizado (GLM)** se define como:

$$Y_i \sim \text{FE}(\theta_i, \phi), \quad i = 1, \dots, n,$$

donde la media μ_i está relacionada con el **predictor lineal** $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$, mediante la **función de enlace** g , como

$$g(\mu_i) = \eta_i,$$

para g una función monótona y diferenciable.



Distribución normal

Considere $Y \sim N(\mu, \sigma^2)$. De este modo, tenemos la función de densidad

$$\begin{aligned} f(y; \mu, \sigma^2) &= (2\pi\sigma^2)^{-1/2} \exp \left\{ -\frac{1}{2\sigma^2}(y - \mu)^2 \right\} \\ &= \exp \left[\left\{ \frac{1}{\sigma^2} \left(y\mu - \frac{\mu^2}{2} \right) - \frac{1}{2} \left\{ \log(2\pi\sigma^2) + \frac{y^2}{\sigma^2} \right\} \right\} \right], \end{aligned}$$

es decir, $\theta = \mu$, $b(\theta) = \theta^2/2$, $\phi = \sigma^{-2}$ y $c(y, \phi) = \frac{1}{2} \log(\phi/2\pi) - \frac{\phi y^2}{2}$.

Distribución Poisson

En el caso que $Y \sim \text{Poi}(\mu)$, tenemos la función de probabilidades,

$$p(y; \mu) = \frac{e^{-\mu} \mu^y}{y!} = \exp(y \log \mu - \mu - \log y!),$$

así, basta considerar $\log \mu = \theta$, $b(\theta) = e^\theta$, $\phi = 1$ y $c(y, \phi) = -\log y!$ para notar que $\text{Poi}(\mu)$ pertenece a la familia exponencial.



Función de log-verosimilitud

Para Y_1, \dots, Y_n siguiendo un GLM tenemos que

$$\begin{aligned}\ell(\boldsymbol{\psi}) &= \sum_{i=1}^n \log f(y_i; \theta_i, \phi) \\ &= \sum_{i=1}^n \{\eta_i \phi(y_i \theta_i - b(\theta_i)) + c(y_i; \phi)\} \\ &= \phi \sum_{i=1}^n (y_i \theta_i - b(\theta_i)) + \sum_{i=1}^n c(y_i; \phi),\end{aligned}$$

donde $\boldsymbol{\psi} = (\boldsymbol{\beta}^\top, \phi)^\top$.

Observación:

Debemos destacar que $\boldsymbol{\beta}$ y ϕ son parámetros ortogonales, y por tanto la inferencia puede ser realizada de manera independiente.



Función score

Para GLM la función score $U(\beta) = \dot{\ell}(\beta)$ adopta la forma:

$$\begin{aligned}\frac{\partial \ell(\psi)}{\partial \beta} &= \phi \sum_{i=1}^n \left\{ Y_i \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \beta} - \frac{db(\theta_i)}{d\theta_i} \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \beta} \right\} \\ &= \phi \sum_{i=1}^n \left\{ Y_i - \frac{db(\theta_i)}{d\theta_i} \right\} \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \beta} \\ &= \phi \sum_{i=1}^n (Y_i - b'(\theta_i)) \left\{ \frac{d\mu_i}{d\theta_i} \right\}^{-1} \frac{d\mu_i}{d\eta_i} \frac{\partial \eta_i}{\partial \beta} \\ &= \phi \sum_{i=1}^n (Y_i - \mu_i) V(\mu_i)^{-1} \frac{d\mu_i}{d\eta_i} \mathbf{x}_i \\ &= \phi \sum_{i=1}^n \omega_i^{1/2} \frac{(Y_i - \mu_i)}{\sqrt{V_i}} \mathbf{x}_i,\end{aligned}$$

donde $\omega_i = (d\mu_i/d\eta_i)^2/V_i$ y $V_i = V(\mu_i)$, para $i = 1, \dots, n$.



Función score

Una manera mucho más compacta de escribir la función score en GLM es:

$$U(\boldsymbol{\beta}) = \phi \mathbf{X}^\top \mathbf{W}^{1/2} \mathbf{V}^{-1/2} (\mathbf{Y} - \boldsymbol{\mu}),$$

donde $\mathbf{W} = \text{diag}(\omega_1, \dots, \omega_n)$, $\mathbf{V} = \text{diag}(V_1, \dots, V_n)$ y $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$.

Vamos a suponer que \mathbf{X} es matriz de rango completo cuya i -ésima fila es dada por \mathbf{x}_i^\top , para $i = 1, \dots, n$. Además, $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)^\top$ con $\mu_i = \mu_i(\boldsymbol{\beta})$.



Matriz de información de Fisher

La matriz Hessiana en GLM es dada por:

$$\begin{aligned}\frac{\partial^2 \ell(\boldsymbol{\psi})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} &= \phi \sum_{i=1}^n (Y_i - \mu_i) \frac{d^2 \theta_i}{d \mu_i^2} \left(\frac{d \mu_i}{d \eta_i} \right)^2 \mathbf{x}_i \mathbf{x}_i^\top \\ &\quad + \phi \sum_{i=1}^n (Y_i - \mu_i) \frac{d \theta_i}{d \mu_i} \left(\frac{d^2 \mu_i}{d \eta_i^2} \right)^2 \mathbf{x}_i \mathbf{x}_i^\top \\ &\quad - \phi \sum_{i=1}^n \frac{d \theta_i}{d \mu_i} \left(\frac{d \mu_i}{d \eta_i} \right)^2 \mathbf{x}_i \mathbf{x}_i^\top\end{aligned}$$

De este modo, la matriz de información de Fisher asume la forma,

$$\mathcal{F}(\boldsymbol{\beta}) = \mathbb{E} \left\{ - \frac{\partial^2 \ell(\boldsymbol{\psi})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} \right\} = \phi \sum_{i=1}^n \frac{(d \mu_i / d \eta_i)^2}{V_i} \mathbf{x}_i \mathbf{x}_i^\top = \phi \mathbf{X}^\top \mathbf{W} \mathbf{X}.$$



Algoritmo Fisher-scoring en GLM

El algoritmo Fisher-scoring para β es dado por:

$$\beta^{(r+1)} = \beta^{(r)} + \mathcal{F}^{-1}(\beta^{(r)})U(\beta^{(r)}),$$

y para el caso de GLM, adopta la forma:

$$\begin{aligned}\beta^{(r+1)} &= \beta^{(r)} + (\mathbf{X}^\top \mathbf{W}^{(r)} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}^{(r)1/2} \mathbf{V}^{(r)-1/2} (\mathbf{Y} - \boldsymbol{\mu}^{(r)}) \\ &= (\mathbf{X}^\top \mathbf{W}^{(r)} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}^{(r)} \mathbf{X} \beta^{(r)} \\ &\quad + (\mathbf{X}^\top \mathbf{W}^{(r)} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}^{(r)1/2} \mathbf{V}^{(r)-1/2} (\mathbf{Y} - \boldsymbol{\mu}^{(r)}) \\ &= (\mathbf{X}^\top \mathbf{W}^{(r)} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W}^{(r)} \mathbf{Z}^{(r)},\end{aligned}$$

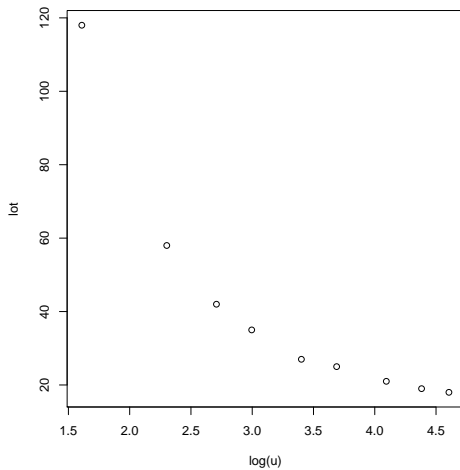
donde

$$\mathbf{Z} = \boldsymbol{\eta} + \mathbf{W}^{-1/2} \mathbf{V}^{-1/2} (\mathbf{Y} - \boldsymbol{\mu}),$$

denota la respuesta de trabajo.



Datos de coagulación



Conjunto de datos: Coagulación inducida por dos lotes de tromboplastina.

```
clotting <- data.frame(  
+   u = c(5,10,15,20,30,40,60,80,100),  
+   lot = c(118,58,42,35,27,25,21,19,18))
```

Exploramos el conjunto de datos por medio del gráfico:

```
> plot(lot ~ log(u), data = clotting)
```



Datos de coagulación de la sangre

```
> fit <- glm(lot ~ log(u), data = clotting, family = Gamma)
> summary(fit)

Call:
glm(formula = lot ~ log(u), family = Gamma, data = clotting)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-0.04008 -0.03756 -0.02637  0.02905  0.08641

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.0165544  0.0009275  -17.85 4.28e-07 ***
log(u)       0.0153431  0.0004150   36.98 2.75e-09 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for Gamma family taken to be 0.002446059)

Null deviance: 3.51283  on 8  degrees of freedom
Residual deviance: 0.01673  on 7  degrees of freedom
AIC: 37.99

Number of Fisher Scoring iterations: 3
```

