

Modelos Discretos

Modelos Lineales Generalizados (GLM)

Mtr. Alcides Ramos Calcina



Modelos Lineales Generalizados (GLM)

Mtr. Alcides Ramos Calcina



Introducción

- El Modelo Lineal Generalizado (**MLG**), que tratamos en este curso, es la extensión natural del Modelo Lineal clásico.
- Expuesto por Nelder y Wedderburn (1972), convirtiéndose en una solución especialmente adecuada para **modelos de dependencia con datos no métricos**.
- En el área de ciencias sociales es frecuente trabajar **atributos**, **actitudes** o **conductas** que, se miden de forma no métrica (discreta, nominal u ordinal), por tanto, no se ajusta, al Modelo Lineal clásico, incumpliendo los supuestos de linealidad y normalidad.
- Por ejemplo, la clasificación binaria de apto - no apto, es una situación que requieren de modelos que trabajen con datos dicotómicos, ordinales, categóricos o de elecciones discretas, es decir, de modelos de probabilidad de un evento.

1. Modelo Estadístico

- Un modelo pretende explicar la variación de una respuesta a partir de la relación conjunta de dos fuentes de variabilidad, una de carácter determinista y otra aleatoria, lo que responde a la expresión:

$$\text{Respuesta} = \text{componente sistemático} + \text{componente aleatorio}$$

- Otros autores toman la siguiente expresión:

$$\text{DATOS} = \text{MODELO} + \text{ERROR}$$

Donde:

MODELO: está asociado a la parte sistemática. es la función que se introduce con objeto de explicar los datos.

DATOS: corresponde a las observaciones que se quieren analizar (la variable de respuesta o variable dependiente).

ERROR: es la que contiene la discrepancia o falta de ajuste entre *Datos* y *Modelo*. Mtr. Alcides Ramos Calcina

1. Modelo Estadístico

1. Criterios para la construcción, formulación y ajuste de modelos

Criterio estadístico o principio de bondad de ajuste:

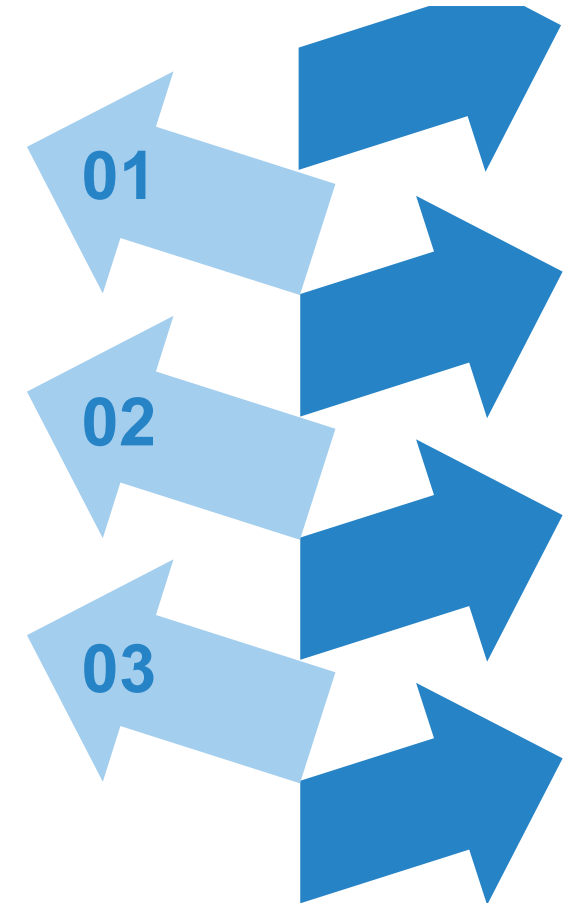
La inclusión de parámetros en el MODELO en beneficio de una mejor representación de los DATOS con la correspondiente disminución del ERROR

Criterio lógico o principio de parsimonia:

la selección de los parámetros que formen parte del modelo de tal modo que éste se convierta en una representación simple y sobria de la realidad.

Criterio sustantivo o integración teórica:

del modelo en la red conceptual que lo generó.



1. Modelo Estadístico

2. Etapas de la construcción de modelos



2. Modelo Lineal (LM)

- La fórmula general del Modelo Lineal es:

$$Y = f(X) + g(\varepsilon)$$

donde toda observación sobre la variable de respuesta es la suma de:

- a) los **efectos de un grupo de factores** o **componentes sistemáticos**, $f(X)$, que implican un conjunto de parámetros de una población y un conjunto de variables independientes relevantes medidas sobre cada uno de los sujetos con los que se trabaja.
- a) **función** $g(\varepsilon)$, que representa el **efecto de los componentes aleatorios** y es resultado de una o más distribuciones de probabilidad dependientes de un pequeño número de parámetros.

2. Modelo Lineal (LM)

- Se habla, entonces, de un **Modelo Lineal** de primer orden para k variables explicativas y $k + 1$ parámetros si el modelo es lineal en sus parámetros y en sus variables explicativas, respondiendo a la siguiente fórmula general:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k + \varepsilon$$

$$Y = \beta_0 + \sum_{j=1}^k \beta_j X_j + \varepsilon$$

- Si el modelo es lineal en sus parámetros, pero no en las variables explicativas sería un Modelo Lineal de m -ésimo orden (cuadrático, cúbico, etc.) con km variables independientes y $km + 1$ parámetros. Su formulación es:

$$Y = \beta_0 + \sum_{j=1}^k \beta_j X_j + \sum_{j=1}^k \beta_{j1} X_{j1}^2 + \cdots + \sum_{j=1}^k \beta_{jm} X_{jm}^m + \varepsilon$$

2. Modelo Lineal (LM)

2.1. Limitaciones en el Modelo Lineal

- Sea cual sea el tipo de datos que pretendemos analizar, para una modelización adecuada siempre es básico preguntar y contestar una serie de cuestiones:

¿Cuál es la variable de respuesta, esto es, la variable a explicar?

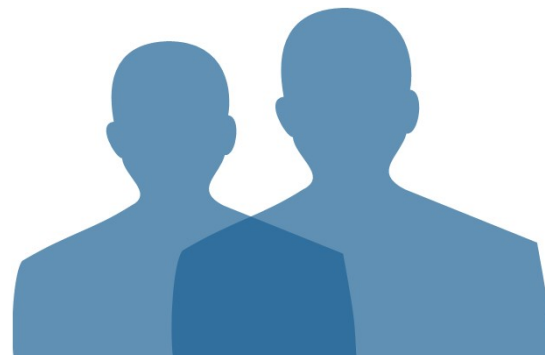
¿Son siempre homogéneas las condiciones de observación?

¿Cómo cambia la respuesta media cuando se modifican las condiciones de experimentación?

¿Qué variables pueden explicar algo sobre la respuesta?

Así mismo, se debe corroborar que las respuestas verifiquen:

- Normalidad
- Homogeneidad de varianzas
- Linealidad (aditividad) de los efectos sistemáticos



2. Modelo Lineal (LM)

- Se tiene el Modelo Lineal normal en su forma matricial: $Y = X\beta + \varepsilon$

Donde $\varepsilon = N(0, \sigma^2 I)$, de forma que,

$$E(Y) - \mu = X\beta \qquad \text{Var}(Y_i) = \sigma^2$$

- Sin embargo, el mundo de los datos no es siempre **NORMAL**.
- Cuando no se verifican las hipótesis del modelo lineal, luego de intentar algunas transformaciones de las variables, surgen algunos inconvenientes:
 - ✓ Algunas transformaciones no están definidas en las fronteras del espacio muestral (como la logit).
 - ✓ No es siempre directa la interpretación de la variable transformada.
 - ✓ Generalmente no existe una técnica única que corrija a la vez todos los defectos.

2. Modelo Lineal (LM)

2.2. Generalización del Modelo Lineal

- Los Modelos Lineales Generalizados (GLM) son una alternativa a transformaciones de la respuesta, justificadas por:
 - ❖ La falta de **linealidad**
 - ❖ La falta de **homogeneidad de la varianza**.
- Las hipótesis básicas de un modelo lineal generalizado (GLM) son:
 - ❑ **Independencia entre las respuestas**
 - ❑ La “respuesta media” cambia con las condiciones, pero no la “forma funcional” de la distribución.
 - ❑ La respuesta media, o alguna transformación de ella, cambia de modo lineal cuando las condiciones cambian.

3. Modelo Lineal Generalizado (GLM)

3.1. Definición de un GLM

- Los modelos lineales generalizados (GLM) amplían el concepto del modelo de regresión lineal. El modelo lineal supone que,

$$E(Y/X) = X' \beta$$

- Su equivalente, $Y = X' \beta + \varepsilon$
- Desafortunadamente, la **restricción a la linealidad** no puede tener en cuenta una variedad de situaciones prácticas.

3. Modelo Lineal Generalizado (GLM)

- En los GLM se demuestra que, si la distribución de la variable dependiente Y es un miembro de la familia exponencial, entonces la clase de modelos que conecta la esperanza de Y a una combinación lineal de las variables $X'\beta$ puede tratarse de manera unificada.
- Por consiguiente, denotamos la función que relaciona:

$$\mu = E(Y | X) \quad \text{y} \quad \eta = X' \beta$$

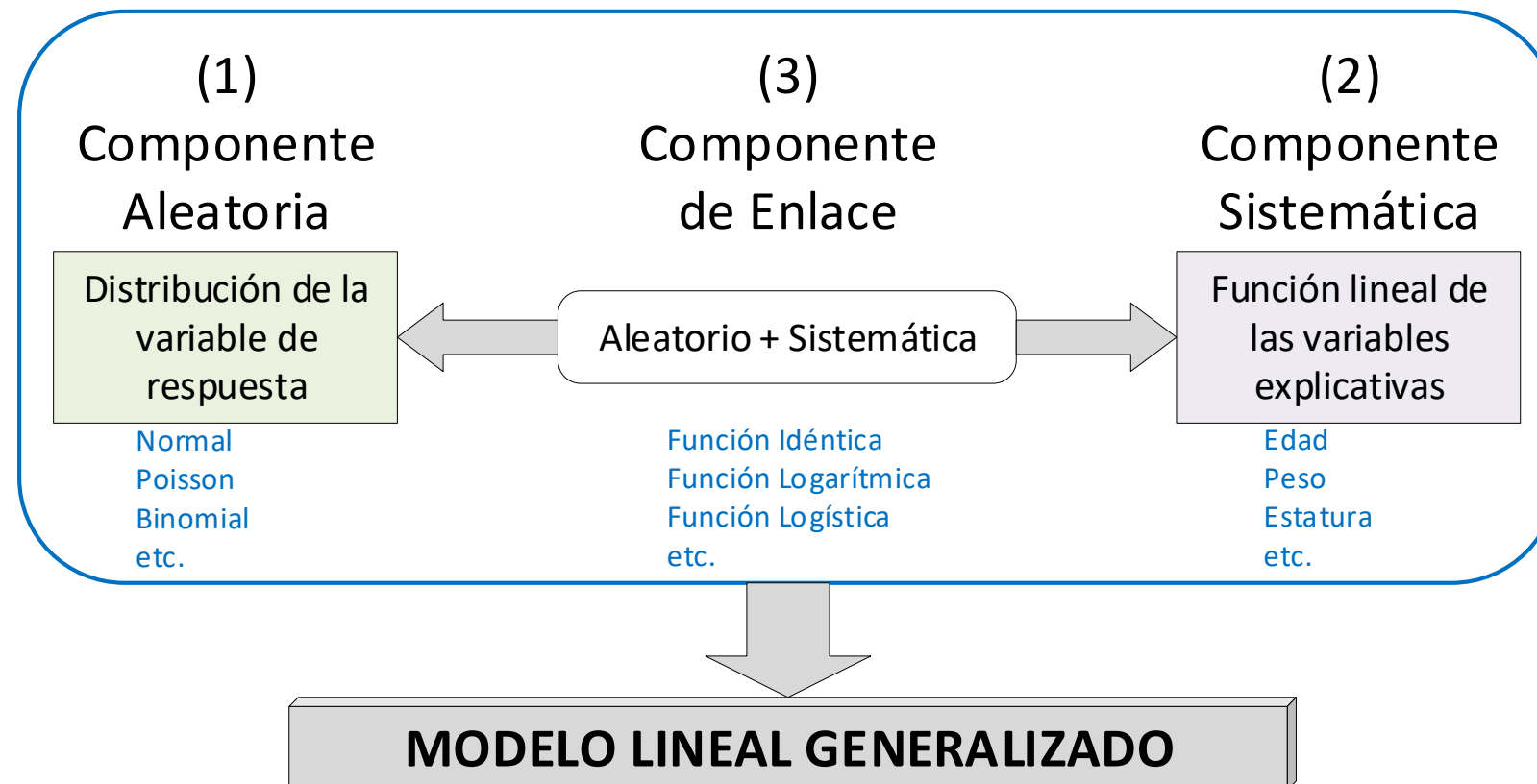
por

$$\eta = G(\mu) \quad \text{o} \quad E(Y | X) = G^{-1}(X' \beta)$$

3. Modelo Lineal Generalizado (GLM)

3.2. Componentes de un GLM

- Un modelo lineal generalizado tiene tres componentes básicos:



3. Modelo Lineal Generalizado (GLM)

a) Componente aleatoria

- En muchas aplicaciones, las observaciones de Y son binarias y se identifican como éxito y fracaso, y se modeliza como una **distribución binomial**.
- En otras ocasiones cada observación es un recuento, con lo que se puede asignar a Y una distribución de **Poisson** o una distribución **binomial negativa**.
- Esta componente, identifica la **distribución de probabilidad** de la variable **dependiente**. Esto es, cada variable muestral Y_i tiene función de densidad de la forma:

$$f(y_i, \theta_i, \phi) = \exp\left[\frac{y_i \theta_i - b(\theta_i)}{\phi} + c(y_i, \phi)\right] = G(y_i, \phi) \exp\left[\frac{y_i \theta_i - b(\theta_i)}{\phi}\right]$$

3. Modelo Lineal Generalizado (GLM)

b) Componente sistemática

- Especifica una función lineal η de los valores fijados x_{1i}, \dots, x_{ki} de las variables explicativas X_1, X_2, \dots, X_k dada por:

$$\eta_i := \delta + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} \quad i = 1, 2, 3, \dots, n$$

donde los β_k son los llamados parámetros del modelo lineal generalizado, incluyendo el llamado intercepto como $\delta = \beta_0$, siendo

3. Modelo Lineal Generalizado (GLM)

- Si se reúnen los valores observados de las variables explicativas en la llamada matriz de diseño:

$$C = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix}$$

de tamaño $n \times (1+k)$, los parámetros del modelo en el vector

$$\alpha = (\delta, \beta_1, \beta_2, \dots, \beta_k)'$$

entonces, la ecuación lineal puede ser escrita en forma vectorial como $\eta = C.\alpha$

3. Modelo Lineal Generalizado (GLM)

c) Función link o de enlace

- Sea μ_i la esperanza condicional de Y_i dada la condición x_{i1}, \dots, x_{ik} , es decir,
 $\mu_i := E(Y | x_{i1}, x_{i2}, \dots, x_{ik})$, $i = 1, \dots, n$. Entonces, este enlace está dado por una llamada función de enlace:

$$g(\mu_i) = \theta_i$$

en cuyo caso resultan $\theta_i = \eta_i$, y el enlace está descrito por la expresión

$$\theta_i = \delta + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}$$

4. Tipos de GLMs

4.1. Lineal

- Supongamos que la variable Y_i , $i = 1, \dots, n$ está normalmente distribuida con esperanza μ_i y varianza σ^2 . La función de densidad en los valores y_i viene dada por:

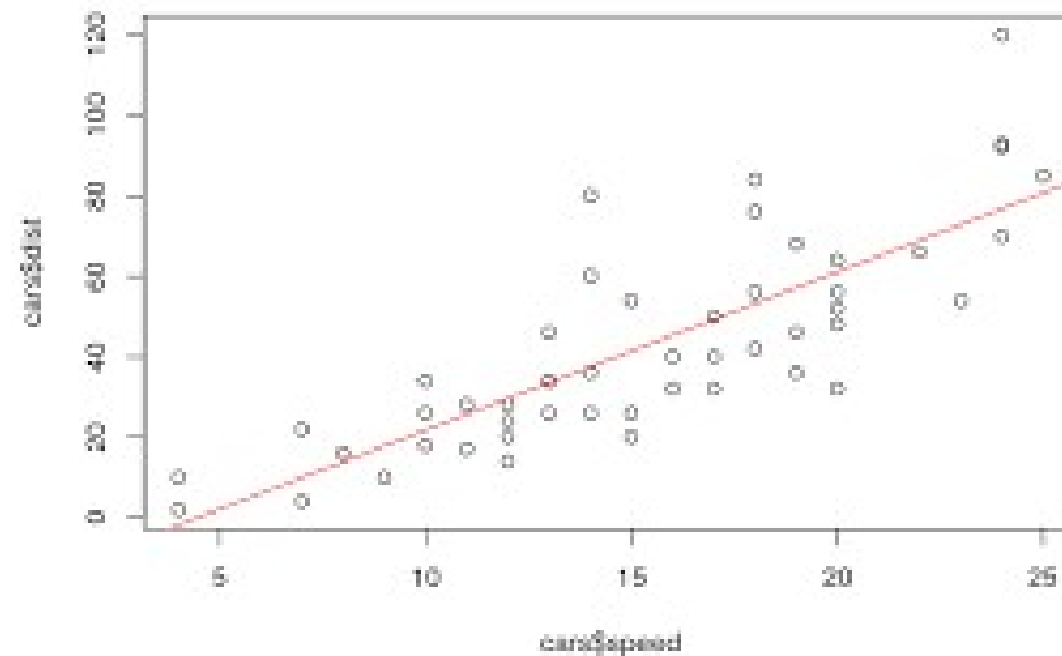
$$\begin{aligned} f(y_i, \theta_i, \phi) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(y_i - \mu_i)^2\right] \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{y_i^2}{2\sigma^2}\right] \cdot \exp\left[\frac{y_i\mu_i - \frac{\mu_i^2}{2}}{\sigma^2}\right] \end{aligned}$$

- Aquí, $\theta_i = \mu_i$ son los parámetros naturales y se tiene $\phi = \sigma^2$ como parámetro de dispersión. Además,

$$b(\theta_i) = \frac{\mu_i^2}{2} \qquad G(y_i, \phi) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{y_i^2}{2\sigma^2}\right]$$

4. Tipos de GLMs

- El enlace canónico está dado por la función identidad $g(\mu_i) = \mu_i$. Los GLMs que usan el enlace identidad son llamados modelos lineales.
- En la siguiente figura se muestra una representación gráfica de un modelo de regresión lineal.



4. Tipos de GLMs

4.2. Log-linear

- Supongamos que la Y_i , $i = 1, \dots, n$, es una variable de Poisson con parámetro $\lambda_i > 0$. Su función de densidad viene dada por

$$f(y_i, \lambda_i) = \frac{\lambda_i^{y_i} e^{-\lambda_i}}{y_i!}$$

- para todo $y_i \in \{0, 1, 2, \dots\}$. Escribiendo esta función en la forma:

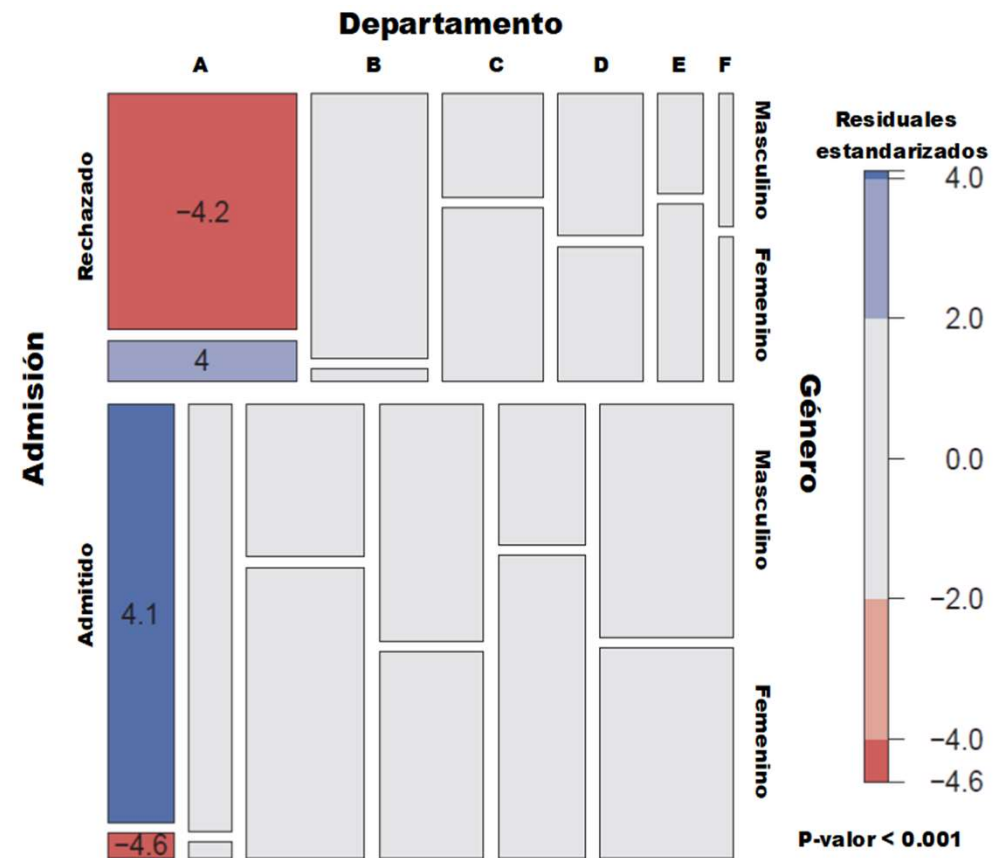
$$f(y_i, \lambda_i) = \frac{\exp[y_i \ln \lambda_i - \lambda_i]}{y_i!} = \exp[y_i \ln \lambda_i - \lambda_i - \ln(y_i!)]$$

- Aquí, $\theta_i = \ln \lambda_i$ son los parámetros naturales y el parámetro de dispersión es $\phi = 1$. Además,

$$b(\theta_i) = e^{\theta_i} = \lambda_i \qquad G(y_i, \phi) = \frac{1}{y_i!}$$

4. Tipos de GLMs

- El enlace canónico es $g(\lambda_i) = \ln \lambda_i$. Los GLMs que usan el enlace log son llamados modelos loglineales.
- En la figura se muestra una representación gráfica (en mosaico) de un modelo de regresión loglineal, que muestra los residuos estandarizados de las contribuciones de las celdas en el valor del estadístico de prueba correspondiente.



4. Tipos de GLMs

4.3. Logístico (Y es Bernoulli)

- Muchas variables categóricas tienen únicamente dos categorías. La observación para cada caso puede ser clasificada como éxito o fracaso.
- En este caso, la variable Y_i tiene distribución de **Bernoulli** con parámetro p_i . Su función de densidad es

$$f(y_i, p_i) = \exp \left[y_i \ln \left(\frac{p_i}{1 - p_i} \right) - \ln \left(\frac{1}{1 - p_i} \right) \right]$$

para todo $0 < p_i < 1$ y $y_i \in \{0, 1\}$. Los parámetros naturales son $\theta_i = \ln \left(\frac{p_i}{1 - p_i} \right)$

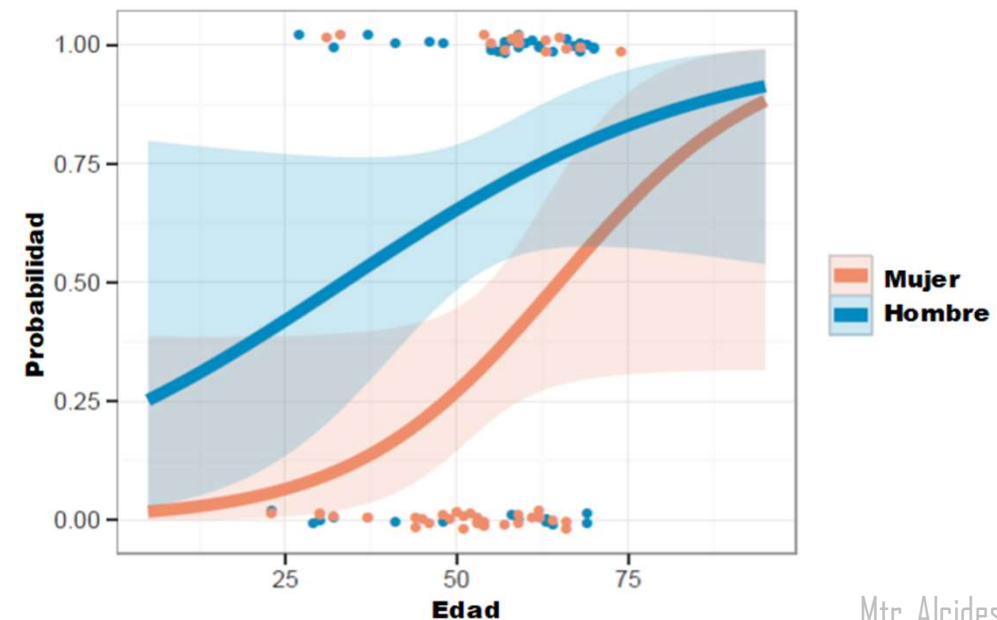
y el parámetro de dispersión es $\phi = 1$.

4. Tipos de GLMs

- Además, $b(\theta_i) = \ln(1 + e^{\theta_i}) = \ln\left(\frac{1}{1 - p_i}\right)$, $G(y_i, \phi) = 1$
- El enlace canónico $g(p_i) = \ln\left(\frac{p_i}{1 - p_i}\right)$

es llamado el logit de p_i . Los GLMs que usan el enlace logit son llamados modelos logit o logísticos.

- En la figura se muestra un ejemplo de un gráfico logit condicional. En él se muestran los puntos separados de las edades y las curvas ajustadas, ambas estratificadas por género.



4. Tipos de GLMs

4.4. Logístico (Y es Binomial)

- La variable Y_i tiene distribución Binomial con parámetros n y p_i . Su función de densidad es

$$\begin{aligned} f(y_i, n, p_i) &= \binom{n}{y_i} p^{y_i} (1-p)^{n-y_i} \\ &= \binom{n}{y_i} \exp \left[y_i \ln \left(\frac{p_i}{1-p_i} \right) - \ln \left(\frac{1}{1-p_i} \right)^n \right] \end{aligned}$$

para todo $0 < p_i < 1$ y $y_i \in \{0, 1, \dots, n\}$. Los parámetros naturales son $\theta_i = \ln \left(\frac{p_i}{1-p_i} \right)$ y el parámetro de dispersión es $\phi = 1$. Además,

$$b(\theta_i) = \ln(1 + e^{\theta_i})^n = \ln \left(\frac{1}{1-p_i} \right)^n, \quad G(y_i, \phi) = \binom{n}{y_i}$$

4. Tipos de GLMs

- El enlace canónico

$$g(p_i) = \ln\left(\frac{p_i}{1-p_i}\right)$$

es llamado el logit de p_i .

Los GLMs que usan el enlace logit son llamados modelos logit o logísticos.

FINESI

Modelos Discretos

IV Semestre



GRACIAS

<https://aulavirtual2.unap.edu.pe/>

