Modelos lineales y modelos lineales generalizados

Rolando Gonzales Martinez, PhD

Fellow postdoctoral Marie Skłodowska-Curie

Universidad de Groningen (Países Bajos)

Investigador (researcher)

Iniciativa de Pobreza y Desarrollo Humano de la Universidad de Oxford (UK)

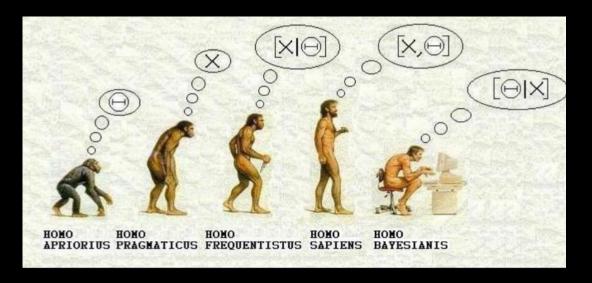
Contenido del curso

(4) Estimación Bayesiana

- Fundamentos de la inferencia Bayesiana.
- Teorema de Bayes.
- Métodos de MCMC (Markov Chain Monte Carlo) para estimación Bayesiana.
- Laboratorio: Estimación Bayesiana de modelos lineales y modelos lineales generalizados

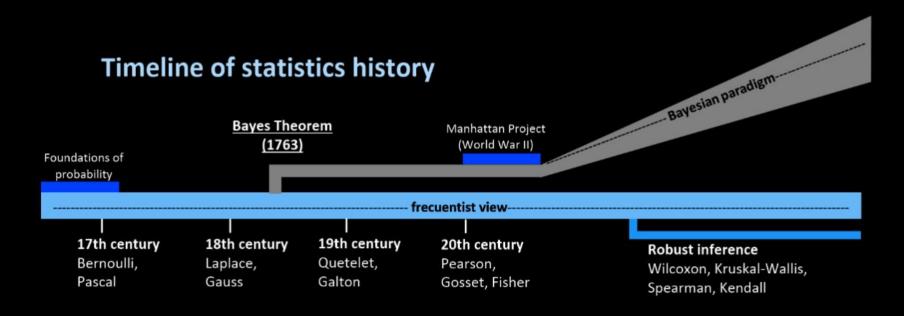
Estimación Bayesiana de modelos lineales

- La estimación Bayesiana no es simplemente un "método adicional o diferente" de estimación:
- El enfoque Bayesiano es un paradigma estadístico diferente.
- Epistemológicamente, un paradigma científico diferente.



Enfoque Bayesiano

Antes de la segunda guerra mundial se utilizaban conjugados naturales

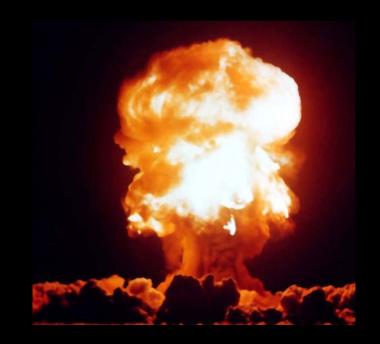


Enfoque Bayesiano

Antes de la segunda guerra mundial se utilizaban conjugados naturales

- Métodos de Monte Carlo

 desarrollados durante el El
 Proyecto Manhattan— permitieron aproximar la integrales multidimensionales del análisis Bayesiano.
- El crecimiento exponencial del software y hardware computacional han hecho el uso de los métodos de integración de Monte Carlo más accesible.



Enfoque Bayesiano

La regla de Bayes surge de los axiomas de probabilidad y no es una materia de controversia.

La división trata sobre la interpretación filosófica de probabilidad P(A):

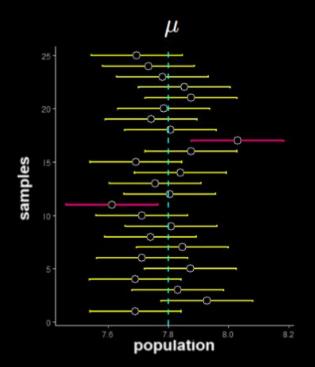
Para un frecuentista, P(A) es una frecuencia de largo plazo:

$$\mathbb{P}(A) := \lim_{n \to \infty_+} \frac{n_A}{n}$$

 Para un Bayesiano, P(A) es cualquier conocimiento/información sobre el evento A, además del contenido en los datos, incluyendo la incertidumbre sobre A.

Enfoque frecuentista:

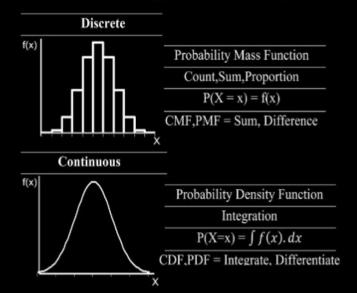
- Los datos son aleatorios
- Los parámetros son (puntos) fijos



Enfoque Bayesiano:

- Los datos son fijos
- Los parámetros son aleatorios

$$\mu \sim \mathcal{D}(\theta_1, \theta_2, ... \theta_d)$$
, e.g. $\mu \sim \mathcal{N}(\theta_\mu, \sigma_\mu^2)$, $\theta_\mu \sim \mathcal{E}(\lambda_{\theta_\mu})$



Estimación Bayesiana de modelos lineales y modelos lineales no generalizados

- Basada en el teorema de Bayes
- Estimación con conjugados naturales
- Estimación con MCMC (Markov Chain Monte Carlo: Monte Carlo con Cadenas de Markov)
- Estimación MCMCMC (MC3)

Función de verosimilitud y teorema de Bayes

$$L(\boldsymbol{\theta}|\mathbf{X}) = \prod_{i=1}^{n} p(\mathbf{X}|\boldsymbol{\theta})$$

$$\ell(\boldsymbol{\theta}|\mathbf{X}) = \log(L(\boldsymbol{\theta}|\mathbf{X})) \qquad p(\boldsymbol{\theta}|\mathbf{X}) = p(\mathbf{X}|\boldsymbol{\theta}) \frac{p(\boldsymbol{\theta})}{p(\mathbf{X})}$$

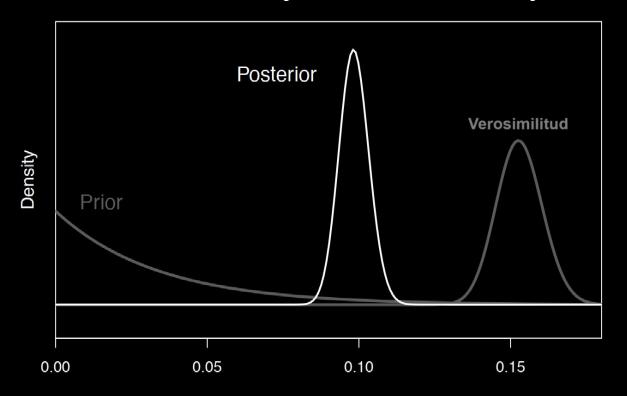
$$\dot{\ell}(\boldsymbol{\theta}|\mathbf{X}) = \frac{\partial}{\partial \boldsymbol{\theta}} \ell(\boldsymbol{\theta}|\mathbf{X})$$

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \stackrel{\mathcal{P}}{\to} \mathcal{N}(\mathbf{0}, \Sigma_{\boldsymbol{\theta}}) \qquad \pi(\boldsymbol{\theta}|\mathbf{X}) = \frac{p(\boldsymbol{\theta})L(\boldsymbol{\theta}|\mathbf{X})}{p(\mathbf{X})}$$

$$\pi(\boldsymbol{\theta}|\mathbf{X}) \propto p(\boldsymbol{\theta})L(\boldsymbol{\theta}|\mathbf{X})$$

Probabilidad posterior ∝ Probabilidad prior × función de verosimilitud

Función de verosimilitud y teorema de Bayes



Probabilidad posterior ∝ Probabilidad prior × función de verosimilitud

Equivalencia asintótica entre los estimadores puntuales Bayesianos y los estimadores máximo verosímiles

Los estimadores bayesianos son asintóticamente equivalente a estimadores máximo verosímiles, si se emplean priors difusos (con varianza muy grande varianza) o uniformes (no informativos)

$$\sqrt{n}(\tilde{\theta} - \hat{\theta}) \underset{n \to \infty}{\longrightarrow} 0$$

$$M(\theta) = \underset{\theta}{\operatorname{argmax}} \pi(\theta | \mathbf{D})$$

Sintetizando la distribución posterior

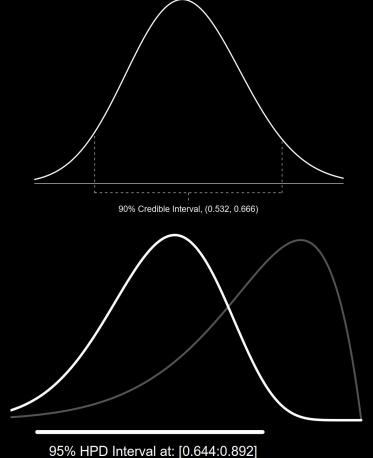
La media de la distribución posterior, intervalos de credibilidad y regiones de mayor densidad posterior (highest posterior density regions, HPD):

$$E[\theta|\mathbf{X}] = \int_{-\infty}^{\infty} \theta p(\theta|\mathbf{X}) d\theta$$

$$1 - \alpha = \int_{C} p(\boldsymbol{\theta}|\mathbf{X}) d\theta$$

$$1 - \alpha = \int_{\theta:\pi(\theta|\mathbf{x}) > k} \pi(\theta|\mathbf{x}) d\theta$$

$$C = \{\theta: \pi(\theta|\mathbf{x}) \ge k\}$$



Beta(15,2) prior in grey

Ejemplo para una distribución exponencial

$$p(X|\theta) = \theta e^{-\theta X}$$

$$p(\theta) = 1/\theta$$

$$\pi(\theta|\mathbf{X}) \propto p(\theta)L(\theta|\mathbf{X}) = \left(\frac{1}{\theta}\right)\theta^n \exp\left[-\theta \sum_{i=1}^n x_i\right]$$

$$= \theta^{n-1} \exp\left[-\theta \sum_{i=1}^n x_i\right]$$

$$\pi(\theta|\mathbf{X}) = \frac{(\sum x_i)^n}{\Gamma(n)}\theta^{n-1} \exp\left[-\theta \sum x_i\right]$$

$$\frac{\alpha}{2} = \int_0^L \pi(\theta|\mathbf{X})d\theta \qquad \frac{\alpha}{2} = \int_0^\infty \pi(\theta|\mathbf{X})d\theta$$

Ejemplo de estimación conjugada Beta-Binomial

$$egin{aligned} heta &\sim \mathrm{Beta}(lpha,eta) & x \mid heta &\sim \mathrm{Binomial}(n, heta) \ p(heta \mid \sum_{i=1}^n x_i) \propto p(\sum_{i=1}^n x_i \mid heta) \cdot p(heta) \end{aligned}$$

$$p(heta\mid \sum_{i=1}^n x_i) \propto heta^{\sum_{i=1}^n x_i + lpha - 1} (1- heta)^{n-\sum_{i=1}^n x_i + eta - 1}$$

$$heta \mid \sum_{i=1}^n x_i \sim \operatorname{Beta}(lpha + \sum_{i=1}^n x_i, eta + n - \sum_{i=1}^n x_i)$$

Algoritmos Markov Chain Monte Carlo

Los algoritmos MCMC son métodos para generar una cadena de muestras de un espacio de parámetros que, después de un tiempo suficiente (conocido como el período de "quemado" o burn-in), se pueden considerar como muestras de la distribución de interés:

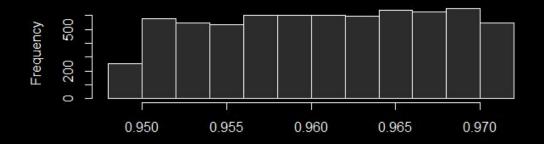
- Los algoritmos MCMC se basan en una cadena de Markov en la que la probabilidad de moverse a un nuevo estado depende solamente del estado actual y no de los estados anteriores
- La cadena de Markov se diseña de tal manera converge a una distribución de equilibrio o estacionaria.
- La cadena debe ser ergódica, lo que significa que debe ser posible alcanzar cualquier estado desde cualquier otro estado en un número finito de pasos, garantizando que las muestras eventualmente representen la distribución objetivo.

Algoritmos Markov Chain Monte Carlo

Integración de Monte Carlo:

$$I[e,\pi] = \int_{e}^{\pi} \arctan(\theta^{\frac{1}{3}}) \mathcal{C}(\theta|\mu=3,\sigma=2) d\theta$$

$$\mathcal{C}(\theta|\mu,\sigma) = \frac{1}{\pi\sigma} \frac{1}{1+(\frac{\theta-\mu}{\sigma})^2}, -\infty < \theta, \mu < \infty, 0 < \sigma$$



Algoritmos Markov Chain Monte Carlo: Gibbs sampling

Gibbs Sampling es un algoritmo de MCMC (Markov Chain Monte Carlo) que se utiliza para generar muestras de una distribución conjunta cuando las distribuciones condicionales son conocidas y pueden ser muestreadas fácilmente

- 1. Inicializar $X_1^{(0)}, X_2^{(0)}, \ldots, X_n^{(0)}$ con algún valor.
- 2. Iterar para t = 1, 2, ..., T:
 - Muestrear $X_1^{(t)} \sim P(X_1|X_2^{(t-1)},X_3^{(t-1)},\dots,X_n^{(t-1)})$
 - Muestrear $X_2^{(t)} \sim P(X_2|X_1^{(t)}, X_3^{(t-1)}, \dots, X_n^{(t-1)})$
 - •
 - Muestrear $X_n^{(t)} \sim P(X_n|X_1^{(t)}, X_2^{(t)}, \dots, X_{n-1}^{(t)})$

Algoritmos Markov Chain Monte Carlo: Gibbs sampling

Con priors no informativos:

$$egin{aligned} y_i &= eta_0 + eta_1 x_{i1} + eta_2 x_{i2} + \ldots + eta_p x_{ip} + \epsilon_i \ \epsilon_i &\sim N(0,\sigma^2) \end{aligned}$$

1. Distribución condicional de eta dado σ^2 y y, X: $eta|\sigma^2, y, X \sim N\left((X'X)^{-1}X'y, \sigma^2(X'X)^{-1}
ight)$

2. Distribución condicional de σ^2 dado β y y, X: $\sigma^2 | \beta, y, X \sim \text{Inv-Gamma}\left(\frac{n}{2}, \frac{1}{2}(y - X\beta)'(y - X\beta)\right)$

Algoritmos Markov Chain Monte Carlo: Gibbs sampling

1. Priors Informativos:

- Prior para eta: $eta \sim N(\mu_{eta}, \Sigma_{eta})$
- Prior para σ^2 : $\sigma^2 \sim \text{Inv-Gamma}(\alpha_0, \beta_0)$

2. Distribuciones Condicionales:

• Condicional de β dado σ^2 , y, y X:

$$eta|\sigma^2,y,X\sim N(\mu_{eta|y},\Sigma_{eta|y}) ~~~~ \Sigma_{eta|y} = \left(rac{1}{\sigma^2}X^TX + \Sigma_{eta}^{-1}
ight)^{-1}$$

• Condicional de σ^2 dado β , y, y X:

$$\sigma^2 | eta, y, X \sim ext{Inv-Gamma} \left(lpha_n, eta_n
ight)$$

$$eta_n = eta_0 + rac{1}{2}(y-Xeta)^T(y-Xeta)$$

 $\mu_{eta|y} = \Sigma_{eta|y} \left(rac{1}{\sigma^2} X^T y + \Sigma_{eta}^{-1} \mu_{eta}
ight)$

 $\alpha_n = \alpha_0 + \frac{n}{2}$

Algoritmos Markov Chain Monte Carlo: Metropolis

- 1. Inicialización: Escoger un valor inicial para θ , $\theta^{(0)}$.
- 2. **Propuesta**: Generar un valor propuesto θ^* desde una distribución de propuesta simétrica $q(\theta^*|\theta^{(t)})$, donde $\theta^{(t)}$ es el valor actual de la cadena. En el caso del algoritmo de Metropolis, se suele utilizar una distribución normal centrada en el valor actual: $\theta^* \sim N(\theta^{(t)}, \sigma^2)$
- 3. Cálculo de la Probabilidad de Aceptación: $lpha = \min\left(1, rac{\pi(heta^*)}{\pi(heta^{(t)})}
 ight)$

Dado que la propuesta es simétrica, $q(heta^*| heta^{(t)})=q(heta^{(t)}| heta^*)$, y estos términos se cancelan en la razón de aceptación.

4. Aceptación o Rechazo: Aceptar θ^* con probabilidad α :

$$heta^{(t+1)} = egin{cases} heta^* & ext{con probabilidad } lpha \ heta^{(t)} & ext{con probabilidad } 1-lpha \end{cases}$$

5. Iteración: Repetir los pasos 2-4 por el número de iteraciones deseado.

Algoritmos Markov Chain Monte Carlo: Metropolis-Hastings

- 1. Inicialización: Escoger un valor inicial para θ , $\theta^{(0)}$.
- 2. **Propuesta**: Generar un valor propuesto θ^* desde una distribución de propuesta $q(\theta^*|\theta^{(t)})$, donde $\theta^{(t)}$ es el valor actual de la cadena.
- 3. Cálculo de la Probabilidad de Aceptación:

$$lpha = \min\left(1, rac{\pi(heta^*)q(heta^{(t)}| heta^*)}{\pi(heta^{(t)})q(heta^*| heta^{(t)})}
ight)$$
 $egin{aligned} \pi(heta) ext{ es la distribución objetivo.} \ q(heta^*| heta^{(t)}) ext{ es la distribución de propuesta} \end{aligned}$

4. Aceptación o Rechazo: Aceptar θ^* con probabilidad α :

$$heta^{(t+1)} = egin{cases} heta^* & ext{con probabilidad } lpha \ heta^{(t)} & ext{con probabilidad } 1-lpha \end{cases}$$

5. Iteración: Repetir los pasos 2-4 por el número de iteraciones deseado.

Estimación Bayesiana del MLRM

Con priors no informativos los estimadores Bayesianos coinciden con los estimadores máximo verosímiles:

$$\mathbb{P}(\beta) \propto c \text{ y } \mathbb{P}(\sigma^2) \propto \sigma^{-1}$$

$$\mathbb{P}(\beta, \sigma^2) \propto \mathcal{L}(\beta, \sigma^2 | \mathbf{X}, \mathbf{y}) \mathbb{P}(\beta) \mathbb{P}(\sigma^2)$$

$$\propto \sigma^{-n-1} \exp \left[-\frac{1}{2\sigma^2} (\hat{\sigma}^2(n-k) + (\beta - \hat{\beta})' \mathbf{X}' \mathbf{X} (\beta - \hat{\beta})) \right]$$

Estimación Bayesiana del MLRM

Con priors informativos:

$$\beta \sim \mathcal{N}_k(\beta_0, B_0), \quad \sigma^2 \sim \mathcal{IG}(\alpha_0/2, \delta_0/2),$$

$$\beta | \sigma^2, y \sim \mathcal{N}(\overline{\beta}, B_1), \qquad \sigma^2 | \beta, y \sim \mathcal{IG}(\alpha_1/2, \delta_1/2)$$

$$B_1 = [s^{-2}X'X + B_0^{-1}]^{-1},$$

$$\overline{\beta} = B_1[\sigma^{-2}X'y + B_0^{-1}\beta_0],$$

$$\alpha_1 = \alpha_0 + n,$$

$$\delta_1 = \delta_0 + (y - X\beta)'(y - X\beta)$$

Algoritmos Markov Chain Monte Carlo: Metropolis-Hastings para el modelo logit

- 1. **Inicialización**: Escoger un valor inicial para β .
- 2. **Propuesta**: Generar un valor propuesto β^* a partir de una distribución propuesta $q(\beta^*|\beta)$.

$$egin{aligned} P(y_i = 1|X_i,eta) &= rac{1}{1+\exp(-X_ieta)} \ \operatorname{logit}(P(y_i = 1|X_i,eta)) &= X_ieta \ L(eta|X,y) &= \prod_{i=1}^n \left(rac{1}{1+\exp(-X_ieta)}
ight)^{y_i} \left(1 - rac{1}{1+\exp(-X_ieta)}
ight)^{1-y_i} \end{aligned}$$

- 3. Cálculo de la Probabilidad de Aceptación: $\alpha = \min\left(1, \frac{L(\beta^*|X,y)\pi(\beta^*)q(\beta|\beta^*)}{L(\beta|X,y)\pi(\beta)q(\beta^*|\beta)}\right)$ donde $\pi(\beta)$ es la distribución a priori de β (pued
 - donde $\pi(eta)$ es la distribución a priori de eta (puede ser no informativa).
- 4. **Aceptación o Rechazo**: Aceptar β^* con probabilidad α , de lo contrario, mantener el valor actual de β .
- 5. Iteración: Repetir los pasos 2-4 por el número de iteraciones deseado.

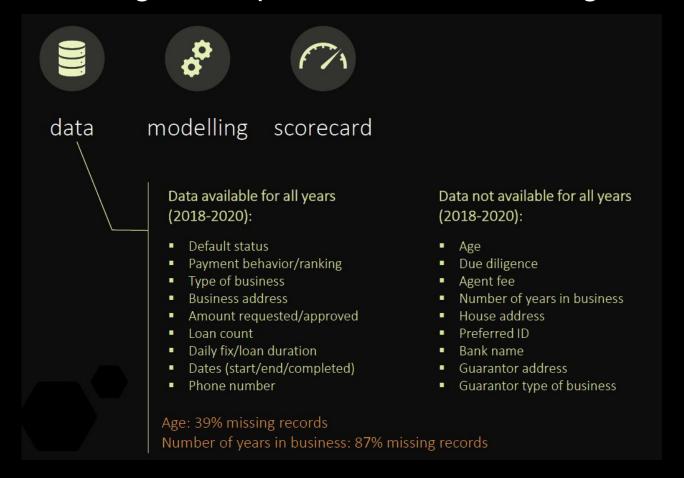
El modelo logit en el contexto de credit scoring

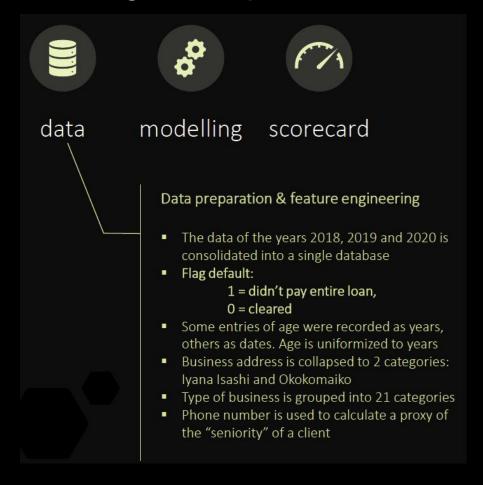
En el contexto de credit scoring, es usual estimar los modelos logit transformando las variables explicativas a WOE:

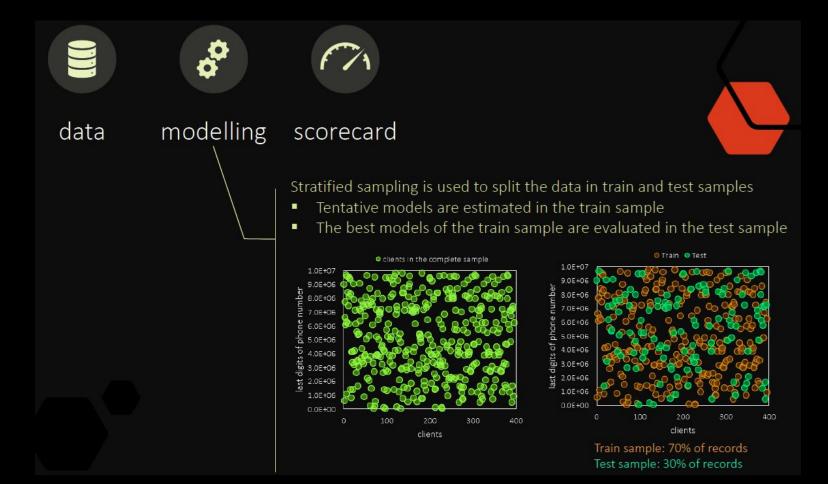
$$WoE_{x=i} = \ln\left(\frac{\% \ of \ y = 0 \ where \ x = i}{\% \ of \ y = 1 \ where \ x = i}\right)$$

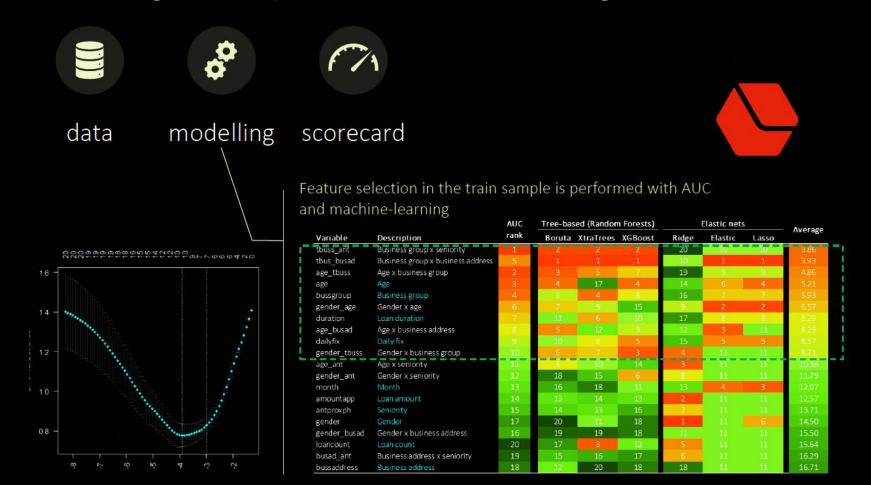
Los valores de WOE proporcionan una escala logarítmica que ayuda a interpretar la fuerza y la dirección de la relación entre las categorías de la variable predictora y la probabilidad del evento objetivo.

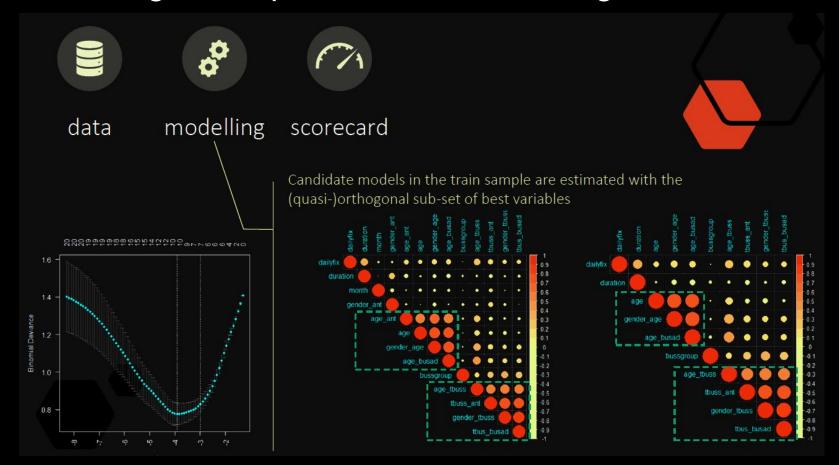
- Valores Positivos de WOE: la categoría tiene una mayor proporción de no defaults que defaults (menos riesgo).
- Valores Negativos de WOE: la categoría tiene una mayor proporción de defaults que no defaults (más riesgo).

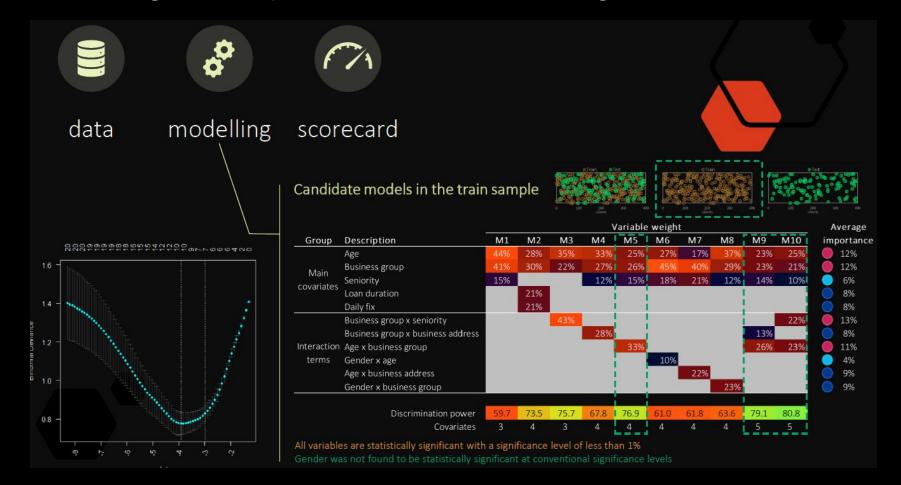


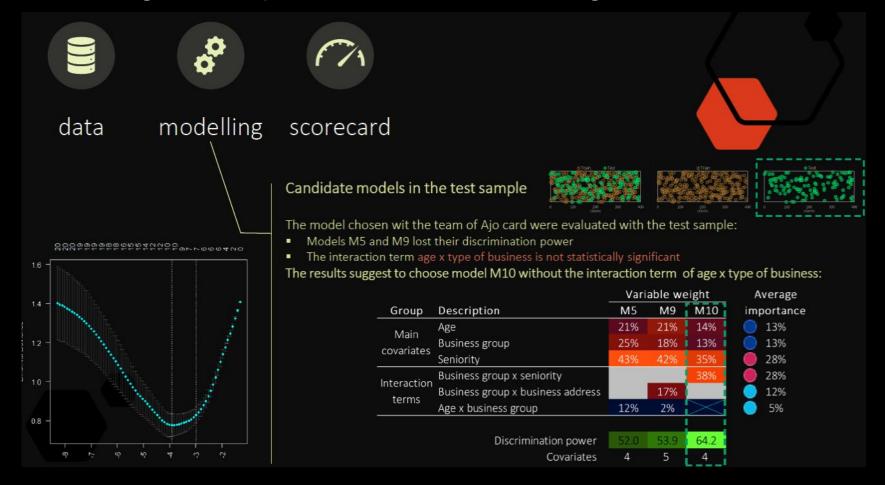












Conversión de resultados del modelo logit a scores crediticios

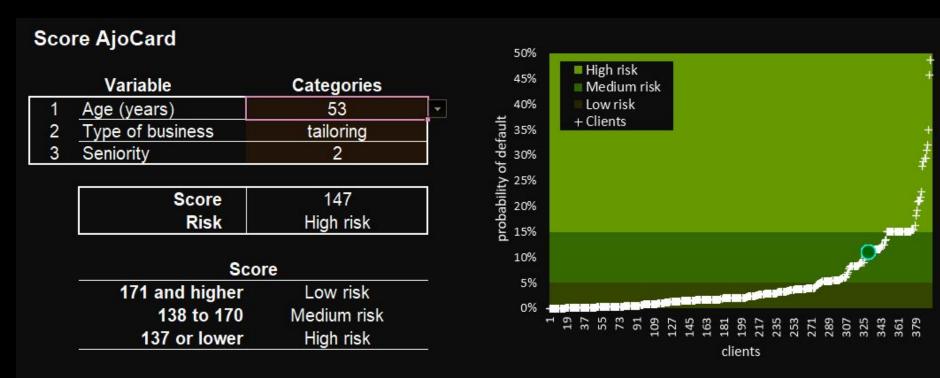
El resultado de la estimación Xβ del modelo logit se transforma a scores crediticios con la fórmula:

Scores = offset + factor \mathbf{X} (-X β)

Con odds (momios) de 50:1 es decir, 50 créditos de no default por un default, y se doblan dichos momios cada 20 puntos del score, el offset y el factor son iguales a 600 y 20, respectivamente.

Vease: Siddiqi, N. (2017). Intelligent credit scoring: Building and implementing better credit risk scorecards. John Wiley & Sons.

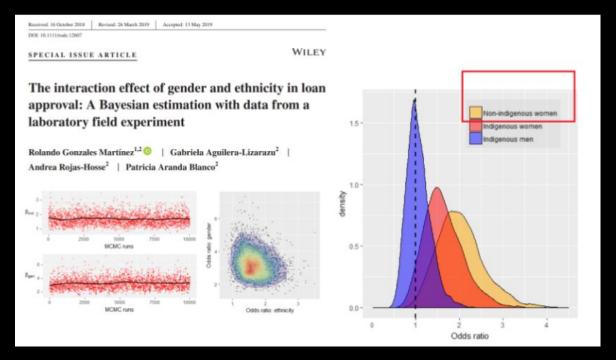
Conversión de resultados del modelo logit a scores crediticios



Laboratorio

- MLMLG_400: Estimación con conjugados naturales Beta-Binomial e intervalos de credibilidad
- MLMLG_401: Integración de Monte Carlo de una distribución Cauchy
- MLMLG_402: Estimación MCMC con Gibbs Sampling del modelo lineal
- MLMLG_403: Estimación MCMC con el algoritmo de Metropolis de modelos lineales generalizados: modelo logit
- MLMLG_404: Estimación MCMC con el algoritmo Metropolis-Hastings de modelos lineales generalizados: modelo logit
- MLMLG_405: Estimación MCMC de un modelo logit aplicado a credit scoring

Ejemplo adicional de aplicación de métodos Bayesianos en MLG



Experimento de laboratorio-campo y modelo Bayesiano logit mixto

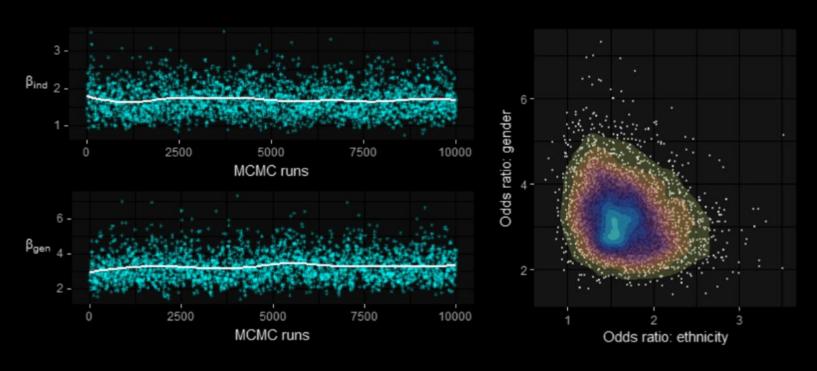
Discriminacion estadistica o discriminacion basada en preferencias (taste-based discrimination).



$$y_{i(h)} \begin{vmatrix} \mathbf{x}_{i(h)} \sim (y_{i(h)} | \mathbf{x}_{i(h)}), & y \in \{0,1\}; & i = 1, \dots, n_h \\ (y_{i(h)} | \mathbf{x}_{i(h)}) = F \left(\mathbf{x}_{i(h)}^{\mathsf{T}} \beta_{N_h} \right)^{\mathsf{y}} \left[1 - F \left(\mathbf{x}_{i(h)}^{\mathsf{T}} \beta_{N_h} \right) \right]^{\mathsf{y}-1} \\ F \begin{vmatrix} (\mathbf{x}^{\mathsf{T}} \beta) = P \left(Y = 1 | \mathbf{x}^{\mathsf{T}} \beta \right) = \frac{\exp(\mathbf{x}^{\mathsf{T}} \beta)}{1 + \exp(\mathbf{x}^{\mathsf{T}} \beta)} \\ \mathbf{x}^{\mathsf{T}} \beta_{N_h} = \mathbf{x}^{\mathsf{T}} \beta + u_{0h} \\ \beta_k \sim N \left(0, \nu_{\beta_k} \right), & k = 1, \dots, p \\ u_{0h} \begin{vmatrix} \tau_0 \sim N \left(0, \nu_0 \right), & h = 1, \dots, N_h \\ \tau_0 \sim IG \left(\nu_0, \nu 1 \right) \end{vmatrix}$$

Experimento de laboratorio-campo y modelo Bayesiano logit mixto

Estimadores para sexo y etnicidad



Experimento de laboratorio-campo y modelo Bayesiano logit mixto

Diferencias entre mujeres no indígenas y mujeres indígenas

