

46114: Estadística Multivariada y Datos  
Categoricos  
Sesión 2: Modelos lineales de regresión

Juan Carlos Martínez-Ovando

`juan.martinez.ovando@itam.mx`

Enero 21, 2016

## Modelos lineales de regresión: Introducción

Los modelos lineales de regresión son motivados por la necesidad de modelar la realización de una variable de interés,  $Y$ , con un conjunto de variables auxiliares,  $X_1, \dots, X_p$ , a través de la distribución condicional,

$$\mathbb{P}(Y|X_1, \dots, X_p).$$

Formas de definir la distribución condicional anterior:

- **Regresión en media** tal que,

$$\mathbb{E}(Y|X_1, \dots, X_p) = \beta_1 X_1 + \dots + \beta_p X_p,$$

y  $\text{var}(Y|X_1, \dots, X_p) = \sigma^2$  (constante en las  $X_j$ s).

- **Regresión en varianza** tal que,

$$\text{var}(Y|X_1, \dots, X_p) = \beta_1 X_1 + \dots + \beta_p X_p,$$

y  $\mathbb{E}(Y|X_1, \dots, X_p) = \mu$  (constante en las  $X_j$ s).

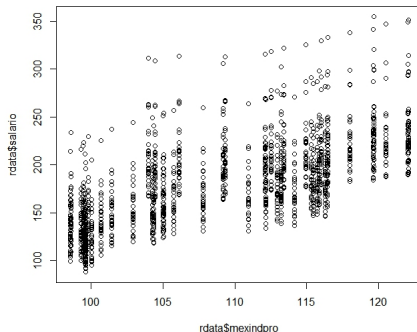
# Modelos lineales de regresión: Introducción

¿Cómo elegir el modelo?

- ▶ Inspección gráfica.
- ▶ Intuición... :)

## Ejemplo

- ▶ Salario base de cotización en México.
- ▶ Nivel de producción industrial (periodo anterior).



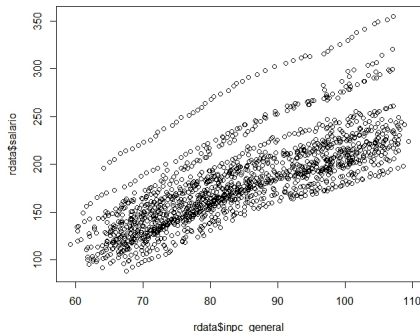
# Modelos lineales de regresión: Introducción

¿Cómo elegir el modelo?

- ▶ Inspección gráfica.
- ▶ Intuición... :)

## Ejemplo

- ▶ Salario base de cotización en México.
- ▶ Nivel del INPC (periodo anterior).



# Modelos lineales de regresión: Introducción

La distribución condicional de  $Y$  en  $X_1, \dots, X_p$  proviene de,

$$\mathbb{P}(Y, X_1, \dots, X_p) = \mathbb{P}(Y|X_1, \dots, X_p) \times \mathbb{P}(X_1, \dots, X_p).$$

Bajo lo anterior, el enfoque bayesiano de inferencia invoca a la noción de intercambiabilidad (conjunta) en  $Y, X_1, \dots, X_p$ , i.e. existen  $\theta$  y  $\gamma$  tales que

$$\begin{aligned}\mathbb{P}(\mathbf{y}\mathbf{x}_1, \dots, \mathbf{y}\mathbf{x}_n) &= \int \prod_{i=1}^n \mathbb{P}(\mathbf{y}\mathbf{x}_i | \theta, \gamma) \pi(\theta, \gamma) d\theta d\gamma \\ &= \int \prod_{i=1}^n \mathbb{P}(y_i | x_{i1}, \dots, x_{ip}, \theta) \mathbb{P}(x_{i1}, \dots, x_{ip} | \gamma) \\ &\quad \times \pi(\theta | \gamma) d\theta \pi(\gamma) d\gamma,\end{aligned}$$

donde  $\mathbf{y}\mathbf{x}_j = (y_j, x_{j1}, \dots, x_{jp})$ .

# Modelos lineales de regresión: Introducción

- ▶ En los modelos lineales de regresión, de lo anterior sólo es de interés modelar la relación de dependencia

$$(X_1, \dots, X_p) \rightarrow Y.$$

- ▶ Los modelos que involucran la modelación conjunta de  $Y$  y las  $X_j$ s corresponden a *modelos gráficos*.

Considerando los modelos de regresión convencionales, sólo es de interés estudiar  $\mathbb{P}(Y|X_1, \dots, X_p, \theta)$  y  $\pi(\theta|\gamma)$ , dando origen a la noción de *intercambiabilidad parcial*,

$$\mathbb{P}(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n) = \int \prod_{i=1}^n \mathbb{P}(y_i | X_{i1}, \dots, X_{ip}, \theta) \pi(\theta | \gamma) d\theta,$$

donde  $\mathbf{x}_j = (x_{j1}, \dots, x_{jp})$ .

# Modelos lineales de regresión: Introducción

Concentrándonos en  $\mathbb{P}(y_i|X_{i1}, \dots, X_{ip}, \theta)$ , tenemos los siguientes modelos (entre otros):

- **Modelo lineal en media**

$$\prod_{i=1}^n \mathbb{P}(y_i|X_{i1}, \dots, X_{ip}, \theta) = N_n(\mathbf{Y}|\mathbf{X}'\theta, \sigma^2 \mathbf{I}).$$

- **Modelo lineal en varianza**

$$\prod_{i=1}^n \mathbb{P}(y_i|X_{i1}, \dots, X_{ip}, \theta) = N_n(\mathbf{Y}|\boldsymbol{\mu}, \text{diag}\{\mathbf{x}'_1\theta, \dots, \mathbf{x}'_n\theta\}).$$

- **Modelo lineal generalizado**

$$\prod_{i=1}^n \mathbb{P}(y_i|X_{i1}, \dots, X_{ip}, \theta) = \mathbb{P}_n(\mathbf{Y}|g(\mathbf{X}'\theta), \gamma).$$

# Modelos lineales de regresión: Inferencia bayesiana

Nos concentraremos ahora en el **modelo de regresión lineal en media**, con la parametrización en términos de la *precisión* en lugar de la *varianza*, i.e.

$$N_n(\mathbf{Y}|\mathbf{X}'\theta, \sigma^2\mathbf{I}) = N_n(\mathbf{Y}|\mathbf{X}'\theta, \lambda\mathbf{I}),$$

donde  $\lambda = 1/\sigma^2$ .

Los **parámetros** del modelo son  $\theta \in \mathbb{R}^p$  y  $\lambda > 0$ .

La **distribución inicial** sobre los parámetros se elige como la distribución conjugada, con

$$\pi(\theta, \lambda) = N_p(\theta|\theta_0, \lambda S_0) \times \text{Ga}(\lambda|\alpha_0, \beta_0).$$

Los **hiper-parámetros**  $\theta_0 \in \mathbb{R}^p$ ,  $S_0$  matriz  $p \times p$  positivo definida y simétrica,  $\alpha_0, \beta_0 > 0$  son *determinados a priori*.



# Modelos lineales de regresión: Inferencia bayesiana

## Distribución posterior

$$\begin{aligned}\pi(\theta, \lambda | \mathbf{y}\mathbf{x}_1, \dots, \mathbf{y}\mathbf{x}_n) &= \pi(\theta | \mathbf{y}\mathbf{x}_1, \dots, \mathbf{y}\mathbf{x}_n, \lambda) \times \pi(\lambda | \mathbf{y}\mathbf{x}_1, \dots, \mathbf{y}\mathbf{x}_n) \\ &= \mathcal{N}_p(\theta | \theta_n, \lambda S_n) \times \text{Ga}(\lambda | \alpha_n, \beta_n)\end{aligned}$$

con

$$S_n = S_0 + \mathbf{X}'\mathbf{X},$$

$$\theta_n = S_n^{-1}(S_0\theta_0 + \mathbf{X}'\mathbf{Y}),$$

$$\alpha_n = \alpha_0 + n/2,$$

$$\beta_n = \beta_0 + 1/2(\mathbf{Y} - \mathbf{X}'\theta_n)'\mathbf{Y} + 1/2(\theta_0 - \theta_n)'S_0\theta_0.$$

# Modelos lineales de regresión: Inferencia bayesiana

## Distribución predictiva

$$\begin{aligned}\mathbb{P}(\mathbf{Y}^f | \mathbf{X}^f, \mathbf{y}\mathbf{x}_1, \dots, \mathbf{y}\mathbf{x}_n) &= \int \int N_{n^f}(\mathbf{Y}^f | \mathbf{X}^{f'}\theta, \lambda I_{n^f}) \\ &\quad \times \Pi(d\theta, d\lambda | \mathbf{y}\mathbf{x}_1, \dots, \mathbf{y}\mathbf{x}_n) \\ &= \text{St}_p(\mathbf{Y}^f | \mathbf{X}^{f'}\theta_n, C(\mathbf{X}^f)\alpha_n/\beta_n, 2\alpha_n),\end{aligned}$$

donde

$$C(\mathbf{X}^f) = 1 + \mathbf{X}^f(\mathbf{X}^{f'}\mathbf{X}^f + S_n)^{-1}\mathbf{X}^{f'},$$

con  $\theta_n$ ,  $S_n$ ,  $\alpha_n$  y  $\beta_n$  dadas como antes.