

## ESTADISTICA 2

### MUESTREO

- Pasos:

- 1) Definir problema
- 2) Relacion problema - Estimacion [M, T, P, A]
- 3) Definir: Elemento, Poblacion, Unidad de muestreo y Marco
- 4) Definir Metodologia de Muestreo [MAS, MAE, MS, MC, MVE]
- 5) Calcular tamaño de muestra con la muestra piloto. NC
- 6) Seleccionar Muestra
- 7) Hacer Inferencia

- Definicion: Parte de la estadística que se encarga del estudio de métodos para la obtencion de una muestra de una poblacion objeto de estudio

Muestra = Parte de la poblacion a partir de la cual hacemos inferencia (sin remplazo)

↓  
Estimamos

↓  
M = Medida

T = Total  $\mu$  suma

P = Proporción

A = Total  $p$  atributo

↳ sacar conclusiones sobre toda la poblacion a partir de la Muestra

↳ Coleccion de unidades de muestreo obtenidas a partir de un marco o marcos (representativa)

**Poblacion**: Conjunto de elementos acerca de los cuales se realizara la inferencia. Estudiantes de primer semestre Fac Minas, tiendas visites en medellin

**Unidad de Muestreo**: Colecciones no traslapadas de elementos de la poblacion que cubren la poblacion completa. Asignaturas introductorias, Manzanas, Parcelas, intervalos de tiempo

**Elemento**: Objeto sobre el que se toma medicion o conteo. Cada estudiante de primer semestre Fac. Minas, tiendas de barrio

**Marco Muestral**: Lista de unidades de muestreo, Listado asignaturas introductorias, Listado de manzanas

### METODOLOGIAS DE MUESTREO

1) Muestreo Aleatorio Simple (MAS) Independiente

\* Poblacion homogénea con respecto a la característica de interés

\* Principio de equiprobabilidad

## 2) Muestreo Aleatorio Estratificado (MAE):

\* Población heterogénea

\* La población es dividible en subpoblaciones homogéneas (Estratos) no traslapados

\* De cada estrato se selecciona una MAS

\* Aumenta la precisión.



Muestra por estrato



MAE

## 3) Muestreo Sistemático (MS)

\* Para poblaciones en movimiento

\* Estimaciones similares a las del MAS

Control de calidad industrial

- La no respuesta causa el muestreo

Aleatorio (S)

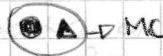
↳ equidistantes

↳ 10, 15, 20, 25

## 4) Muestreo por Conglomerados (MC)

\* La población está conformada por subpoblaciones (conglomerados) homogéneas internamente, distintas de Unidades de Muestreo

\* Se selecciona una MAS de conglomerados y se muestrean todas las unidades de los conglomerados seleccionados



MC

## 5) Muestreo en Varias Etapas (MVE)

\* Disminuye error de MC

\* Aumenta Varianza

Estratificado → Conglomerado → Simple

## PLANIFICACION ENCUESTA

1) Declaración de Objetivos: Concisos, Simples

2) Población objetivo: Población factible seleccionada cuidadosamente

3) Marco: La lista y la población objetivo coinciden lo máximo, La eficiencia puede aumentar con marcos múltiples

4) Diseño de: # de elementos de la muestra que proporcione suficiente muestreo ~~información~~ información para los objetivos (representativos)

5) Método de Medición: Entrevistas, cuestionarios personales, telefónicos, correos u observación

6) Instrumento de medición: Minimizar la no respuesta y el sesgo por respuesta incorrecta

7) Selección y formación de personal de campo.

8) Pre-test: Muestreo piloto, Calificación entrevistadores, Operatividad

9) Organización: Establecer líneas de autoridad

10) Administración esquemática de los datos

11) Análisis de datos.

## MUESTREO PILOTO

- \* Validar preguntas
- \* Estimar varianza
- \* Determinar tamaño de muestra
- \* Da idea para el límite de error de estimación definitivo

## TAMAÑO DE MUESTRA

- \* Depende del tipo de muestreo seleccionado y del muestreo piloto
- \* Es el número mínimo de elementos que estarán en la muestra
- \* Es el número de elementos que permite hacer inferencia

## NIVEL DE CONFIANZA Y LIMITE DE ERROR DE ESTIMACION (NC) (B)

- \* Son fijados por el investigador (los dos)
- \* B, dos desviaciones estándar del estimador

$$\text{Error de estimación} = |\theta - \hat{\theta}| < B$$

$$P[\text{Error de estimación}] = 1 - \alpha$$

$$NC = 1 - \alpha$$

## OBSERVACIONES GENERALES

- Cada elemento del muestreo provee información de los parámetros de interés
- Poca información impide realizar buenas estimaciones
- Mucha información desperdicio de recursos
- La encuesta controla la variación de los datos
- En la encuesta se debe recolectar información adicional que ayude a la pregunta de interés.
- El mejor diseño da la precisión necesaria en términos de B con un coste mínimo

## MUESTREO ALEATORIO SIMPLE (MAS)

Procedimiento estadístico a partir del cual se selecciona una muestra de tamaño  $n$  de una población de  $N$  unidades, garantizando que cada muestra posible de tamaño  $n$  tenga la misma probabilidad de ser seleccionada (Principio de equiprobabilidad)

- En la práctica una MAS es seleccionada unidad por unidad
- El muestreo se realiza sin remplazo. (No se repiten las unidades)

El MAS se utiliza cuando la población es homogénea con respecto a la característica de interés y cuando las estimaciones se refieren a toda la población y no a subgrupos

### Estadísticos e intervalos

- Media Muestral ( $\bar{y}$ )

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \quad \rightarrow \quad y_i = \text{Observaciones de la variable de interés}$$

$n = \text{tamaño de la muestra}$

$$\text{Var}[\bar{y}] = \frac{\sigma^2}{n} \left( \frac{N-n}{N-1} \right) \rightarrow \text{Factor de corrección para Poblaciones Finitas}$$

$$\hat{\sigma}^2 = \frac{N-1}{N} \cdot s^2 \quad N = \text{tamaño de la población}$$

- Varianza Estimada ( $\hat{\sigma}^2$ )

$$\hat{\text{Var}}[\bar{y}] = \frac{\hat{\sigma}^2}{n} \cdot \frac{N-n}{N-1} = \frac{N-1}{N} \cdot s^2 \cdot \frac{N-n}{N-1}$$

$$\hat{\text{Var}}[\bar{y}] = \frac{s^2}{n} \left( \frac{N-n}{N} \right) \rightarrow \text{Factor de corrección para Poblaciones Finitas}$$

- Intervalos de confianza ( $\mu$ )

↳ Límite en el error de estimación  $B$

$$\bar{y} \pm B_{\alpha}$$

Donde  $B$  depende de  $n$ :

Si  $n \geq 30$

$$B = z_{\frac{\alpha}{2}} \sqrt{\frac{s^2}{n} \left( \frac{N-n}{N} \right)}$$

Entonces

$$\bar{y} \pm z_{\frac{\alpha}{2}} \sqrt{\frac{s^2}{n} \left( \frac{N-n}{N} \right)}$$

Si  $n \leq 30$

$$B = t_{\frac{\alpha}{2}, n-1} \sqrt{\frac{s^2}{n} \left( \frac{N-n}{N} \right)}$$

Entonces

$$\bar{y} \pm t_{\frac{\alpha}{2}, n-1} \sqrt{\frac{s^2}{n} \left( \frac{N-n}{N} \right)}$$

- Total de la población ( $\tau$ )

$$\tau = NM$$

$$\hat{\tau} = N\hat{\mu} = N\bar{y}$$

- Varianza Estimada ( $\tau$ )

$$\text{Var}[\hat{\tau}] = N^2 \frac{s^2}{n} \left( \frac{N-1}{N} \right)$$

- Intervalo de confianza ( $\tau$ )

$$\hat{\tau} \pm B_{\tau}$$

$$N(\bar{y} \pm B_{\mu})$$

Donde  $B_{\tau}$  depende de  $n$

Si  $n > 30$

$$B_{\tau} = z_{\frac{\alpha}{2}} \sqrt{N^2 \frac{s^2}{n} \left( \frac{N-1}{N} \right)}$$

Entonces

$$\hat{\tau} \pm z_{\frac{\alpha}{2}} \sqrt{N^2 \frac{s^2}{n} \left( \frac{N-1}{N} \right)}$$

Si  $n \leq 30$

$$B_{\tau} = t_{\frac{\alpha}{2}, n-1} \sqrt{N^2 \frac{s^2}{n} \left( \frac{N-1}{N} \right)}$$

Entonces

$$\hat{\tau} \pm t_{\frac{\alpha}{2}, n-1} \sqrt{N^2 \frac{s^2}{n} \left( \frac{N-1}{N} \right)}$$

- Proporción Poblacional ( $p$ )

$y_i = \begin{cases} 1, & \text{tiene el atributo} \\ 0, & \text{no tiene el atributo} \end{cases}$

$$\hat{p} = \frac{\sum y_i}{n} = \frac{\# \text{ Unidades con el atributo}}{\# \text{ unidades en la muestra}}$$

- Varianza Estimada ( $p$ )

$$\text{Var}[\hat{p}] = \frac{p(1-p)}{n-1} \left( \frac{N-1}{N} \right)$$

- Intervalo de confianza

$$\hat{p} \pm B_p$$

Donde  $B_p$  depende de  $n$

Si  $n \geq 100$

$$B_p = z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n-1} \left(\frac{N-n}{N}\right)}$$

Entonces

$$\hat{p} \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n-1} \left(\frac{N-n}{N}\right)}$$

Si  $n < 100$

$$B_p = t_{\frac{\alpha}{2}, n-1} \sqrt{\frac{\hat{p}(1-\hat{p})}{n-1} \left(\frac{N-n}{N}\right)}$$

Entonces

$$\hat{p} \pm t_{\frac{\alpha}{2}, n-1} \sqrt{\frac{\hat{p}(1-\hat{p})}{n-1} \left(\frac{N-n}{N}\right)}$$

• Total de la Población (A)

$$\hat{A} = N\hat{p}$$

• Varianza estimada (A)

$$\text{Var}[\hat{A}] = N^2 \cdot \frac{\hat{p}(1-\hat{p})}{n-1} \left(\frac{N-n}{N}\right)$$

• Intervalo de confianza (A)

$$\hat{A} \pm B_A$$

$$N(\hat{p} \pm B_p)$$

Donde  $B_A$  depende de  $n$

Si  $n \geq 100$

$$B_A = z_{\frac{\alpha}{2}} \sqrt{N^2 \frac{(1-\hat{p})\hat{p}}{n-1} \left(\frac{N-n}{N}\right)}$$

Entonces

$$\hat{A} \pm z_{\frac{\alpha}{2}} \sqrt{N^2 \frac{\hat{p}(1-\hat{p})}{n-1} \left(\frac{N-n}{N}\right)}$$

Si  $n < 100$

$$B_A = t_{\frac{\alpha}{2}, n-1} \sqrt{N^2 \frac{(1-\hat{p})\hat{p}}{n-1} \left(\frac{N-n}{N}\right)}$$

Entonces

$$\hat{A} \pm t_{\frac{\alpha}{2}, n-1} \sqrt{N^2 \frac{\hat{p}(1-\hat{p})}{n-1} \left(\frac{N-n}{N}\right)}$$

Conocidos  $B$  y el nivel de confianza  $1 - \alpha$  fijados por el investigador para estimar  $N$ , se busca hallar cual debe ser el tamaño de muestra  $n$ . Para ello partimos del límite de error de estimación  $B$

$$B = z_{\frac{\alpha}{2}} \sqrt{\frac{\sigma^2}{n} \cdot \frac{N-n}{N-1}}$$

$$n = \frac{\sigma^2 N}{\frac{\sigma^2}{2} (N-1) + \sigma^2}$$

$$n = \frac{1}{\frac{1}{N} + \frac{N-1}{N} \cdot \frac{1}{n_0}}$$

$$n_0 = \frac{2 \frac{\sigma^2}{2} \sigma^2}{\sigma^2}$$

↳ Tamaño de muestra cuando la población es infinita

Para una muestra para la media (M)

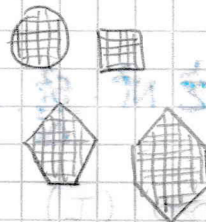
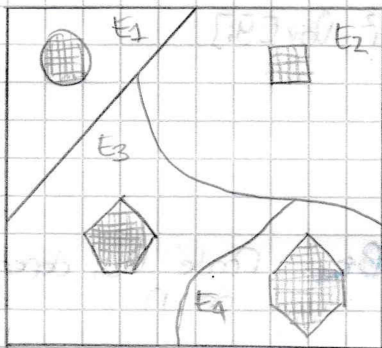
$$\sigma^2 = \hat{s}^2$$

Para una muestra para la Proporción (P)

$$\sigma^2 = \hat{p}(1-\hat{p})$$

### MUESTREO ALEATORIO ESTRATIFICADO (MAE)

Técnica de muestreo mediante la cual se obtiene una muestra a partir de la separación de los elementos de una población en grupos no traslapados, llamados estratos y la selección posterior de una muestra aleatoria simple, independiente, en cada estrato



MUESTRA ALEATORIA ESTRATIFICADA

- Maximiza la información obtenida
- Minimiza el límite del error de estimación (B)
- La variabilidad de la población es heterogénea
- La variabilidad dentro de los estratos es homogénea
- Disminuye costos
- Permite estimaciones para los subgrupos (estratos)

• Tamaño de la población (N)

$$N = \sum_{i=1}^L N_i$$

N = Tamaño de la población

N<sub>i</sub> = Unidades muestrales en el i-ésimo estrato

L = # de estratos

• Tomado de la muestra (n)

$$n = \sum_{i=1}^L n_i \quad n_i = \text{\# unidades muestrales seleccionadas del estrato } i$$

• Medidas de resumen por estrato

- $\bar{y}_i$  = Media muestral de la MAS seleccionada del estrato i
- $s_i^2$  = Varianza muestral de la MAS seleccionada del estrato i
- $p_i = n_i/n$  = \# unidades muestrales en la muestra que tienen el atributo

• Total de toda la población (T)

$$\hat{T}_{st} = \sum_{i=1}^L \hat{t}_i = \sum_{i=1}^L N_i \bar{y}_i$$

$$\hat{T}_{st} = \sum_{i=1}^L N_i \bar{y}_i$$

• Varianza estimada de T

$$\hat{Var}[\hat{T}_{st}] = \sum_{i=1}^L \hat{Var}[N_i \bar{y}_i] = \sum_{i=1}^L N_i^2 \cdot \hat{Var}[\bar{y}_i]$$

$$\hat{Var}[\hat{T}_{st}] = \sum_{i=1}^L N_i^2 \cdot \frac{s_i^2}{n_i} \cdot \left(\frac{N_i - n_i}{N_i}\right)$$

• Intervalo de confianza (I)  $\hat{T}_{st} \pm B_{T_{st}}$ , donde  $B_{T_{st}}$  depende de n

Si  $n > 30$

$$B_{T_{st}} = Z_{\alpha/2} \sqrt{\sum_{i=1}^L N_i^2 \cdot \frac{s_i^2}{n_i} \left(\frac{N_i - n_i}{N_i}\right)}$$

Entonces

$$\hat{T}_{st} \pm z_{\alpha/2} \sqrt{\sum_{i=1}^L N_i^2 \cdot \frac{s_i^2}{n_i} \left(\frac{N_i - n_i}{N_i}\right)}$$

Si  $n \leq 30$

$$B_{T_{st}} = t_{\alpha/2, n-L} \sqrt{\sum_{i=1}^L N_i^2 \cdot \frac{s_i^2}{n_i} \left(\frac{N_i - n_i}{N_i}\right)}$$

Entonces

$$\hat{T}_{st} \pm t_{\alpha/2, n-L} \sqrt{\sum_{i=1}^L N_i^2 \cdot \frac{s_i^2}{n_i} \left(\frac{N_i - n_i}{N_i}\right)}$$

• Media Muestral (M)

$$T = N \cdot M \Rightarrow M = \frac{T}{N} ; M = \bar{y}_{st} \Rightarrow \bar{y}_{st} = \frac{\hat{T}_{st}}{N}$$

$$\bar{y}_{st} = \frac{1}{N} \cdot \sum_{i=1}^L N_i \bar{y}_i$$

$\frac{N_i}{N}$  = Proporción de Unidades Muestrales en el estrato



- Varianza estimada (M)

$$\text{Var}[\bar{y}_{st}] = \text{Var}\left[\frac{1}{N} \cdot \sum N_i y_i\right] = \frac{1}{N^2} \cdot \text{Var}[\hat{t}_{st}]$$

$$\text{Var}[\bar{y}_{st}] = \sum \frac{N_i^2}{N^2} \cdot \frac{s_i^2}{n_i} \left(\frac{N_i - n_i}{N_i}\right)$$

- Intervalo de confianza (M)

$$\bar{y}_{st} \pm B_M ; \text{ Donde } B_M \text{ depende de } n \quad \frac{1}{N} (\hat{t}_{st} \pm B_{t_{st}})$$

$n > 30$

$$B_M = z_{\frac{\alpha}{2}} \sqrt{\sum \frac{N_i^2}{N^2} \cdot \frac{s_i^2}{n_i} \left(\frac{N_i - n_i}{N_i}\right)}$$

$n \leq 30$

$$B_M = t_{\frac{\alpha}{2}, n-L} \sqrt{\sum \frac{N_i^2}{N^2} \cdot \frac{s_i^2}{n_i} \left(\frac{N_i - n_i}{N_i}\right)}$$

$$\bar{y}_{st} \pm z_{\frac{\alpha}{2}} \sqrt{\sum \frac{N_i^2}{N^2} \cdot \frac{s_i^2}{n_i} \left(\frac{N_i - n_i}{N_i}\right)}$$

$$\bar{y}_{st} \pm t_{\frac{\alpha}{2}, n-L} \sqrt{\sum \frac{N_i^2}{N^2} \cdot \frac{s_i^2}{n_i} \left(\frac{N_i - n_i}{N_i}\right)}$$

- Total de la población (A)

$$\hat{A} = \sum \hat{A}_i = \sum N_i \hat{p}_i$$

$$\hat{A}_{st} = \sum N_i \hat{p}_i$$

- Varianza estimada (A)

$$\text{Var}[\hat{A}_{st}] = \text{Var}\left[\sum N_i \hat{p}_i\right] = \sum N_i^2 \cdot \text{Var}[\hat{p}_i]$$

$$\text{Var}[\hat{A}_{st}] = \sum N_i^2 \cdot \frac{\hat{p}_i(1-\hat{p}_i)}{n_i-1} \left(\frac{N_i - n_i}{N_i}\right)$$

- Intervalo de confianza (A)

$$\hat{A}_{st} \pm B_{\hat{A}_{st}} ; \text{ Donde } B_{\hat{A}_{st}} \text{ depende de } n$$

Si  $n > 30$

$$B_{\hat{A}_{st}} = z_{\frac{\alpha}{2}} \sqrt{\sum N_i^2 \cdot \frac{\hat{p}_i(1-\hat{p}_i)}{n_i-1} \left(\frac{N_i - n_i}{N_i}\right)}$$

Entonces

Si  $n \leq 30$

$$B_{\hat{A}_{st}} = t_{\frac{\alpha}{2}, n-L} \sqrt{\sum N_i^2 \cdot \frac{\hat{p}_i(1-\hat{p}_i)}{n_i-1} \left(\frac{N_i - n_i}{N_i}\right)}$$

Entonces

$$\hat{A}_{st} \pm z_{\frac{\alpha}{2}} \sqrt{\sum N_i^2 \frac{\hat{P}_i(1-\hat{P}_i)(N_i-n_i)}{n_i-1}} = \hat{A}_{st} \pm t_{\frac{\alpha}{2}, n-L} \sqrt{\sum N_i^2 \frac{\hat{P}_i(1-\hat{P}_i)(N_i-n_i)}{n_i-1}}$$

• Proporción poblacional (P)

$$\hat{A}_{st} = N \hat{P}_{st} \Rightarrow \hat{P}_{st} = \frac{\hat{A}_{st}}{N} \Rightarrow \hat{P}_{st} = \frac{\sum A_i}{N} = \frac{\sum N_i \hat{P}_i}{N}$$

$$\hat{P}_{st} = \frac{1}{N} \sum N_i \hat{P}_i$$

• Varianza estimada (P)

$$\begin{aligned} \text{Var}[\hat{P}_{st}] &= \text{Var}\left[\frac{1}{N} \sum N_i \hat{P}_i\right] = \sum \frac{N_i^2}{N^2} \text{Var}[\hat{P}_i] \\ &= \frac{1}{N^2} \cdot \text{Var}[\hat{A}_{st}] \end{aligned}$$

$$\text{Var}[\hat{P}_{st}] = \sum \frac{N_i^2}{N^2} \cdot \frac{\hat{P}_i(1-\hat{P}_i)}{n_i-1} \cdot \left(\frac{N_i-n_i}{N_i}\right)$$

• Intervalo de confianza

$$\hat{P} \pm B_{P_{st}}, \text{ Donde } B_{P_{st}} \text{ depende de } n \quad \frac{1}{N} (\hat{A}_{st} \pm B_{A_{st}})$$

Si  $n > 30$

$$B_{P_{st}} = z_{\frac{\alpha}{2}} \sqrt{\sum \frac{N_i^2}{N^2} \cdot \frac{\hat{P}_i(1-\hat{P}_i)}{n_i-1} \cdot \left(\frac{N_i-n_i}{N_i}\right)}$$

Entonces

Si  $n \leq 30$

$$B_{P_{st}} = t_{\frac{\alpha}{2}, n-L} \sqrt{\sum \frac{N_i^2}{N^2} \cdot \frac{\hat{P}_i(1-\hat{P}_i)}{n_i-1} \cdot \left(\frac{N_i-n_i}{N_i}\right)}$$

Entonces

$$\hat{P}_{st} \pm z_{\frac{\alpha}{2}} \sqrt{\sum \frac{N_i^2}{N^2} \cdot \frac{\hat{P}_i(1-\hat{P}_i)}{n_i-1} \cdot \left(\frac{N_i-n_i}{N_i}\right)}$$

$$\hat{P}_{st} \pm t_{\frac{\alpha}{2}, n-L} \sqrt{\sum \frac{N_i^2}{N^2} \cdot \frac{\hat{P}_i(1-\hat{P}_i)}{n_i-1} \cdot \left(\frac{N_i-n_i}{N_i}\right)}$$

Conocidos B y el nivel de confianza  $1-\alpha$  fijados por el investigador para  $\mu$ , se busca hallar cual debe ser el tamaño de muestra  $n$ , para ello partimos de la formula del limite del error de estimacion B

$$B = \frac{z_{\alpha/2}}{2} \sqrt{\text{Var}(\bar{y}_{sc})}; \text{Entonces } V(\bar{y}_{sc}) = \frac{B^2}{\frac{z_{\alpha/2}^2}{2}}$$

La varianza disminuye conforme aumenta  $n$

$$\frac{B^2}{\frac{z_{\alpha/2}^2}{2}} = \sum \left( \frac{N_i^2}{N} \right) \frac{\sigma_i^2}{n_i} \frac{N_i - n_i}{N_i}; \quad n_i = w_i \cdot n$$

$$n = \frac{\sum N_i^2 \sigma_i^2 / w_i}{D N^2 + \sum N_i \cdot \sigma_i^2}, \quad D = \frac{B^2}{\frac{z_{\alpha/2}^2}{2}}$$

Donde  $w_i$  depende de la afijacion

Para $\mu$	Para $P$	Para totales (T, A) lo unico que cambia con respecto a (M, P) es D
$\sigma^2 = S^2$	$\sigma^2 = P(1-P)$	
$\sigma_i = \frac{R_i}{6} = \frac{y_{i,max} - y_{i,min}}{6}$		$D = \frac{B^2}{\frac{z_{\alpha/2}^2}{2} \cdot N^2}$

• Asignacion Optima con costos variables

$C_i =$  diferentes  
 $\sigma_i =$  diferentes  
 $w_i \propto \sigma_i$   
 $w_i \propto 1/\sqrt{C_i}$   
 $w_i \propto N_i$

$$C = C_0 + \sum C_i n_i$$

$$w_i = \frac{N_i \cdot \sigma_i / \sqrt{C_i}}{\sum N_k \cdot \sigma_k / \sqrt{C_k}} \quad \text{Estandarizar}$$

• Asignacion de Neyman

$$C_1 = C_2 = C_3 = \dots = C_n =$$

$\sigma_i =$  diferentes

$$w_i = \frac{N_i \cdot \sigma_i}{\sum N_k \cdot \sigma_k}$$

• Asignacion Proporcional

$$C_1 = C_2 = C_3 = \dots = C_n$$

$\sigma_i =$  Iguales

$$w_i = \frac{N_i}{N}$$

# REGRESIÓN LINEAL

↳ Técnica para investigar y modelar la relación existente entre variables.

Y = Variable respuesta → Dependiente

X = Variable regresora → Independiente

## Origen

↳ Sir Francis Galton (1869)

\* El análisis de regresión permite estudiar las relaciones de asociación entre  $y$  y  $x$

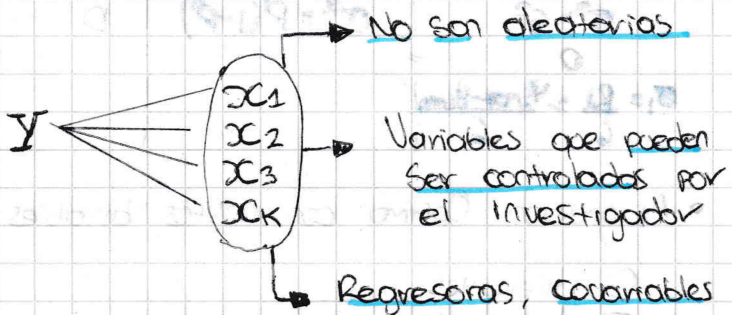
Relación estatura de niños con estatura de padres

- ESTA ASOCIACIÓN NO ES RELACION CAUSA EFECTO

↳ NO IMPLICA CAUSALIDAD

- ↳ A padres Altos, los hijos generalmente lo son también.
- A padres bajos, los hijos son de menor estatura
- A padres muy altos o muy bajos se percibe una regresión hacia la estatura media de la población

- La relación Causa-Efecto, solo aplica para diseño de experimentos



## OTROS MODELOS DE REGRESION:

- Lineales generalizados
- Logística (Sigmoides)
- Prueba de Signos
- No paramétricas

## • MODELO DE REGRESION LINEAL SIMPLE

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad , i=1, 2, \dots, n$$

Donde

$n_i$  = # de datos en el que se supone que existe una relación de tipo lineal

$\beta_1$ : Pendiente de la recta de regresión (Cambio en la respuesta media por un cambio unitario en  $x$ )

$y_i$  = Observación

$\epsilon_i$  = Error aleatorio del modelo

$\beta_0$  = Intercepto de la recta de regresión  
Se interpreta si  $0 \in [x_{min}, x_{max}]$

$\beta_0, \beta_1$  = Parámetros

$\epsilon_i$  = Variable aleatoria

Para estimar los parámetros  $(\beta_1, \beta_0)$  del modelo utilizamos dos métodos: (1) Máxima Verosimilitud (Maximum likelihood) y (2) Mínimos Cuadrados

1) Máxima Verosimilitud = Maximizar la probabilidad de ocurrencia

↳ Requiere la verificación sobre los residuos del modelo (Supuestos): Resid (Modelo)

- Distribuyen Normal (Shapiro-Wilk)
- Son independientes (Durbin-Watson)
- Tienen Varianza Constante (Prueba F)

$$E[\epsilon_i] = 0$$

$$\text{Var}[\epsilon_i] = \sigma^2$$

$\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$  → Los residuos distribuyen normal  
(Independientes e Idénticamente distribuidos)

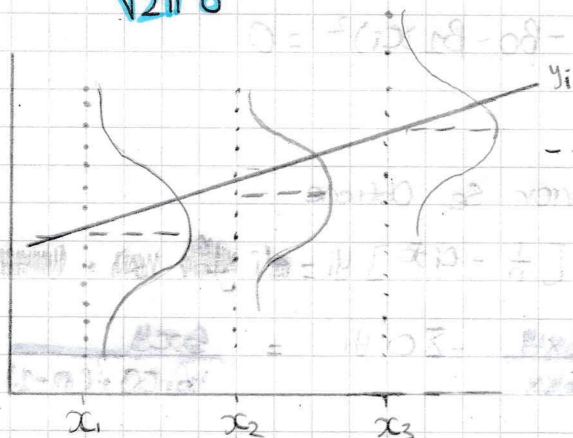
↳ Consecuencias de los supuestos

$y_i \sim NI(\beta_0 + \beta_1 x_i, \sigma^2)$  → Distribución Normal (No idénticas)  
Independientes

Fdp

$$f(y_i) = \frac{e^{-\frac{1}{2\sigma^2}(y_i - \beta_0 - \beta_1 x_i)^2}}{\sqrt{2\pi}\sigma}$$

↳  $E[y_i] = E[y_i | x_i] = \beta_0 + \beta_1 x_i$   
 $\text{Var}[y_i] = \sigma^2$



- OTRAS PRUEBAS DE NORMALIDAD
- Kolmogorov-Smirnov
  - Cramer-von Mises
  - Anderson-Darling

- FUNCION DE VEROSIMILITUD

$L(\beta_0, \beta_1, \sigma^2)$  = Verosimilitud

$$L = \prod_{i=1}^n f(y_i) = \prod_{i=1}^n \frac{e^{-\frac{1}{2\sigma^2}(y_i - \beta_0 - \beta_1 x_i)^2}}{\sqrt{2\pi}\sigma}$$

$l(\beta_0, \beta_1, \sigma^2) = \text{Log Verosimilitud}$

$$l = \text{Ln } L(\beta_0, \beta_1, \sigma^2) = \text{Ln} \left[ (2\pi\sigma^2)^{-n/2} * e^{-\frac{1}{2\sigma^2} \sum (y_i - \beta_0 - \beta_1 x_i)^2} \right]$$

$$l = -\frac{n}{2} \text{Ln}(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

$$l = -\frac{n}{2} \text{Ln}(2\pi) - \frac{n}{2} \text{Ln}(\sigma^2) - \frac{1}{2\sigma^2} \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

Para hallar los valores que maximizan la función derivamos e igualamos a cero.

$$\frac{\partial l}{\partial \beta_0} = \frac{1}{2\sigma^2} \sum 2(y_i - \beta_0 - \beta_1 x_i)(-1) = 0$$

$$\frac{\partial l}{\partial \beta_1} = \frac{1}{2\sigma^2} \sum 2(y_i - \beta_0 - \beta_1 x_i)(-x_i) = 0$$

$$\frac{\partial l}{\partial \sigma^2} = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{(\sigma^2)^{-2} \sum (y_i - \beta_0 - \beta_1 x_i)^2}{2} = 0$$

Resolviendo el sistema anterior se obtiene

$$\tilde{\beta}_0 = \bar{y} - \tilde{\beta}_1 \bar{x} = \sum \left[ \frac{1}{n} - c_i \bar{x} \right] y_i = d_i y_i$$

$$\tilde{\beta}_1 = \frac{(\sum x_i y_i) - (n \bar{x} \bar{y})}{\sum x_i^2 - n \bar{x}^2} = \frac{S_{xy}}{S_{xx}} = \sum c_i y_i = \frac{S_{xy}}{\text{Var}(x) \cdot (n-1)}$$

$$\tilde{\sigma}^2 = \frac{\sum (y_i - \tilde{\beta}_0 - \tilde{\beta}_1 x_i)^2}{n}$$

Donde  $c_i = \frac{x_i - \bar{x}}{S_{xx}}$  ;  $d_i = \frac{1}{n} - c_i \bar{x}$

Los estimadores  $\tilde{\beta}_0$  y  $\tilde{\beta}_1$  son insesgados, esto se verifica conforme a los siguientes casos:

$$a) \sum c_i = \sum \frac{(x_i - \bar{x})}{S_{xx}} = \frac{1}{S_{xx}} \cdot \sum (x_i - \bar{x}) \rightarrow 0$$

$$b) \sum c_i x_i = \sum \frac{x_i - \bar{x}}{S_{xx}} (x_i) = \frac{1}{S_{xx}} \cdot \sum (x_i - \bar{x}) x_i \rightarrow S_{xx}$$

$$\boxed{\sum c_i = 0}$$

$$\boxed{\sum c_i x_i = 1}$$

1. Para  $\tilde{\beta}_1$

$$E[\tilde{\beta}_1] = E\left[\sum c_i y_i\right] = \sum E[c_i y_i] = \sum c_i E[y_i]$$

$$= \sum c_i (\beta_0 + \beta_1 x_i) = \beta_0 \cdot \sum c_i + \beta_1 \cdot \sum c_i x_i \rightarrow 0$$

$$= \beta_1 \quad \checkmark$$

$$\text{Var}(\tilde{\beta}_1) = V\left[\sum c_i y_i\right] \rightarrow y_i \text{ independientes}$$

$$= \sum c_i^2 \cdot \text{Var}[y_i] = \sum \frac{(x_i - \bar{x})^2}{S_{xx}} \cdot \sigma^2$$

$$= \frac{\sum (x_i - \bar{x})^2}{(S_{xx})^2} \cdot \sigma^2 = \frac{\sigma^2}{S_{xx}}$$

2. Para  $\tilde{\beta}_0$

$$E[\tilde{\beta}_0] = E\left[\bar{y} - \tilde{\beta}_1 \bar{x}\right] = E\left[\frac{1}{n} \sum y_i - \hat{\beta}_1 \bar{x}\right] = \frac{1}{n} \sum E[y_i] - \bar{x} E[\hat{\beta}_1]$$

$$= \frac{1}{n} \sum (\beta_0 + \beta_1 x_i) - \bar{x} \beta_1 = \beta_0 + \beta_1 \bar{x} - \bar{x} \beta_1 = \beta_0 \quad \checkmark$$

$$\text{Var}(\tilde{\beta}_0) = V\left[\sum d_i y_i\right] = \sum d_i^2 \cdot V[y_i] = \sum d_i^2 \sigma^2 = \sigma^2 \sum d_i^2$$

$$= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)$$

El estimador  $\tilde{\sigma}^2$  es sesgado para  $\sigma^2$

$E[\hat{\sigma}^2] = \left(\frac{n-2}{n}\right) \sigma^2$  Pero es posible construir un estimador insesgado de  $\sigma^2$

$$\hat{\sigma}^2 = \left(\frac{n}{n-2}\right) \tilde{\sigma}^2 = \frac{\sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}{n-2} = \frac{\sum e_i^2}{n-2} = \frac{\sum (y_i - \hat{y}_i)^2}{n-2}$$

$\hat{\sigma}^2 = \text{MSE} = \text{Error Cuadrático Medio}$

$e_i = \text{Residuos}$

$$\text{MSE} = \frac{1}{n-2} (S_{yy} - \hat{\beta}_1^2 * S_{xx})$$

Dados los estimadores concluimos además

$$\hat{y}_i = E[y_i | x_i] = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

↳ que  $\hat{y}_i$  es un estimador insesgado para  $y_i$ , esto debido al principio de Invariancia

$$\begin{aligned} E[\hat{y}_i] &= E[\hat{\beta}_0 + \hat{\beta}_1 x_i] \\ &= E[\hat{\beta}_0] + E[\hat{\beta}_1] x_i \\ &= \beta_0 + \beta_1 x_i \end{aligned}$$

En consecuencia

$$\hat{\beta}_1 \sim N(\beta_1, \sigma^2 / S_{xx})$$

$$\hat{\beta}_0 \sim N(\beta_0, \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right])$$

2) Minimos Cuadrados = Minimizan la suma de los cuadrados de los errores

↳ No requiere supuestos

$$\begin{aligned} \text{Funcion MC: } S(\beta_0, \beta_1) &= \sum e_i^2 = \sum (y_i - \hat{y}_i)^2 \\ &= \sum [y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i]^2 \end{aligned}$$



Los estimadores de mínimos cuadrados los obtendremos de la siguiente manera

$$\begin{aligned} \frac{\partial S}{\partial \hat{\beta}_1} &= \sum [y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i]^2 \\ &= \sum 2[y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i](-x_i) \\ &= \sum 2[x_i y_i - \hat{\beta}_0 x_i - \hat{\beta}_1 x_i^2] \end{aligned}$$

$$\begin{aligned} \bar{y} &= \frac{\sum y_i}{n} \\ \bar{y}n &= \sum y_i \end{aligned}$$

Iguales a cero

$$\begin{aligned} 0 &= \sum x_i y_i - \hat{\beta}_0 \sum x_i - \hat{\beta}_1 \sum x_i^2 \\ \sum x_i y_i &= \hat{\beta}_0 \sum x_i + \hat{\beta}_1 \sum x_i^2 \end{aligned} \quad (1)$$

$$\frac{\partial S}{\partial \hat{\beta}_0} = \sum 2[y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i](-1)$$

Iguales a cero

$$\sum y_i - \sum \hat{\beta}_0 - \sum \hat{\beta}_1 x_i = 0$$

$$n\bar{y} - n\hat{\beta}_0 - n\bar{x}\hat{\beta}_1 = 0$$

$$\bar{y} - \hat{\beta}_0 - \bar{x}\hat{\beta}_1 = 0$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (2)$$

Reemplazamos (2) en (1)

$$\sum x_i y_i = \bar{y}n\bar{x} - n\hat{\beta}_1 \bar{x}^2 + \hat{\beta}_1 \sum x_i^2$$

$$\sum x_i y_i - n\bar{x}\bar{y} = \hat{\beta}_1 [-n\bar{x}^2 + \sum x_i^2]$$

$$\hat{\beta}_1 = \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2}$$

Tras realizar el procedimiento notamos que los estimadores de Máxima Verosimilitud y de mínimos cuadrados son iguales, sin embargo los EMV facilitan las pruebas de significancia y el cálculo de intervalos de confianza y de predicción

## SIGNIFICANCIA DE LOS ESTIMADORES ( $\beta_0, \beta_1$ )

- Para  $\beta_1$  (pendiente): Por un aumento de 1 unidad en -- se estima que hay un incremento promedio de ( $\beta_1$ ) en --

$$\hat{\beta}_1 \sim N(\beta_1, \sigma^2 / S_{xx})$$

$\beta_{10}$  = valor fijo

↳ generalmente cero

↳ sig. regresión

Prueba de hipótesis

$$H_0: \beta_1 = \beta_{10}$$

$$H_a: \beta_1 \neq \beta_{10}$$

Estadístico de Prueba

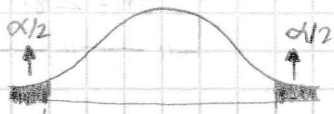
$$t_{\hat{\beta}_1} = \frac{\hat{\beta}_1 - \beta_{10}}{ee(\hat{\beta}_1)} \sim t_{n-2}$$

$$\text{Donde } ee(\hat{\beta}_1) = \sqrt{\hat{Var}(\hat{\beta}_1)}$$

$$Var(\hat{\beta}_1) = \frac{MSE}{S_{xx}}$$

Criterio de decision = Significancia de la regresion  $[(t_{01})^2 \geq F_0]$

- $RR_1 = \{ |t_{01}| > t_{\frac{\alpha}{2}, n-2} \}$  // Rechazamos si  $t_{01}$  esta en la RR
- $P\text{-value} = P [ |t_{01}| > t_{\frac{\alpha}{2}, n-2} ]$  //  $P\text{-Value} < \alpha \rightarrow$  Rechazamos  $H_0$



Si Rechazo  $H_0$  la regresion es significativa de lo contrario las variables son independientes ( $\beta_1=0$ )

- Para  $\beta_0$  (Intercepto) = Es interpretable si  $0 \in [X_{min}, X_{max}]$

$$\beta_0 \sim N \left( \beta_0, \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum x^2} \right] \right)$$

$\beta_{00}$  = Valor fido

Generalmente cero

Prueba de hipotesis

$$H_0: \beta_0 = \beta_{00}$$

$$H_1: \beta_0 \neq \beta_{00}$$

Estadistico de prueba

$$t_{00} = \frac{\beta_0 - \beta_{00}}{ee(\hat{\beta}_0)} \sim t_{n-2}$$

$$\text{Var}[\hat{\beta}_0] = \text{MSE} \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum x^2} \right]$$

Donde  $ee(\hat{\beta}_0) = \sqrt{\text{var}[\hat{\beta}_0]}$

Criterio de decision

- $RR_0 = \{ |t_{00}| > t_{\frac{\alpha}{2}, n-2} \}$  // Rechazamos si  $t_{00}$  esta en la RR
- $P\text{-value} = P [ |t_{00}| > t_{\frac{\alpha}{2}, n-2} ]$  //  $P\text{-Value} < \alpha \rightarrow$  Rechazamos  $H_0$

### INTERVALOS DE CONFIANZA DE LOS ESTIMADORES ( $\beta_0, \beta_1$ )

- Para  $\beta_1$

$$\beta_1 \pm t_{\frac{\alpha}{2}, n-2} * ee(\hat{\beta}_1)$$

- Si el IC no contiene al cero se dice que hay una relacion lineal en la cual la variable  $x$  es adecuada para predecir el comportamiento de  $y$

- Para  $\beta_0$

$$\beta_0 \pm t_{\frac{\alpha}{2}, n-2} * ee(\hat{\beta}_0)$$

- Si  $x=0$  tiene sentido y el IC no contiene al cero entonces  $E[y|x=0] = \beta_0$
- Si contiene al cero y  $x=0$  tiene sentido entonces se recomienda abstenerse en modelo de la forma  $y = \beta_1 x + \epsilon$

## ANÁLISIS DE VARIANZA (ANOVA)

Recordemos que evaluamos la Significancia de la regresión con una prueba de hipótesis para  $\beta_1 = 0$ , retomaremos este concepto para el análisis de Varianza, y evaluaremos de forma similar la Significancia de la regresión.

- Se analiza la variabilidad de  $(y)$ , analizando las desviaciones de cada observación  $y_i$  respecto de su medio  $\bar{y}$

$$(y_i - \bar{y})$$

La medida de la variación total en  $y$  es la suma de los cuadrados de las desviaciones de los  $y_i$ 's con respecto a su media Muestral

$$S_{yy} = SST = \sum (y_i - \bar{y})^2 = \text{Var}[y] * (n-1) = \sum y_i^2 - n\bar{y}^2$$

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2 \rightarrow \text{Identidad suma de cuadrados}$$

$\downarrow$                        $\downarrow$                        $\downarrow$   
SST                      SSR                      SSE  
(n-1)                      (1)                      (n-2)

Donde

SST = Suma de cuadrados total

SSR = Suma de Cuadrados de la regresión

SSE = Suma de cuadrados del error

$$SSR = \hat{\beta}_1 S_{xy} = \hat{\beta}_1^2 S_{xx} \rightarrow \text{Cantidad de variación que alcanza a explicar con el modelo RLS}$$

$$SSE = SST - SSR = S_{yy} - \hat{\beta}_1^2 S_{xx}$$

Cada una de las sumas tiene asociados grados de libertad y con ellos construiremos estimaciones independientes de  $\sigma^2$ , de la siguiente manera:

$$MSR = \frac{SSR}{1} ; \quad MSE = \frac{SSE}{n-2}$$

Para realizar la Prueba de Significancia de la regresión partiendo de el anova utilizamos  $F_{\alpha}$  para  $\beta_1$ ; utilizando  $F_{\alpha}$  como estadístico de Prub.

$$\begin{aligned} H_0 &: \beta_1 = 0 \\ H_a &: \beta_1 \neq 0 \end{aligned}$$

$$F_0 = \frac{MSR}{MSE} \sim F_{1, n-2}$$

$$RR = Y \quad F_{\alpha} > F_{1-\alpha, 1, n-2}$$

$$P\text{-value} = P[F_{1, n-2} > F_{\alpha}]$$

• Tabla ANOVA



Fuente de Variación	SS	gl	MS	F
Regresión o Modelo	SSR	1	MSR = SSR/1	F <sub>01</sub> = MSR / MSE
Error o Residual	SSE	n-2	MSE = SSE/n-2	$\sim F_{1, n-2}$
Total	SST	n-1		

Una medida de bondad de ajuste es el coeficiente de determinación denominado  $R^2$ , interpretado como la proporción de variación en la respuesta que alcanza a ser explicada por el modelo de regresión lineal

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} = r^2 \quad \text{Donde}$$

$r$ : Coeficiente de correlación de Pearson

Consideraciones

- Pendiente Cercana a Cero
- Si  $R^2$  tiende a cero la relación entre  $x$  y  $y$  es muy pobre
- Si  $R^2$  tiende a uno la recta ajustada se aproxima bien a los puntos
- Pendiente Significativa
- $R^2$  grande no implica pendiente ( $\beta_1$ ) Grande
- $R^2$  grande no garantiza que el modelo de RLS ajustado sea el adecuado para los datos

RESPUESTA MEDIA

Se permite <sup>solo</sup> interpolación es decir buscar un  $x_0 \in [x_{\min}, x_{\max}]$  y este dato por la expresión  $\hat{y}_0 = E[y|x=x_0] = \beta_0 + \beta_1 x_0$  para una estimación puntual

Para la estimación por intervalo tenemos  $\beta_0$  y  $\beta_1$  son LD

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-2} ee(\hat{y}_0)$$

Donde

$$ee(\hat{y}_0) = \sqrt{\hat{Var}(\hat{y}_0)}$$

$$\hat{Var}(\hat{y}_0) = MSE \left( \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)$$

$$\begin{aligned} \hat{Var}(\hat{y}_0) &= \hat{Var}[\hat{\beta}_0 + \hat{\beta}_1 x_0] \\ \hat{Var}(\hat{y}_0) &= \hat{Var}\left[\sum d_i y_i + \sum c_i y_i x_0\right] \\ \hat{Var}(\hat{y}_0) &= \hat{Var}\left[\sum (d_i + c_i x_0) y_i\right] \\ &= \sum (d_i + c_i x_0)^2 * \hat{Var}(y_i) \\ &= \sum \left[ \frac{1}{n} - \bar{x} c_i \right] + \frac{(x_i - \bar{x}) x_0}{S_{xx}} \sigma^2 \\ &= MSE \left[ \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right] \end{aligned}$$

El intervalo de longitud  $2 ee(\hat{y}_0)$  muestra se da cuando  $x_0 = \bar{x}$

## PREDICCIÓN

↳ Interesa una nueva observación  $y_0$  que corresponda a un nivel específico de  $x$ , denominado  $x_0$ .

El estimador puntual está dado al igual que el anterior por la expresión  $\hat{y}_0 = E[y|x=x_0] = \beta_0 + \beta_1 x_0$

↳ solo se permite interpolar

Para la estimación del intervalo tenemos el error en la predicción estará dado por  $(y_0 - \hat{y}_0)$  por lo que el intervalo quedará:

$$\hat{y}_0 \pm t_{\alpha/2, n-2} * ee(y_0 - \hat{y}_0)$$

Donde

$$ee(y_0 - \hat{y}_0) = \sqrt{\hat{var}(y_0 - \hat{y}_0)}$$

$$\hat{var}(y_0 - \hat{y}_0) = MSE \left[ 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]$$

$$\begin{aligned} \hat{var}(y_0 - \hat{y}_0) &= \hat{var}(y_0) + \hat{var}(\hat{y}_0) \\ &= \sigma^2 + \sigma^2 \left[ \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right] \\ &= \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right] \\ &= MSE \left[ 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right] \end{aligned}$$

## VALIDACION DE LOS SUPUESTOS

↳ Los supuestos son los ya nombrados y los hacemos sobre los residuales → Se hace antes se ajusta el modelo

- Distribución normal
- Media cero
- Varianza Constante
- Inconrelacionados

Ayuda a detectar si

- La función de regresión no es lineal
- Faltan variables predictoras al modelo

• Normalidad

1. Shapiro-Wilk :  $W = \frac{(\sum a_i e_i)^2}{\sum (e_i - \bar{e})^2}$

Donde se concluye con el Criterio del P-value

H<sub>0</sub>: Distribuye Normal

H<sub>a</sub>: NO Distribuye Normal

2. Q-Q Plot

3. P-P pbt

4. Histograma

• Media cero e Inconrelacion

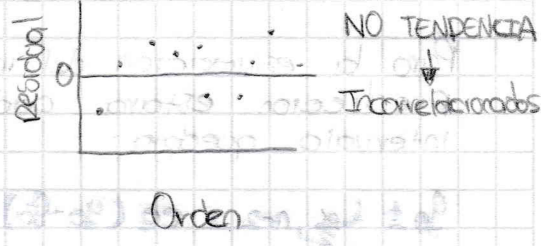
↳ Usualmente no se valida por la forma en que se estiman los parámetros

Si se tiene el orden de registro ( $t_i$ ) se grafica  $e_i$  vs  $t_i$ . SI NO HAY TENDENCIA SE CONCLUYE QUE NO HAY CORRELACION ENTRE LOS ERRORES DEL MODELO

$$\bar{e} = \frac{1}{n} \sum e_i = 0, \text{ lo cual se cumple si } \sum e_i = 0$$

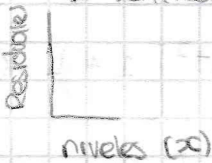
Lo que a su vez implica

$$\sum y_i = \sum \hat{y}_i$$

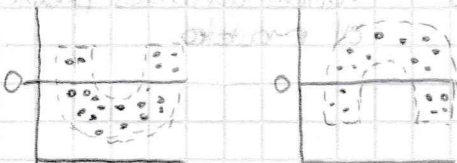


• Varianza Constante

↳ Verificamos graficamente = NO TENDENCIA



FALTA VARIABLE CUADRATICA



↳ NO podemos decir si es Constante o no

Transformacion Box-cox Coords no es normal

PRUEBA DE CARENANCIA DE AJUSTE

↳ Violacion del supuesto de linealidad

Consideraciones:

- Se realiza para la validacion del supuesto de linealidad
- Se hace despues de validar Normalidad, Independencia y  $\sigma^2$  Constante
- Solo se puede realizar cuando se tiene al menos en un nivel de  $x$  dos o mas valores de  $y$  (Replicas)

$$H_0: E[Y|X=x] = \beta_0 + \beta_1 x$$

$$H_a: E[Y|X=x] \neq \beta_0 + \beta_1 x$$

$H_0$ : El modelo de primer orden es apropiado

$H_a$ : El modelo de primer orden no es apropiado

En esta prueba las replicas son utilizadas para estimar  $\sigma^2$  independiente del modelo abstracto

$X$	$Y$			$n_i$	
$x_1$	$y_{11}$	$y_{12}$	...	$y_{1n_1}$	$n_1$
$x_2$	$y_{21}$	$y_{22}$	...	$y_{2n_2}$	$n_2$
$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$
$x_m$	$y_{m1}$	$y_{m2}$		$y_{mn_m}$	$n_m$

$n=1$	$n=2$	$n=m$
$\vdots$	$\vdots$	$y_{m1}$
$\vdots$	$y_{12}$	$y_{m2}$
$\vdots$	$y_{21}$	$y_{m1}$
$x_1$	$x_2$	$x_m$

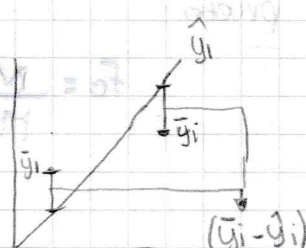
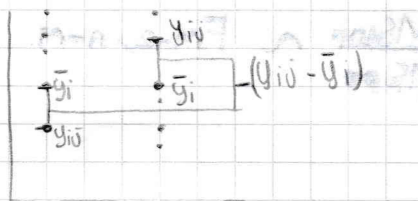
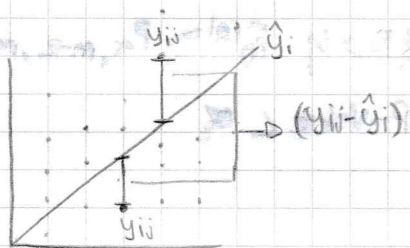
$m$  = Distintos valores de la variable  $X$   $\sum n_i = \#$  total de observaciones

$n_i$  = Observaciones de la variable  $Y$   $y_{ij}$  =  $j$ -ésima observación de la variable asociada al nivel  $i$ -ésimo de la variable regresora  $X$ .

En la prueba de falta de ajuste descomponemos la Suma de Cuadrados del error (SSE)

$$SSE = \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \hat{y}_i)^2 = \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 + \sum_{i=1}^m \sum_{j=1}^{n_i} (\bar{y}_i - \hat{y}_i)^2$$

$\downarrow$   $\downarrow$   $\downarrow$   
 SSE  $SS_{pe}$   $SS_{Lof}$   
 $(n-2)$   $(n-m)$   $(m-2)$



$SS_{pe}$  = Suma de cuadrados asociada al error puro (dentro)

$SS_{Lof}$  = Suma de cuadrados asociada a la falta de ajuste (Aduste)

$MS_{pe}$  = Cuadrados Medios debido al error puro

$MS_{Lof}$  = Cuadrados Medios debido a la falta de ajuste

$$MS_{SPE} = \frac{SS_{SPE}}{n-M}$$

$$E[MS_{SPE}] = \sigma^2$$

$$MS_{SLOF} = \frac{SS_{SLOF}}{M-2}$$

$$E[MS_{SLOF}] = \sigma^2 + \frac{\sum_{i=1}^M (E(y_i) - \beta_0 - \beta_1 x_i)^2}{M-2}$$

Si  $V(\epsilon_i) = \sigma^2$  entonces  $SS_{SPE}$  es un estimador de  $\sigma^2$  independiente del modelo, solo se usa la variabilidad de los  $y_i$ 's en cada nivel de  $x$

Observemos que si la función verdadera es lineal, entonces se cumple que

$$E[y_i] = E[y_i | x = x_i] = \beta_0 + \beta_1 x_i$$

Y por tanto se obtiene que la esperanza del  $MS_{SLOF}$  es:

$$E[MS_{SLOF}] = \sigma^2 + \sum_{i=1}^M [(\beta_0 + \beta_1 x_i) - (\beta_0 + \beta_1 x_i)]^2$$

$$E[MS_{SLOF}] = \sigma^2$$

Si la función verdadera no es lineal, entonces  $E[MS_{SLOF}] > \sigma^2$ .

Para la Ph. de falta de ajuste tenemos el estadístico de prueba:

$$F_0 = \frac{MS_{SLOF}}{MS_{SPE}} \sim F_{M-2, n-M}$$

$$RR = Y \quad F_0^{cal} > F_{\alpha, m-2, m-n}$$

P-Value  $< \alpha$

### CONSIDERACIONES

1. El  $SS_{SPE}$  refleja la variación aleatoria o error experimental puro
2. El  $SS_{SLOF}$  es una medida de la variación sistemática introducida por términos de distorsión lineal o de primer orden
3. Si el modelo de RLS no se ajusta a los datos de manera apropiada, entonces la  $SSE$  estará inflada y producirá un estimador sesgado de  $\sigma^2$  [ $MS_{SLOF} = \hat{\sigma}^2$ ]



TRANSFORMACIONES :  $y = \beta_0 + \beta_1 x + \epsilon$

Usamos transformaciones cuando hay falta de ajuste o cuando no se cumplen los supuestos de normalidad o homogeneidad de varianzas:

- Falta de ajuste
  - No se cumple el supuesto de normalidad
  - No se cumple el supuesto de homocedasticidad
- TRASFORMACION

### Notas

- Tenga en cuenta siempre la transformación a la que se sometieron los datos
- Al hacer inferencia se deben convertir los resultados a la escala original e interpretarlos en el contexto del problema

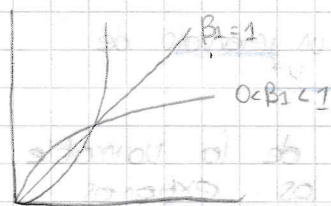
Las transformaciones los podemos aplicar en  $x$ , o en  $y$ , o en  $x$  y  $y$ . Al realizar la transformación tenemos un modelo de la forma

$$y^* = \beta_0^* + \beta_1^* x^* + \epsilon^*$$

Donde \* indica las variables transformadas

### 1. MODELO POTENCIA

$y = \beta_0 \times \beta_1^x$  Aplicamos log de tal manera que obtenemos:



$$\text{Log}(y) = \text{Log}(\beta_0) + \beta_2 \text{Log}(x)$$

$$y^* = \text{Log}(\beta_0) + \beta_2 x^*$$

### 2. MODELO EXPONENCIAL



$$y = \beta_0 e^{\beta_2 x}$$

Aplicamos Log de tal manera que obtenemos:

$$\text{Log}(y) = \text{Log} \beta_0 + \beta_2 x$$

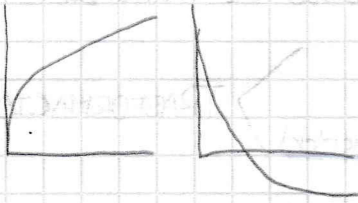
$$y^* = \text{Log} \beta_0 + \beta_2 x$$

### 3. MODELO LOG EN X

$$y = \beta_0 + \beta_1 \text{Log}(x)$$

Aplicamos log sobre la variable x:

$$y = \beta_0 + \beta_1 x^*$$



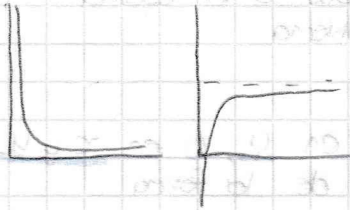
### 4. MODELO RECÍPROCO

$$y = \frac{x}{\beta_0 x - \beta_1}$$

Aplicamos el inverso de tal manera que obtenemos

$$\frac{1}{y} = \beta_0 - \beta_1 \frac{1}{x}$$

$$y^* = \beta_0 - \beta_1 x^*$$



### CONSIDERACIONES DE LAS TRANSFORMACIONES

La transformación de forma general se expresa:

$$y^* = g(y) \quad \text{con } g \text{ invertible}$$

$(u, v)$  es una expresión para un intervalo de la variable transformada  $y^*$

Es decir, que para concluir en términos de la variable Original debemos invertir cada uno de los extremos

- Para intervalos de predicción el valor de la respuesta futura  $y$  en el valor de  $x = x_0$  es de la forma:

$$(g^{-1}(u), g^{-1}(v)) \quad \text{con un nivel de confianza } 1 - \alpha$$

- Para intervalos de confianza el valor de la respuesta media  $E[y|x_0]$  es de la forma

$$(K \cdot g^{-1}(u), K \cdot g^{-1}(v)) \quad \text{Donde } K \text{ es el factor de corrección que depende de la transformación aplicada}$$

En el caso de  $y^* = \text{Log}(y)$  el factor de corrección  $K$  es  $[e^{MSE/2}]$

# • MODELO DE REGRESION LINEAL MULTIPLE

Estos modelos con frecuencia se escriben de forma Matricial:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}_{n \times p} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}_{p \times 1} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}_{n \times 1}$$

El MODELO DESCRIBE UN HIPERPLANO CON COEFICIENTES DE REGRESION  $\beta_0, \beta_1, \dots, \beta_p$  Y ES LINEAL POR SER UNA FUNCION LINEAL DE  $\beta$

$$\overset{VA}{y} = \overset{FID}{X} \overset{VA}{\beta} + \overset{VA}{\epsilon}$$

$n \times 1$       $n \times p$       $p \times 1$       $n \times 1$

Donde

- $y$  = Vector Aleatorio (respuesta)
- $\beta$  = Vector de parametros
- $\epsilon$  = Vector error aleatorio

$$\epsilon \sim N_n(0, \sigma^2 I_n)$$

Los errores distribuyen Normal n-varrida, con vector de medias cero y varianza constante

$$E[\epsilon] = \begin{bmatrix} E[\epsilon_1] \\ \vdots \\ E[\epsilon_n] \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$COV[\epsilon] = \begin{bmatrix} Var[\epsilon_1] & Cov[\epsilon_1, \epsilon_2] & \dots \\ \vdots & Var[\epsilon_2] & \dots \\ \vdots & \vdots & \ddots \\ Cov[\epsilon_1, \epsilon_2] & \dots & Var[\epsilon_n] \end{bmatrix}$$

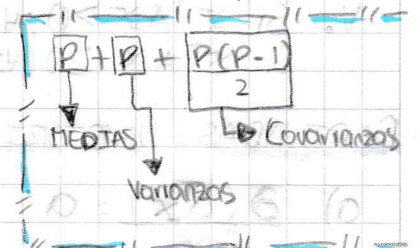
COV[X<sub>1</sub>, X<sub>2</sub>] = 0  
 ↳ NO implica independencia si correlacion  
 Independencia  
 ↳ implica COV[X<sub>1</sub>, X<sub>2</sub>] = 0

$$COV[\epsilon] = \begin{bmatrix} \sigma^2 & 0 & \dots \\ 0 & \sigma^2 & \dots \\ \vdots & \vdots & \ddots \\ 0 & \dots & \sigma^2 \end{bmatrix} = \sigma^2 \begin{bmatrix} 1 & 0 & \dots \\ 0 & 1 & \dots \\ \vdots & \vdots & \ddots \\ 0 & \dots & 1 \end{bmatrix}$$

$$COV[\epsilon] = \sigma^2 I_n$$

La esperanza y la varianza del modelo estan dadas por:

$$E[y] = H y = \begin{bmatrix} E[y_1] \\ \vdots \\ E[y_n] \end{bmatrix} = \text{Vector de Medias}$$



$$Var[y] = \sum y = COV[y] = \begin{bmatrix} Var[y_1] & Cov[y_1, y_2] & \dots & Cov[y_1, y_n] \\ Cov[y_2, y_1] & Var[y_2] & \dots & Cov[y_2, y_n] \\ \vdots & \vdots & \ddots & \vdots \\ Cov[y_n, y_1] & \dots & \dots & Var[y_n] \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \dots & \dots & \sigma_n^2 \end{bmatrix}$$

Matriz Varianzas - Covarianzas

Matriz Simetrica y Cuadrada

La Matriz de Varianzas-Covarianzas asociada a un vector aleatorio es una matriz cuadrada y simétrica de orden igual al tamaño del vector aleatorio donde la diagonal principal contiene las varianzas asociadas a cada elemento del vector y por fuera de la diagonal están las covarianzas entre pares de elementos del vector aleatorio.

### PROPIEDADES DEL VALOR ESPERADO Y LA VARIANZA DE UN VECTOR ALEATORIO

i.  $E[AY] = AY$

$Y$  = vector de medias

ii.  $E[AY + b] = AY + b$

$\Sigma Y$  = Matriz var-cov

iii.  $V[AY] = A \Sigma Y A'$

$A$  = Matriz  $m \times n$

iv.  $V[AY + b] = A \Sigma Y A'$

$b$  = vector  $m \times 1$  (constante)

### DEFINICIONES BASICAS DE TEORIA MATRICIAL

i.  $(BA)^T = A^T B^T$

$A$  = Matriz constante ( $n \times n$ )

$B$  = Matriz Constante ( $m \times n$ )

ii.  $A$  es simétrica si  $A^T = A$

$X$  = vector de variables ( $n \times 1$ )

$A$  es idempotente si  $AA = A$

$a$  = vector de constantes ( $n \times 1$ )

$I$  = Matriz identidad ( $n$ )

iii. Si  $A$  es simétrica e idempotente entonces  $(I - A)$  también lo es

iv. FORMA CUADRÁTICA: La función  $X^T A X = \sum \sum a_{ij} x_i x_j$  es la forma cuadrática de  $X$  con  $a_{ij}$  la  $ij$ -ésima componente de  $A$

MODELOS MULTIREGRESIA

v. Si  $X^T A X > 0, \forall X$ ,  $A$  es definida positiva  
Si  $X^T A X \geq 0, \forall X$ ,  $A$  es definida semipositiva

### DERIVADAS VECTORIALES O MATRICIALES

a)  $\frac{\partial (a^T X)}{\partial X} = a$

$a$  = vector de constantes ( $n \times 1$ )

$A$  = Matriz de constantes ( $n \times n$ )

$X$  = vector de variables ( $n \times 1$ )

b)  $\frac{\partial (X^T X)}{\partial X} = 2X$

c)  $\frac{\partial (X^T A X)}{\partial X} = AX + A^T X$  - Si  $A$  es simétrica el resultado es  $= 2AX$

## RESULTADOS DISTRIBUCIONALES PARA VECTORES ALEATORIOS

Sea  $y$  un vector aleatorio normal  $n$ -variado con media  $\mu_y$  y matriz no singular de var-cov  $\Sigma_y$ , es decir

$$y \sim N_n(\mu_y, \Sigma_y)$$

Sea  $A$  una matriz  $n \times n$  de constantes y sea  $U$  una forma cuadrática de  $y$  definida por

$$U = y^T A y$$

Entonces:

i. Si  $A \Sigma_y$  o  $\Sigma_y A$  es una matriz idempotente de rango  $p$ , entonces

$$U \sim X_{p, k}^2 \quad \text{donde } k = \mu^T A \mu \text{ es el parámetro de no centralidad de la distribución chi-cuadrado}$$

ii. Sea  $\Sigma_y = \sigma^2 I$ , lo cual es una suposición típica, si  $A$  es idempotente y de rango  $p$  entonces

$$\frac{U}{\sigma^2} \sim X_{p, k}^2 \quad \text{donde } k = \mu^T A \mu / \sigma^2$$

iii. Sea  $B$  una matriz  $m \times n$  y  $W$  la forma lineal definida por:  $W = B y$ , entonces la forma cuadrática  $U = y^T A y$  son independientes entre sí

$$B \Sigma_y A = 0 \quad 0 = \text{matriz nula } (m \times n)$$

iv. Sea  $B$  una matriz  $n \times n$  y sea  $V = y^T B y$ , entonces las dos formas cuadráticas  $U$  y  $V$  son independientes entre sí

$$A \Sigma_y B = 0$$

NOTA: Si  $\Sigma_y = \sigma^2 I$ , entonces  $U$  y  $V$  son independientes si  $AB = 0$

## CONSIDERACIONES MRLM

1. Se dice lineal ya que la parte determinística es una función lineal de los parámetros desconocidos ( $\beta_0, \beta_1, \dots, \beta_p$ ) denominados coeficientes de regresión
2. Describe un hiperplano en el espacio de  $k$ -dimensiones de las variables regresoras  $X_j$ 's

3. Bj: Representa el cambio esperado en la variable respuesta y, por un cambio unitario en xj, cuando las demás variables regresoras xks (con  $k \neq j$ ) se mantienen constantes.

### ENFOQUE MATRICIAL DEL MODELO RLS

$$\begin{array}{l}
 y_1 = \beta_0 + \beta_1 x_1 + \varepsilon_1 \\
 y_2 = \beta_0 + \beta_1 x_2 + \varepsilon_2 \\
 \vdots \\
 y_n = \beta_0 + \beta_1 x_n + \varepsilon_n
 \end{array}
 \quad \Leftrightarrow \quad
 \begin{array}{l}
 \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \\
 n \times 1
 \end{array}
 \quad \mathbf{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \\
 2 \times 1$$

$$\mathbf{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix} \\
 n \times 1$$

$$\mathbf{\varepsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 x_1 \\ \beta_0 + \beta_1 x_2 \\ \vdots \\ \beta_0 + \beta_1 x_n \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix} \Leftrightarrow \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

AHORA

$$\mathbf{y} = \mathbf{X} \mathbf{\beta} + \mathbf{\varepsilon}$$

$n \times 1 \quad n \times 2 \quad 2 \times 1 \quad n \times 1$

$$\hat{\mathbf{\beta}} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}$$

$\bar{x} = \frac{\sum x_i}{n}$   
 $S_{xx} = \sum x_i^2 - n \bar{x}^2$

$$(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{bmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{bmatrix}$$

$$\mathbf{X}^T \cdot \mathbf{y} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \frac{1}{n \sum x_i^2 - n \bar{x}^2} \begin{bmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{bmatrix} \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \frac{1}{n [\sum x_i^2 - n \bar{x}^2]} \begin{bmatrix} (\sum x_i^2)(\sum y_i) - (\sum x_i)(\sum x_i y_i) \\ -(\sum x_i)(\sum y_i) + (\sum x_i y_i)(n) \end{bmatrix}$$

$$\hat{\beta}_1 = \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{n [\sum x_i^2 - n \bar{x}^2]} = \frac{n \sum x_i y_i - n \bar{x} n \bar{y}}{n [\sum x_i^2 - n \bar{x}^2]}$$

$$\hat{\beta}_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}$$

$$\hat{\beta}_0 = \frac{(\sum x_i^2)(\sum y_i) - (\sum x_i)(\sum x_i y_i)}{n [\sum x_i^2 - n \bar{x}^2]} = \frac{n \bar{y} \cdot \sum x_i^2 - n \bar{x} \sum x_i y_i}{n S_{xx}} \dots$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Como mencionamos anteriormente el MRLM se considera lineal siempre y cuando este en funcion lineal de los parametros, por lo que es valido encontrar modelos de la forma

### 1. Polinomial de tercer orden

POLINOMICA  $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + E$   
NO RECOMEN- Si hacemos  $x_1 = x$ ,  $x_2 = x^2$  y  $x_3 = x^3$  tenemos que  
DADA NO ES  
PARCIMONIOSA  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + E \rightarrow$  EL CUAL ES UN MRLM

### 2. MODELO CON Interaccion

El AUMENTO DE  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 \cdot x_2 + E$   
VARIABLES Si hacemos  $x_1 = x_1$ ,  $x_2 = x_2$ ,  $x_3 = x_1 \cdot x_2$  tenemos que  
DISMINUYE  
LA PRECISION  
EN LA  
ESTIMACION  
DEL MSE  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + E \rightarrow$  EL CUAL ES UN MRLM

Para estimar los parámetros  $\beta$  del MRLM utilizamos dos métodos: (1) Máxima Verosimilitud y (2) Mínimos Cuadrados los cuales resultan en los mismos estimadores.

1) Mínimos Cuadrados = Minimiza la suma de cuadrados de los errores

$$\begin{aligned} S(\beta) &= S(\beta_0, \beta_1, \dots, \beta_p) = \sum_{i=1}^n \epsilon_i^2 = \epsilon^T \epsilon \\ &= (y - X\beta)^T (y - X\beta) \\ &= y^T y - \beta^T X^T y - y^T X\beta + \beta^T X^T X\beta \end{aligned}$$

Para obtener una estimación de  $\beta$  derivamos vectorialmente respecto a  $\beta$  y obtenemos:

$$\frac{\partial S(\beta)}{\partial \beta} \Big|_{\hat{\beta}} = -2X^T y + 2X^T X \hat{\beta}$$

Igualando a cero obtenemos

$$X^T X \hat{\beta} = X^T y \quad \rightarrow \text{Ecuaciones Normales}$$

Si " $X^T X$ " es invertible o definida positiva entonces podemos hallar  $\beta$  y por tanto el vector estimado es

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

NOTA = Otra forma de hallar los estimadores  $\beta$ 's es utilizando la función  $S(\beta_0, \beta_1, \dots, \beta_p)$  y derivar con respecto a cada parámetro  $\beta_p$

$$\begin{aligned} S(\beta_0, \beta_1, \dots, \beta_p) &= \sum_{i=1}^n \epsilon_i^2 \\ &= \sum_{i=1}^n (y_i - E(y_i))^2 \\ &= \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 \end{aligned}$$

Cuando hay problemas de multicolinealidad  
REGRESION RIDGE

$(X^T X)^{-1} + \text{Kerces identidad}$



CONSIDERACIONES EMC → Estimadores de Mínimos Cuadrados

1.  $(X^T X)^{-1}$  = Existe, siempre que los regresores  $X_j$ 's sean LI
2. Los elementos de la diagonal de  $X^T X$  son las sumas de cuadrados de los elementos de las columnas de  $X$
3. Los elementos fuera de la diagonal de  $X^T X$  son las sumas de los productos cruzados de los elementos de las columnas de  $X$
4. Los elementos de  $X^T y$  son las sumas de los productos cruzados de los elementos de las columnas de  $X$  con los elementos de  $y$

INTERPRETACION  $\hat{\beta}_j$ 's

El valor de  $\hat{\beta}_j$  con  $j=1, 2, \dots, K$  se interpreta como el cambio esperado (o efecto parcial) en la respuesta promedio debido a un incremento unitario en  $x_j$ , cuando los demás predictores se mantienen fijos.

El valor de  $\hat{\beta}_0$  solo tiene sentido práctico interpretarlo si  $x_j=0 \in [\min x_j, \max x_j]$   $\forall j=1, 2, \dots, K$  y en ese caso se interpreta como la respuesta media cuando  $x_j=0$ ,  $\forall j=1, 2, \dots, K$ .

PROPIEDADES de  $\hat{\beta}$ 

i  $\hat{\beta} = (X^T X)^{-1} X^T y$  es insesgado para  $\beta$ , es decir  $E[\hat{\beta}] = \beta$

$$\begin{aligned}
 E[\hat{\beta}] &= E[(X^T X)^{-1} X^T y] \\
 &= (X^T X)^{-1} X^T E[y] \\
 &= (X^T X)^{-1} X^T E[X\beta + E] \\
 &= (X^T X)^{-1} X^T [X\beta + E(E)] \\
 &= (X^T X)^{-1} X^T X \beta \\
 &= I \beta = \beta
 \end{aligned}$$

ii  $\hat{\beta}$  es el mejor estimador lineal-insesgado de  $\beta$ , en el sentido de que  $\hat{\beta} = (X^T X)^{-1} X^T y$  tiene mínima varianza entre todos los estimadores insesgados de  $\beta$

$$\begin{aligned}
 \text{cov}(\hat{\beta}) &= \text{cov}((X^T X)^{-1} X^T y) \\
 &= (X^T X)^{-1} X \text{cov}(y) [X^T X)^{-1} X^T] \quad \left| \begin{array}{l} \text{cov}(y) = \text{cov}(\epsilon) \\ = \sigma^2 I \end{array} \right. \\
 &= (X^T X)^{-1} X^T \sigma^2 I X (X^T X)^{-1} \\
 &= \sigma^2 (X^T X)^{-1} \\
 &= \sigma^2 [c_{ij}]
 \end{aligned}$$

iii La matriz de Varianzas-Covarianzas de  $\hat{\beta}$  está dada por

$$\text{cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1} = \sigma^2 [c_{ij}] \quad \text{En donde los } c_{ij} \text{ son los elementos de la matriz } (X^T X)^{-1}$$

$$C = (X^T X)^{-1} \begin{bmatrix} c_{00} & c_{01} & \dots & c_{0k} \\ c_{01} & c_{11} & \dots & c_{1k} \\ \vdots & \vdots & \dots & \vdots \\ c_{0k} & c_{1k} & \dots & c_{kk} \end{bmatrix}$$

$$\text{cov}(\hat{\beta}) = \sigma^2 C \quad \text{Donde } V(\hat{\beta}_j) = \sigma^2 c_{jj} \text{ y}$$

$$\text{cov}(\hat{\beta}_i, \hat{\beta}_j) = \sigma^2 c_{ij}$$

$$\text{Var}[\hat{\beta}] = \text{Var}[(X^T X)^{-1} X^T y]$$

$$= \text{Var}[A y] \quad \text{con } A = (X^T X)^{-1} X^T$$

$$= A \text{Var}[y] A^T$$

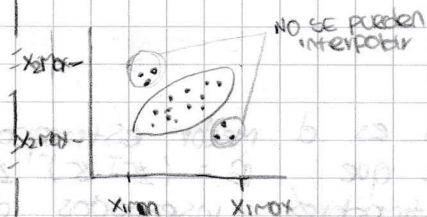
$$= A \sigma^2 I A^T$$

$$= \sigma^2 A A^T$$

$$= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1}$$

$$= \sigma^2 (X^T X)^{-1}$$

PARA INTERPOLAR



NO BASTA CON QUE ESTE EN EL RANGO DE  $x_1, x_2$  SINO QUE ESTE EN LA NUBE

# MODELO DE REGRESION LINEAL MULTIPLE (SEGUNDA PARTE)

En general observamos un modelo del tipo

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \epsilon_i$$

Donde  $\epsilon_i \sim N(0, \sigma^2)$ ,  $i = 1, 2, 3, \dots, n$

Matricialmente lo escribimos

$$y = X\beta + \epsilon$$

Donde

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix}_{n \times p}$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}_{p \times 1}$$

$$\epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}_{n \times 1}$$

Ecuciones Normales

$$X^T X \hat{\beta} = X^T y$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

$p = \#$  de parámetros =  $k+1$   
 $k = \#$  de predictores

$$E[\hat{\beta}] = \beta$$

$$\text{Cov}[\hat{\beta}] = \sigma^2 (X^T X)^{-1}$$

## ECUACIONES NORMALES

$$X^T X \hat{\beta} = X^T y$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

## CURVA DE REGRESION AJUSTADA

$$\hat{y}_i = \hat{E}(y_i | x_{i1}, x_{i2}, \dots, x_{ik})$$

$$= \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \dots + \hat{\beta}_k x_{ik}$$

$$= [1 \ x_{i1} \ \dots \ x_{ik}] \hat{\beta}$$

$$= x_{oi}^T \hat{\beta}$$

$$\hat{y} = \hat{E}(y | X)$$

$$= X\hat{\beta}$$

$$= X(X^T X)^{-1} X^T y = Hy$$

Propiedades de  $H$

(Matriz hat o sombrero)

- Es simétrica

- Es idempotente  $HH = H$

- Es cuadrada ( $n \times n$ )

$$\hat{y} = Hy$$

Residuos

$$e_i = y_i - \hat{y}_i$$

$$e = y - \hat{y}$$

$$e = y - Hy$$

$$e = (I - H)y$$

$\hookrightarrow$  Simétrico  
idempotente

## DISTRIBUCION DE $e$

$$E[e] = E[(I-H)y]$$

$$= (I-H)E[y]$$

$$= (I-H)X\beta$$

$$= X\beta - (HX)X\beta$$

$$= X\beta - X\beta$$

$$= \mathbf{0} \rightarrow [E[\hat{y}_i] = E[y_i]]$$

$$\text{Cov}[e] = \text{Cov}[(I-H)y]$$

$$= (I-H)\text{Cov}[y](I-H)$$

$$= (I-H)\sigma^2 I (I-H)$$

$$= \sigma^2 (I-H)(I-H)$$

$$= \sigma^2 (I-H)$$

Los residuos son correlacionados  
 ↳ hay que estandarizarlos para ver  
 tendencia

$$e \sim N(0, \sigma^2 (I-H))$$

Los residuos distribuyen normal  
 N variado, con vector de  
 medias cero y varianzas  
 constante

$$\text{Var}[e_i] = \sigma^2 (1-h_{ii})$$

Diagonal Matriz  
 Sobrevivo

$$\text{Cov}[e_i, e_j] = \sigma^2 (1-h_{ij}) \rightarrow \neq 0$$

↳ Los  $e_i$  son correlacionados

$$\text{Cov}(e) = \sigma^2 (I-H)$$

EN R se calcula la  
 Inversa generalizada via  
 Métodos Numéricos

## VERIFICACION DE SUPUESTOS

i) Distribucion Normal de los Errores [ $e_i = y_i - \hat{y}_i$ ]

$$H_0: e_i \sim N$$

$$H_a: e_i \neq N$$

Utilizamos la prueba de  
 Shapiro Wilk y  
 Concluimos con p-value

La prueba F no es  
 significativa si no  
 se cumple la  
 Normalidad de los  
 residuos

NO olvidemos que los  
 $e_i$  son correlacionados

Estandarizamos los residuos

$$v_i = \frac{e_i}{\sqrt{\text{MSE}(1-h_{ii})}}$$

Donde  $h_{ii}$  es la  
 diagonal de H

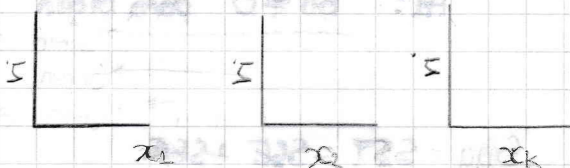
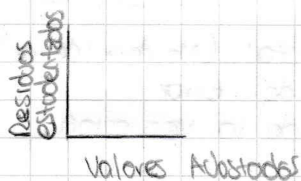
ii) Varianza Constante (Homogeneidad de Varianzas)

$Var(\epsilon_i) = \sigma^2$

Cuando el modelo no tiene falta de ajuste se llama adejado

\* Contrastamos  $\epsilon_i$  vs  $y_i$

\* Contrastamos  $\epsilon_i$  vs  $x_{ij}$ ,  $i=1,2,3...K$



Estos graficos tambien sirven para detectar "Outliers" y nos dicen si necesitamos transformaciones para la respuesta ( $y_i$ ), sobre una variable regresora ( $x_{ij}$ ) o sobre varias variables regresoras.

TENDENCIA IDEAL

NO CONSTANTE (PARABOLICA)

Rombo

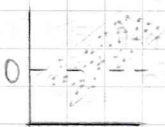


La variabilidad lo aumentamos

La variabilidad lo disminuimos

FALTA VARIABLE LINEAL

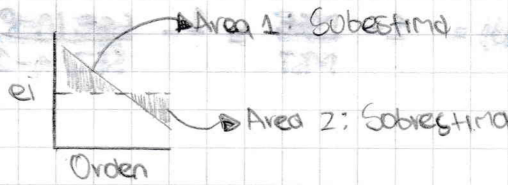
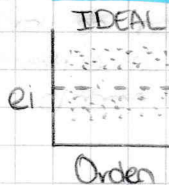
FALTA VARIABLE CUADRATICA



Con estos graficos tambien podemos evidenciar la linealidad del modelo

iii) Incorrelacion de los  $\epsilon_i$

Solo se realiza cuando se tiene el orden en el que se obtuvo cada dato (Observacion)



Lo podemos utilizar (row.names)

Hay dos opciones en este caso

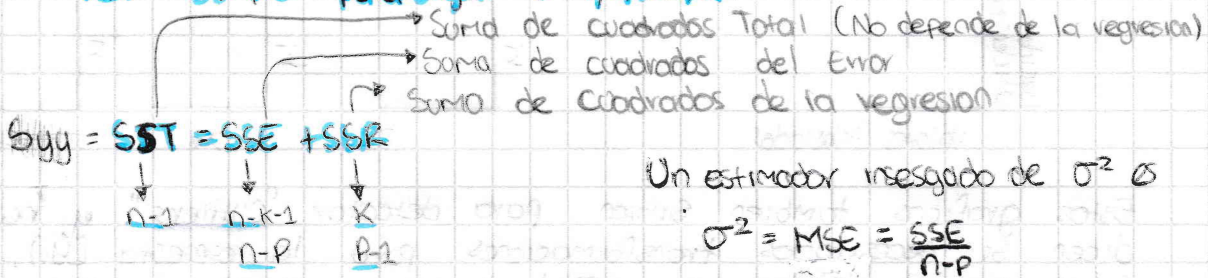
- 1) Incluir Correlacion en el modelo
- 2) Repetir la toma de datos

# PRUEBA DE SIGNIFICANCIA DE LA REGRESION

$H_0: \beta_1 = \beta_2 = \dots = \beta_K = 0$

$H_1: \beta_j \neq 0$  para algún  $j = 1, 2, 3, \dots, K$

Rechaz  $H_0$   
Al menos una de las predictoras es significativa para explicar la variabilidad.



- Tabla anova (Funcion Mi Anova)

Funcion creada para la clase

	gl	SS	MS	F	P-value
Regresion	$k-p-1$	SSR	$MSE = SSR/k$	$F_0 = \frac{MSR}{MSE}$	$P(F_{k,n-p} > F_0)$
Error	$n-p$	SSE	$MSE = SSE/n-p$		
Total	$n-1$	SST			

$F_0 = \frac{MSR}{MSE} \sim F_{k,n-p}$

App. Probability Distributions  
 $\Rightarrow$   
 $d_1 = k$   
 $d_2 = n-p \Rightarrow P(X > x) = P\text{-value}$   
 $X = F_0$

- Coeficiente de determinacion

$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$

Proporcion de varianza explicada por el modelo, en el caso de RLM el  $R^2$  aumenta a medida que aumenta el numero de variables regresoras

$R_{adj}^2 = 1 - \frac{MSE}{MST} = 1 - \frac{SSE/n-p}{SST/n-1}$

No es Proporcion de varianza explicada, es un criterio de seleccion de modelo

# PRUEBA DE SIGNIFICANCIA DE CADA PARAMETRO ( $T_{0j}$ )

$H_0: \beta_j = 0$   
 $\beta_j \neq 0$

$t_{0j} = \frac{\hat{\beta}_j - \beta_j}{ee(\hat{\beta}_j)}$

IC  $\beta_j \pm t_{n-p, \frac{\alpha}{2}} \cdot ee(\hat{\beta}_j)$

Summary (modelo)

$P\text{-value} = 2P(|t_{n-p}| > |t_{0j}|)$

$Var[\beta] = \begin{bmatrix} c_{11} & c_{22} & c_{ij} \end{bmatrix}$

Componente (i,i) de la matriz C

Donde

$ee(\hat{\beta}_j) = \sqrt{Var[\hat{\beta}_j]} = \sqrt{\sigma^2 \cdot c_{jj}} = \sqrt{MSE \cdot c_{jj}}$

$\beta_0 =$  NO significativo  
NO interpretable Significativo  
interpretable Significativo  
NO interpretable DE  $[X_{\min}, X_{\max}]$   $\forall j=1, \dots, k$   
si no solo no cumple ya  
no es interpretable

La interpretación de la prueba la expresamos de la siguiente manera

Por cada unidad de incremento de la variable regresora ( $X_k$ ), se espera que la cantidad de la variable explicada, en promedio, se incrementa de manera significativa, en el (valor del parámetro), cuando las demás variables regresoras se mantienen constantes.

- \* Si el signo del parámetro es positivo = INCREMENTO
- Si el signo del parámetro es negativo = DISMINUYA

En general un parámetro es interpretable solo si es significativo

### INFERENCIA DE UN SUBCONJUNTO DE PARAMETROS

Lo utilizamos para probar la significancia de un subconjunto de parámetros o para ver si el efecto de una variable es igual al de otra variable.

Para la construcción del estadístico de prueba se utiliza la metodología de sumas de cuadrados extra, la cual cuantifica el incremento o disminución en la SSR o SSE al incluir una o varias variables al modelo de regresión.

Un aumento unitario en  $X_j$  se estima hay un cambio promedio de  $\beta_j$  en la respuesta cuando las demás  $X$  están constantes

EN GENERAL

"  $H_0: \beta_2 = 0$   
 $H_a: \beta_2 \neq 0$  "

Donde  $\beta_2$  contiene los componentes de  $\beta$  a los que se desea analizar la no significancia en el modelo de regresión particionado  $\beta$

Por tanto, en esta situación compararemos el modelo completo (FM) contra un modelo reducido (RM)

$\beta_2$  = Contiene las condiciones de interés

$\beta_1$  = contiene  $\beta_0$  y los demás componentes de  $\beta$  que no están en  $\beta_2$

$SS_{extra} = SSR(\beta_2 | \beta_1)$

$= SSR(\beta_2, \beta_1) - SSR(\beta_1)$   
 $= SSR(FM) - SSR(RM)$   
 $= [SST - SSE(FM)] - [SST - SSE(RM)]$   
 $= SSE(RM) - SSE(FM)$   
 $= MSE(RM) * (n - p_{RM}) - MSE * (n - p_{FM})$

Estadístico de prueba

$$F_{\text{parcial}} = \frac{SSR(\beta_2 | \beta_1) / M}{MSE(\beta_1, \beta_2)} \sim F_{M, n-p_{FM}}$$

Donde:  $M = \dim(\beta_2)$

La suma de cuadrados de la regresión se descompone como sumas de cuadrados extra

$$SSR(\beta) = \underbrace{SSRE(\beta_1 | \beta_0, \beta_2, \dots, \beta_k)}_{1 \text{ gl}} + \underbrace{SSRE(\beta_2 | \beta_0, \beta_1, \beta_3, \dots, \beta_k)}_{1 \text{ gl}} + \dots + \underbrace{SSRE(\beta_k | \beta_0, \beta_1, \beta_2, \dots, \beta_{k-1})}_{1 \text{ gl}}$$

Veremos dos tipos de sumas de cuadrados

\* Suma de cuadrados tipo I ( $SS_1$ ) No muy útil

Es una suma de cuadrados secuencial de cada variable, dado que en el modelo están presentes las variables anteriores

En R anova (modelo)

\* Suma de cuadrados tipo II ( $SS_2$ )

Es una suma de cuadrados parcialmente secuencial de cada variable, dado que en el modelo están presentes las demás variables

En R Anova (modelo)

Podemos probar la significancia de un parámetro con  $SS_{\text{extra}}$ , de forma homogénea como lo hacemos con la prueba T. y se llega a la misma conclusión.

$$H_0: \beta_2 = 0$$
$$H_a: \beta_2 \neq 0$$

> El modelo reducido correspondiente sería

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_3 x_{i3} + \dots + \beta_k x_{ik}$$

Estadístico de prueba

$$F_0 = \frac{[SSE(RM) - SSE(FM)] / 1}{MSE(FM)} \sim F_{1, n-p_{FM}}$$

Los grados de libertad del numerador en este caso son el # de parámetros

cuadrados ( $\beta_1$ )



## - Significancia de un subconjunto de parámetros

Buscamos probar la no significancia de un subconjunto de parámetros, en general la hipótesis que subyace es la siguiente

Partiendo de un modelo de 6 variables regresoras podemos probar la significancia de digamos 3 de ellos al mismo tiempo.

$$H_0: \beta_4 = \beta_5 = \beta_6 = 0$$

$$H_a: \text{Al menos un } \beta_j \neq 0$$

Como se ve la prueba es similar a la  $F_{\beta}$  de significancia de la regresión y conviene de la misma manera

Estadístico de prueba

$$F_0 = \frac{[SSE(RM) - SSE(FM)]/r}{MSE(FM)} \sim F_{r, n-p_{FM}}$$

Si rechaza  $H_0$

Donde  $r$  es el número de parámetros que estamos evaluando

↳ Para este caso  $r=3$

↳ Al menos uno de los  $\beta_j$  en estudio es significativo para explicar la variabilidad

$$RR) F_0 / F_{\alpha, r, n-p_{FM}}$$

El modelo reducido (RM) se ve de la siguiente forma

$$V_p = P(F_{r, n-p_{FM}} > F_0)$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \epsilon$$

$\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$

## - Prueba Lineal General

Con esta prueba verificamos si el efecto de una variable es igual al de otra.

Volviendo al modelo con 6 regresoras analicemos

$$H_0: \beta_1 = \beta_2 \wedge \beta_5 = \beta_6$$

↳ Esta escritura dificulta el chequeo de hipótesis, entonces se escribe de forma matricial

$$\begin{matrix} \beta_1 - \beta_2 = 0 \\ \beta_5 - \beta_6 = 0 \end{matrix} \rightarrow \beta$$

$$\begin{bmatrix} 0 & 1 & -2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \\ \beta_6 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

De forma Matricial las hipótesis quedan de la forma

$H_0: \pi\beta = 0$

$H_a: \pi\beta \neq 0$

Forma general  $\rightarrow$  El Modelo reducido (RM) se ve de la siguiente forma.

$$y_i = \beta_0 + \beta_1(x_{i2} + x_{i2}) + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5(x_{i5} + x_{i6})$$

$$F_0 = \frac{[SSE(RM) - SSE(FM)]/m}{MSE(FM)} \sim F_{m, n-pm}$$

Donde:

$m =$  Es el # de filas LI de  $\pi$   
 $\dim(\pi)$

Regresión por mínimos cuadrados ponderados

Recordemos que para la inferencia de un subconjunto de parámetros sea cual sea siempre tenemos un  $\beta_2$  y que los grados de libertad  $r$  o  $m$  siempre son  $\dim(\beta_2)$ , solo que en el caso de la significancia de un subconjunto de parámetros  $\dim(\beta_2)$  es igual al # de parámetros que analizamos

OBSERVACIONES O VALORES ATÍPICOS ("outliers")

Estas observaciones o valores atípicos afectan considerablemente los supuestos del Modelo ya que tienen efecto sobre:

- Estimaciones de los  $\beta_j$ 's
- Aumento considerable de  $ee(\hat{\beta}_j)$
- Supuestos del modelo

Antes de abordar en la identificación de los outliers miremos los tipos de residuos

-  $e_i = y_i - \hat{y}_i \rightarrow$  R. Crudos

-  $r_i = \frac{d_i}{\sqrt{1-h_{ii}}} = \frac{e_i}{\sqrt{MSE(1-h_{ii})}} \sim t_{n-p}$   
 $\rightarrow$  R. Estándarizados

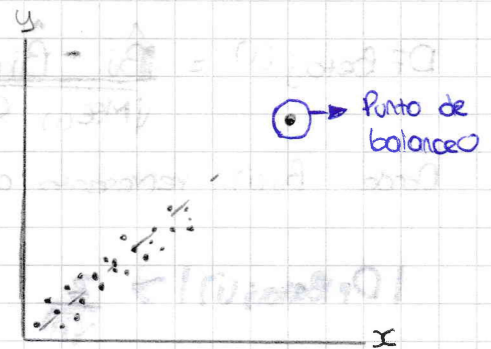
-  $d_i = \frac{e_i}{\sqrt{MSE}} \rightarrow$  R. Estándarizados

Donde  $h_{ii}$ : Es el elemento de la diagonal de  $H = X(X^T X)^{-1} X^T$   
 Representa la distancia al centro de los datos

Las Observaciones atípicas se clasifican en Puntos de balanceo y puntos influyentes

- Punto de balanceo (PB)

Punto alejado del resto de los datos. Posiblemente no afecta los  $\beta_j$ 's pero si el  $R^2$  y los  $ee$  ( $\beta_j$ 's)



Tiene un valor moderadamente inusual de las variables predictoras y en la variable respuesta.

CRITERIO

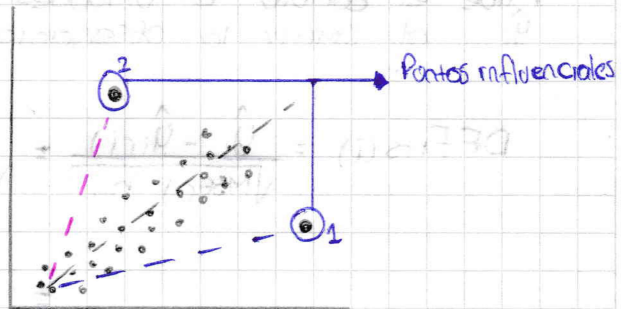
$$h_{ii} > 2 \left( \frac{k+1}{n} \right) \equiv h_{ii} > \frac{2p}{n}$$

Valor inusual en  $x$  y en  $y$

- Punto influyente (PI)

Tiene un impacto importante en los  $\beta_j$ 's ya que el punto influyente "hala" el modelo en su dirección.

Afecta los supuestos de modelo



--- Modelo dado PI1  
- - - Modelo dado PI2

CRITERIO

Para ser candidato a PI se debe cumplir que

- $h_{ii}$  grande ( $h_{ii} \approx 1$ )
  - $|r_i|$  grande ( $|r_i| > 3$ )
- Candidato a PI (99%)

- $h_{ii}$  grande
  - $|r_i| > 2$
- Candidato a PI (95%)

Medidas para detectar observaciones influyentes

1) Distancia de Cook:

Utilizamos  $(D_i)$ , si su valor es grande la observación es influyente

$$D_i = \frac{(\hat{y} - \hat{y}_{(i)})^T (\hat{y} - \hat{y}_{(i)})}{(k+1)MSE} = \frac{(\hat{\beta} - \hat{\beta}_{(i)})^T X^T X (\hat{\beta} - \hat{\beta}_{(i)})}{(k+1)MSE} = \frac{r_i^2}{k+1} \cdot \frac{h_{ii}}{1-h_{ii}}$$

Donde el numerador es una distancia euclidiana y  $\hat{y}_{(i)}$  es la estimación quitando la  $i$ -ésima observación

$D_i$  es alto si:  $D_i > F_{0,90, k+1, n-p-1} \equiv D_i > 1$   
 ↳ Entonces es INFLUYENTE

## 2) DF Betas:

Mide el cambio del  $\hat{\beta}_i$  en unidades estandar al omitir la observacion  $i$

$$DF\text{ Betas}(i) = \frac{\hat{\beta}_i - \hat{\beta}_{(i)}}{\sqrt{MSE(i)} \cdot C_{(i),ii}} = \frac{r_{ii}}{\sqrt{r_{ii} r_{ii}}} \cdot \frac{e_i}{\sqrt{MSE} \cdot (1-h_{ii})}$$

Donde  $\hat{\beta}_{(i)}$  representa al modelo cuando le quite la misma observacion

$$|DF\text{ Betas}(i)| > \frac{2}{\sqrt{n}}$$

Si es influyente en el parametro es influyente en el modelo.

## 3) DF Fits

Mide el cambio en unidades estandar del valor ajustado  $\hat{y}_i$  al omitir la observacion  $i$

$$DF\text{ Fits}(i) = \frac{\hat{y}_i - \hat{y}_{(i)}}{\sqrt{MSE(i)} \cdot h_{ii}} = \sqrt{\frac{h_{ii}}{1-h_{ii}}} \cdot \frac{e_i}{\sqrt{MSE(i)} \cdot (1-h_{ii})}$$

$$= \sqrt{\frac{h_{ii}}{1-h_{ii}}} \cdot r_i$$

$$|DF\text{ Fits}(i)| > 2 \sqrt{\frac{k+1}{n}}$$

Si la observacion se clasifica influyente con alguno de los 3 criterios la observacion sera influyente.

## IC

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-p} * ee(\hat{y}_0)$$

Donde  $\hat{y}_0 = X_0 \cdot \hat{\beta} = y_{ohat}$

$$Se. y_{ohat} = ee(\hat{y}_0) = \sqrt{MSE (X_0 (X^T X)^{-1} X_0^T)}$$

## IP

$$\hat{y}_0 \pm t_{\frac{\alpha}{2}, n-p} * ee(\hat{y}_0 - y_0)$$

Donde

$$ee(\hat{y}_0 - y_0) = \sqrt{MSE + ee^2(\hat{y}_0)}$$

Con una confianza del  $(1-\alpha)\%$  la respuesta media (Promedio) se encuentra entre  $(U, L)$  Para las condiciones dadas

Con una confianza del  $(1-\alpha)\%$  se predice que la respuesta se encuentra entre  $(U, L)$  Para las condiciones dadas

No olvidemos que los IC y los IP solo se pueden calcular si las condiciones dadas son interpolables, lo cual verificamos de la siguiente manera

$$\text{Mín. valor: } X_0^T (X^T X)^{-1} X_0 \leq \max |h|$$

Donde  $X_0$  contiene las condiciones de interés

## MULTICOLINEALIDAD

Esta presente cuando en el modelo de RLM hay dependencias lineales entre las variables regresoras. Si hay multicolinealidad en el modelo se pueden presentar inestabilidad numérica en el cálculo de la inversa asociada a  $H$

$$\rightarrow (X^T X)$$

En R para evitar los problemas de multicolinealidad en la matriz se calcula la inversa generalizada

### - Fuentes de multicolinealidad

- \* Método de recolección de los datos
- \* Restricciones en el modelo o población  $\rightarrow$  Consumo Eléctrico = # habitantes
- \* Modelo sobredefinido ( $K > n$ )
- \* Elección del modelo (Adición de términos polinomiales)  $\rightarrow$  Potencia

### - Detección de multicolinealidad

- \* Varianza de los estimadores inflada (muy grande o muy pequeña con respecto a la estimación de uno de los parámetros)
  - Parámetros resultan ser significativos cuando no lo son
  - Parámetros resultan no significativos cuando realmente lo son
- \* Cuando al realizar la prueba F resulta significativa y los parámetros individuales (usando  $t$ ) resultan no significativos
- \* Signos de los parámetros estimados son contrarios a lo esperado.

### - Diagnóstico de Multicolinealidad

# 1) Analisis de correlacion de las variables regresoras

↳ Funcion  $cor(\text{Datos})$

$$R = [r_{ij}] = cor(x_i, x_j)$$

Si  $r_{ij} \approx 1$ , entonces hay indicios de multicolinealidad

La correlacion de las variables regresoras con la variable explicada debe ser alta

$|r_{ij}| > 0,5 \rightarrow$  Normalmente se espera que entre las variables regresoras sea pequena la correlacion

# 2) VIF (Factor de inflacion de varianza) $\rightarrow$ Funcion $miscoeficientes()$

Se ajusta un modelo tipo combinacion lineal

↳ (modelo, datos)

$$x_j \text{ vs } x_1, x_2, \dots, x_{j-1}, x_{j+1}, \dots, x_k$$

Para  $j=1$

$$x_1 = \alpha_2 x_2 + \alpha_3 x_3 + \dots + \alpha_k x_k + C_0 + e$$

Con el modelo se calcula  $\rightarrow R_j^2 =$  Coeficiente de determinacion el cual se utiliza para determinar el VIF.   
 Múltiple del modelo  $j$

$$VIF = \frac{1}{1 - R_j^2}$$

$VIF_j < 5$  No hay multicolinealidad

$5 \leq VIF_j < 10$  Hay multicolinealidad moderada

$VIF_j \geq 10$  Hay multicolinealidad grave

# 3) Analisis de los valores propios de $(X^T X)$

↳ Funcion  $misDiagnosticos()$  (modelo)

Este analisis se puede realizar sin centrar o centrado.

$$\frac{x_{ij} - \bar{x}_j}{s_{x_j}}$$

↳ Se hace centrado cuando el intercepto no tiene interpretacion y sucede cuando  $0 \notin [x_{j\min}, x_{j\max}]$  para algun  $j$

↳ (Centrar=T)

## i) Numero de condicion

$$\kappa = \frac{\lambda_{\max}}{\lambda_{\min}} \quad \text{Criterio}$$

La funcion  $cond$   $\sqrt{\kappa}$

$\kappa < 100$  no hay multicolinealidad

$100 \leq \kappa < 1000$  Multicolinealidad moderada

$\kappa \geq 1000$  Multicolinealidad grave

ii) Indicadores o índices de condición

cond. index:  $K = \frac{\lambda_{max}}{\lambda_j}$ ,  $j = 1, 2, 3, \dots, K$

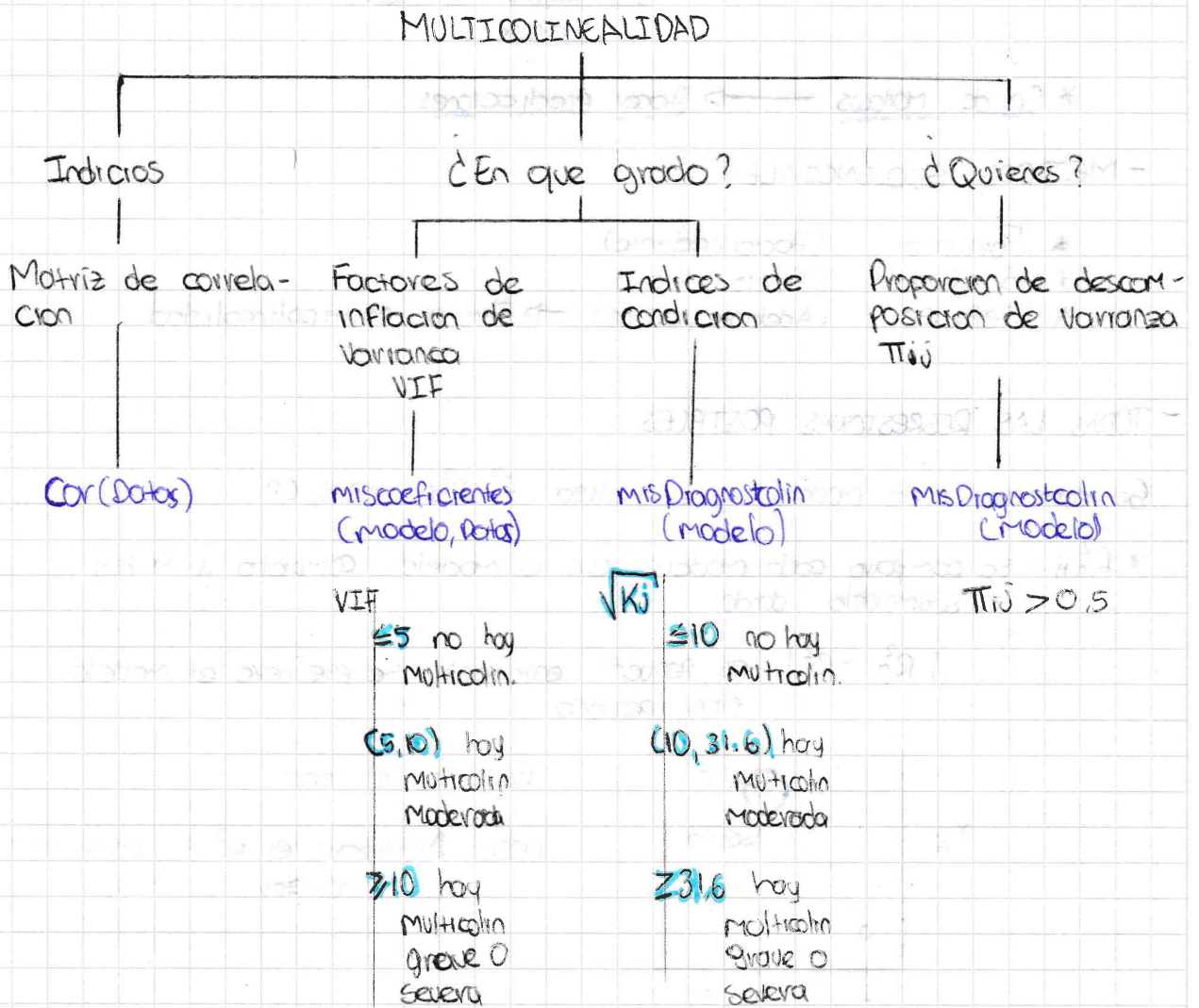
$K_j < 100$   $\forall j = 1, 2, 3, \dots, K$  no hay multicolinealidad  
 $100 \leq K_j < 1000$   $\forall j = 1, 2, 3, \dots, K$  Hay multicolinealidad  
 $K_j \geq 1000$   $\forall j = 1, 2, 3, \dots, K$  Multicolinealidad grave o severa.

iii) Proporción de descomposición de varianza ( $\pi_{ij}$ )

↳ Cuando se ve multicolinealidad grave o severa

Representa la proporción de varianza de cada  $\beta_j^1$  (o de cada VIF) debida al íesimo valor propio

$\pi_i, \pi_j$



## SELECCION DE MODELO

Suponga que se tiene un modelo con  $K$  variables regresoras y desea hallar un modelo reducido que utilice menos de  $K$  variables regresoras.

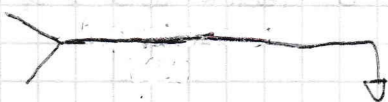
$$\binom{K}{1} + \binom{K}{2} + \binom{K}{3} + \dots + \binom{K}{K} = \sum_{j=1}^K \binom{K}{j} = 2^K - 1$$

Para la selección de modelos dos métodos:

### - METODO DE TODAS LAS REGRESIONES POSIBLES

↳ allregtable (Modelo, Resposta)

- \*  $R^2$
- \*  $R^2_{adj}$



Estimaciones o hacer inferencia

\* Cp de Mallows → Hacer predicciones

### - METODOS SECUENCIALES

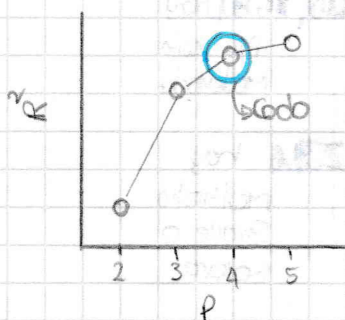
- \* Forward (Hacia adelante)
- \* Backward (Hacia atrás)
- \* Stepwise (Adelante, atrás) → Elimina multicolinealidad

### - TODAS LAS REGRESIONES POSIBLES

Se compara cada modelo y se analiza  $R^2$ ,  $R^2_{adj}$ , MSE, Cp.

\*  $R^2_{adj}$  = se compara cada modelo con el modelo completo y si hay un submodelo donde

$|R^2_p - R^2|$  es pequeño entonces se prefiere el modelo más pequeño



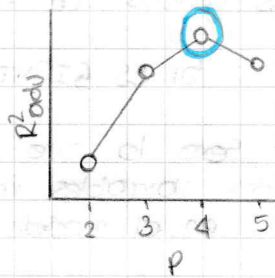
Buscamos el código

NOTA: Analizar el  $R^2$  es equivalente a analizar el MSE



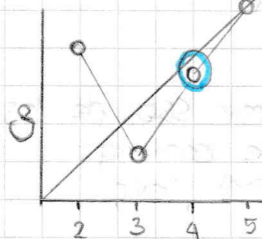
\*  $R^2_{adj}$ : Se busca el que tenga un Mayor valor, ya que el  $R^2_{adj}$  penaliza en función del número de variables

→  $1 - \frac{MSE}{MST}$



Buscamos el modelo que tenga el Mayor  $R^2_{adj}$

\* Cp de Mallows: El mejor modelo es el de cp más pequeño comparado con el p correspondiente



$|Cp - p| \approx 0$

El modelo más cercano a cero es insesgado

p = # de parámetros del modelo analizado

Buscamos el modelo con Cp más cercano a p dentro de todos los modelos.

El cp del modelo completo siempre es igual a p

El Cp es una medida del sesgo del modelo de regresión, es decir

$(E[\hat{y}_i] - E[y_i])$

$Cp = \frac{SSEP}{MSE} - (n - 2p)$

= Modelo completo

Se puede demostrar que  $E[Cp | \text{sesgo} = 0] = p$  por ello se eligen modelos.

$Cp \approx p$

- METODOS SECUENCIALES

\* FORWARD (Hacia adelante)

Se parte de un modelo con intercepto de la forma

$y_i = \beta_0 + \epsilon_i$

Paso 1. El parámetro candidato a entrar será el que tenga Mayor SSR  $\equiv$  Menor SSE dentro de los parámetros con una variable regresora y probamos su significancia en el modelo utilizando  $SSE_{extra}$

$F_0 = \frac{SSR(\beta_{j1} | \beta_0)}{MSE(FM)} \sim F_{1, n-p}$

$H_0: \beta_{j1} = 0$

$H_a: \beta_{j1} \neq 0$

NO olvidemos que el modelo completo en el primer paso tiene  $\beta_0, \beta_1$  por tanto  $P=2$ , además recordemos que para cada paso los grados de libertad del numerador siempre son "1"

SI ES SIGNIFICATIVA ENTRA

Paso 2. La candidata a entrar sera la que tenga menor SSE de los modelos con dos variables regresoras que contenga  $\beta_1$ . Los parámetros en el modelo (FM) seran  $\beta_0, \beta_1, \beta_2$  por tanto  $P=3$

$$H_0: \beta_2 = 0$$

$$H_a: \beta_2 \neq 0$$

$$F_0 = \frac{SSR(\beta_2 | \beta_0, \beta_1)}{MSE(FM)} \sim F_{1, n-3}$$

SI ES SIGNIFICATIVA ENTRA

Paso 3. La candidata a entrar sera la que tenga menor SSE de los modelos con tres variables regresoras que contenga  $\beta_1, \beta_2$ . Los parámetros en el modelo (FM) seran  $\beta_0, \beta_1, \beta_2, \beta_3$  por tanto  $P=4$

$$H_0: \beta_3 = 0$$

$$H_a: \beta_3 \neq 0$$

$$F_0 = \frac{SSR(\beta_3 | \beta_0, \beta_1, \beta_2)}{MSE(FM)} \sim F_{1, n-4}$$

SI ES SIGNIFICATIVA ENTRA

En general se van probando variables una a una y si son significativas entran, esto hasta que ya no hayan mas variables por probar.

PASO 1 = Prueba sobre  $\beta_1$   
 FM:  $\beta_1, \beta_0$   
 $F_{1, n-2}$   
 Significativa: Entra

PASO 2 = Prueba sobre  $\beta_2$   
 FM:  $\beta_1, \beta_2, \beta_0$   
 $F_{1, n-3}$   
 Significativa: Entra

PASO 3 = Prueba sobre  $\beta_3$   
 FM:  $\beta_1, \beta_2, \beta_3, \beta_0$   
 $F_{1, n-4}$   
 Significativa: Entra

PASO 4 = Prueba sobre  $\beta_4$   
 FM:  $\beta_1, \beta_2, \beta_3, \beta_4, \beta_0$   
 $F_{1, n-5}$   
 Significativa: Entra

⋮

La primera candidata por lo general es la de F parcial mas grande

Hasta que no hayan mas variables regresoras (K)

Falle alguno de los pasos, entonces se detendra

## \* Backward (Hacia atras)

Se parte de un modelo de la forma

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \dots + \beta_K X_{iK}$$

↳ El modelo abastado con todas las variables (Modelo Inicial)

PASO 1 El parametro candidato a salir sera el que no este en el modelo con Mayor SSR  $\Rightarrow$  Menor SSE que tiene  $K-1$  Variables regresoras

$$F_0 = \frac{SSR(\beta_{j2} | \beta_0, \beta_1, \dots, \beta_{j1-1}, \beta_{j1+1}, \beta_K)}{MSE(FM)} \sim F_{1, n-p}$$

↳ MODELO INICIAL

$H_0: \beta_{j2} = 0$  Si no es significativa sale  
 $H_a: \beta_{j2} \neq 0$

NO olvidemos que el modelo completo en el primer paso tiene todas las variables regresoras recolectadas por tanto  $P = K+1$ , ademas recordemos que para cada paso los grados de libertad del numerador siempre son "1".

PASO 2 El parametro candidato a salir sera el que no este en el modelo con menor SSE que tiene  $K-2$  Variables regresoras en el cual no esta  $\beta_{j2}$

$$F_0 = \frac{SSR(\beta_{j2} | \beta_0, \beta_1, \beta_{j2-1}, \beta_{j2+1}, \dots, \beta_{j2-1}, \beta_{j2+1}, \dots)}{MSE(FM)} \sim F_{1, n-k}$$

↳ modelo sin  $\beta_{j2}$

$H_0: \beta_{j2} = 0$  Si no es significativa sale  
 $H_a: \beta_{j2} \neq 0$

PASO 3 El parametro candidato a salir sera el que no este en el modelo con menor SSE que tiene  $K-3$  Variables regresoras en el cual no estan  $\beta_{j1}$  y  $\beta_{j2}$

$$F_0 = \frac{SSR(\beta_{j3} | \beta_0, \beta_1, \beta_{j1-1}, \beta_{j1+1}, \beta_{j2-1}, \beta_{j2+1}, \dots, \beta_{j3-1}, \beta_{j3+1}, \beta_{K-2})}{MSE(FM)} \sim F_{1, n-k-1}$$

↳ modelo sin  $\beta_{j1}$  y  $\beta_{j2}$

$H_0: \beta_{j3} = 0$   
 $H_a: \beta_{j3} \neq 0$  Si no es significativa sale.

En general se van probando variables una a una y si no son significativas salen, esto hasta llegar a un modelo solo con intercepto o hasta que falle alguno de los pasos.

Paso 1 = Prueba sobre  $\beta_{j1}$   
 FM: Modelo inicial  
 $F_{1, n-k-1}$   
 no significativa: sale

Paso 2 = Prueba sobre  $\beta_{j2}$   
 FM: Modelo sin  $\beta_{j1}$   
 $F_{1, n-k}$   
 no significativa: sale

Paso 3 = Prueba sobre  $\beta_{j3}$   
 FM: Modelo sin  $\beta_{j1}$  y  $\beta_{j2}$   
 $F_{1, n-k-1}$   
 no significativa: sale

Paso 4 = Prueba sobre  $\beta_{j4}$   
 FM: Modelo sin  $\beta_{j1}, \beta_{j2}, \beta_{j3}$   
 $F_{1, n-k-2}$   
 no significativa: sale

⋮  
 Hasta llegar al  
 modelo solo  
 con  $\beta_0$

o  
 falle algún paso, → que no se elimine  
 entonces se detiene una variable

\* STEPWISE

↳ El mejor = Elimina las que tienen multicolinealidad

En los primeros pasos es igual al método Forward y después analiza significancia de subconjunto de parámetros

PASO 1 = Prueba sobre  $\beta_{j1}$   
 FM:  $\beta_{j1}, \beta_0$   
 $F_{1, n-2}$   
 Significativo = entra

PASO 2 = Prueba sobre  $\beta_{j2}$   
 FM:  $\beta_{j1}, \beta_{j2}, \beta_0$   
 $F_{1, n-3}$   
 Significativo = entra

PASO 3 se determina si el efecto de  $\beta_{j1}$  es significativo en la presencia de  $\beta_{j2}$

Eliminación  $H_0: \beta_{j1} = 0 \mid \beta_{j2}$   $F_0 = \frac{SSR(\beta_{j1} \mid \beta_0, \beta_{j2})}{MSE(FM)} \sim F_{1, n-3}$   
 $H_a: \beta_{j1} \neq 0 \mid \beta_{j2}$   
 ↳  $\beta_0, \beta_{j1}, \beta_{j2}$

SI ES SIGNIFICATIVA NO SALE

PASO 4 Igual a paso 3 en forward: Prueba sobre  $\beta_{j3}$   
 FM:  $\beta_{j3}, \beta_{j2}, \beta_{j1}, \beta_0$   
 $F_{1, n-4}$   
 Significativa = entra

Paso 5. Se analiza la significancia de los parámetros  $\beta_1$  y  $\beta_2$  dado que en el modelo está  $\beta_3$ , se hace Eliminación Individual y para los dos al tiempo.

- Prueba 1.

Siempre son modelos con intercepto

$$H_0: \beta_1 = 0 \mid \beta_2, \beta_3 \quad F_0 = \frac{SSR(\beta_1 \mid \beta_2, \beta_3)}{MSE(FM)} \sim F_{1, n-4}$$

$$H_a: \beta_1 \neq 0 \mid \beta_2, \beta_3$$

↳  $\beta_0, \beta_1, \beta_2, \beta_3$

SI ES SIGNIFICATIVA  
NO SALE

- Prueba 2.

$$H_0: \beta_2 = 0 \mid \beta_1, \beta_3 \quad F_0 = \frac{SSR(\beta_2 \mid \beta_1, \beta_3)}{MSE(FM)} \sim F_{1, n-4}$$

$$H_a: \beta_2 \neq 0 \mid \beta_1, \beta_3$$

↳  $\beta_0, \beta_1, \beta_2, \beta_3$

SI ES SIGNIFICATIVA  
NO SALE

- Prueba 3.

$$H_0: \beta_1 = 0 \wedge \beta_2 = 0 \mid \beta_3 \quad F_0 = \frac{SSR(\beta_1, \beta_2 \mid \beta_3)}{MSE(FM)} \sim F_{2, n-4}$$

$$H_a: \beta_1 \neq 0 \vee \beta_2 \neq 0 \mid \beta_3$$

↳  $\beta_0, \beta_1, \beta_2, \beta_3$

SI SON SIGNIFICATIVAS  
NO SALEN

Si se corroborara la significancia de los dos parámetros dado que entro  $\beta_3$  entonces no elimino ninguna de las dos variables y sigo con el procedimiento [Entro, Elimino, Entro...]

REGRESION CON VARIABLES INDICADORAS

Existen problemas que además de tener un conjunto de variables  $(x_1, x_2, x_3, \dots, x_k)$  tienen diferentes niveles de una variable cualitativa o categorica, por ejemplo:

- Sexo: Hombre  
Mujer

- Nacionalidad:

- Turno: A = Mañana  
B = Tarde  
C = Noche

- Tipo de Maquina

- Estrato =  $E_1, E_2, E_3, E_4, E_5, E_6$

- Tipo de Monta

Otras regresiones  
- Modelo lineal generalizado  
- Logística  
a = Categorías  
a-1 = Indicadoras

Cada una de las variables se incorporan usando variables indicadoras

$$I = \begin{cases} 1 & \text{Categoría 1} \\ 0 & \text{e.o.p} \end{cases}$$

Para variables cualitativas con dos niveles tenemos

$$I = \begin{cases} 1 & \text{Categoría 1} \\ 0 & \text{Categoría 2} \rightarrow \text{Categoría de referencia} \end{cases}$$

En el cual solo es necesario definir una variable indicadora, tomando una categoría de referencia con la que se va a comparar

El modelo para este caso sería de la forma

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 I + \beta_3 I x_{i1} + \epsilon$$

con  $I=1$ ; categoría 1, Modelo =

$$y_i = (\beta_0 + \beta_2) + (\beta_1 + \beta_3) x_{i1} + \epsilon$$

Donde

$\beta_2$  = Mide el cambio en el intercepto con la categoría 1

$\beta_3$  = Mide el cambio en la pendiente con la categoría 1

Con  $I=0$ ; Categoría 2, Modelo =

$$y_i = \beta_0 + \beta_1 x_{i1} + \epsilon$$

Hallando los valores esperados

- categoría 1

$$E[y_i | x] = (\beta_0 + \beta_2) + (\beta_1 + \beta_3) x$$

- categoría 2

$$E[y_i | x] = \beta_0 + \beta_1 x$$

## PREGUNTAS

dependencia de las categorías en la relación  $y$  vs  $x$

- Es necesaria la variable categórica para explicar y con  $x$

↳ Igualdad rectas de regresión

$$H_0: \beta_0 = \beta_0 + \beta_2 \quad \wedge \quad \beta_1 = \beta_1 + \beta_3 \quad = \quad \beta_2 = 0 \quad \wedge \quad \beta_3 = 0$$

$$H_a: \beta_0 \neq \beta_0 + \beta_2 \quad \vee \quad \beta_1 \neq \beta_1 + \beta_3 \quad = \quad \beta_2 \neq 0 \quad \vee \quad \beta_3 \neq 0$$

- Los pendientes de ambas rectas son iguales

$$H_0: \beta_1 = \beta_1 + \beta_3 \quad = \quad \beta_3 = 0$$

$$H_a: \beta_1 \neq \beta_1 + \beta_3 \quad = \quad \beta_3 \neq 0$$

↳ Efecto de  $x$  sobre  $y$  para las categorías

- Los pendientes de ambas rectas son iguales y no significativas

$$H_0: \beta_1 = \beta_1 + \beta_3 = 0 \quad = \quad \beta_1 = 0 \quad \wedge \quad \beta_3 = 0$$

$$H_a: \beta_1 \neq 0 \quad \vee \quad \beta_1 + \beta_3 \neq 0 \quad \vee \quad \beta_1 \neq \beta_1 + \beta_3 \quad = \quad \beta_1 \neq 0 \quad \vee \quad \beta_3 \neq 0$$

Para variables cualitativas con tres niveles tenemos

Es necesario definir 2 variables indicadoras

$$I_1 = \begin{cases} 1 & \text{categoría 1} \\ 0 & \text{e.o.p} \end{cases} \quad I_2 = \begin{cases} 1 & \text{categoría 2} \\ 0 & \text{e.o.p} \end{cases}$$

La categoría de referencia es la 3 y se obtiene cuando  $I_1 = 0$   
 $\wedge I_2 = 0$ , esto es lo visualizamos de la siguiente forma:

$I_1$	$I_2$	
1	0	Categoría 1
0	1	Categoría 2
0	0	Categoría 3

Modelo de segundo orden: con una variable regresora

$$y = \beta_0 + \beta_1 x + \beta_2 I_1 + \beta_3 I_2 + \beta_4 I_1 x + \beta_5 I_2 x + \epsilon$$

## MODELOS

- Categoría 1:  $I_1 = 1 \quad I_2 = 0$

$$y = (\beta_0 + \beta_2) + (\beta_1 + \beta_4) x + \epsilon$$

- Categoría 2:  $I_1=0$   $I_2=1$

$$y = (\beta_0 + \beta_3) + (\beta_1 + \beta_5)x + E$$

- Categoría 3:  $I_1=0$   $I_2=0$

$$y = \beta_0 + \beta_1 x + E$$

## PREGUNTAS

- Igualdad de las rectas de regresión: Necesidad de incorporar la variable categórica en el modelo de regresión.

$$H_0: \beta_0 = \beta_0 + \beta_3 = \beta_0 + \beta_2 \quad \wedge \quad \beta_1 = \beta_1 + \beta_4 = \beta_1 + \beta_5$$



$$H_0: \beta_3 = 0 \quad \wedge \quad \beta_2 = 0 \quad \wedge \quad \beta_4 = 0 \quad \wedge \quad \beta_5 = 0$$

$$H_a: \beta_3 \neq 0 \quad \vee \quad \beta_2 \neq 0 \quad \vee \quad \beta_4 \neq 0 \quad \vee \quad \beta_5 \neq 0$$

- Las pendientes de las rectas son iguales

con esta prueba se puede analizar el efecto de  $x$  sobre  $y$  de acuerdo a las diferentes combinaciones de las categorías

$$H_0: \beta_1 = \beta_1 + \beta_4 = \beta_1 + \beta_5 \quad \equiv \quad \beta_4 = 0 \quad \wedge \quad \beta_5 = 0$$

$$H_a: \beta_1 \neq \beta_1 + \beta_4 \neq \beta_1 + \beta_5 \quad \equiv \quad \beta_4 \neq 0 \quad \vee \quad \beta_5 \neq 0$$

- Las pendientes de las rectas son iguales y no significativas

$$H_0: \beta_1 = \beta_1 + \beta_4 = \beta_1 + \beta_5 = 0$$

$$H_a: \beta_1 \neq 0 \quad \vee \quad \beta_1 + \beta_4 \neq 0 \quad \vee \quad \beta_1 + \beta_5 \neq 0 \quad \vee \quad \beta_1 \neq \beta_1 + \beta_4 \quad \vee \quad \beta_1 \neq \beta_1 + \beta_5 \quad \vee \dots$$



$$H_0: \beta_1 = 0 \quad \wedge \quad \beta_4 = 0 \quad \wedge \quad \beta_5 = 0$$

$$H_a: \beta_1 \neq 0 \quad \vee \quad \beta_4 \neq 0 \quad \vee \quad \beta_5 \neq 0$$