

# Fundamentos de Estadística



*Silvia N. Pérez*  
*Especialización en Ciencia de*  
*Datos - UNO*

# Pruebas de comparación

POBLACIÓN DE ESTUDIO	CONDICION	DATOS CUANTITATIVOS <i>pruebas : cuanti vs cuanti</i>		DATOS <u>CUALITATIVOS</u> <i>cuanti vs cuanti</i>
		PRUEBA PARAMETRICA	PRUEBA NO PARAMETRICA	PRUEBA NO PARAMETRICA
DOS GRUPOS	INDEPENDIENTES	T-student para muestras indepenientes	U- Mann Whitney <i>(Wilcoxon-MW)</i>	$\chi^2$ (Independencia)
	APAREADOS	T-student para muestras relacionadas <i>(apareados)</i>	Wilcoxon <i>(apareados)</i>	Mc. Nemar (dicotómicas)
MÁS DE DOS GRUPOS	INDEPENDIENTES	ANOVA (Análisis de varianza)	Kruskal Wallis	$\chi^2$

# Análisis de Varianza (ANOVA)

# Qué datos tenemos?

*Quiero analizar*  
*cuanti ~ cuali*

- Un factor (tratamientos) **Cualitativa**
- Una variable de respuesta **Cuantitativa**
- Pregunta de interés:

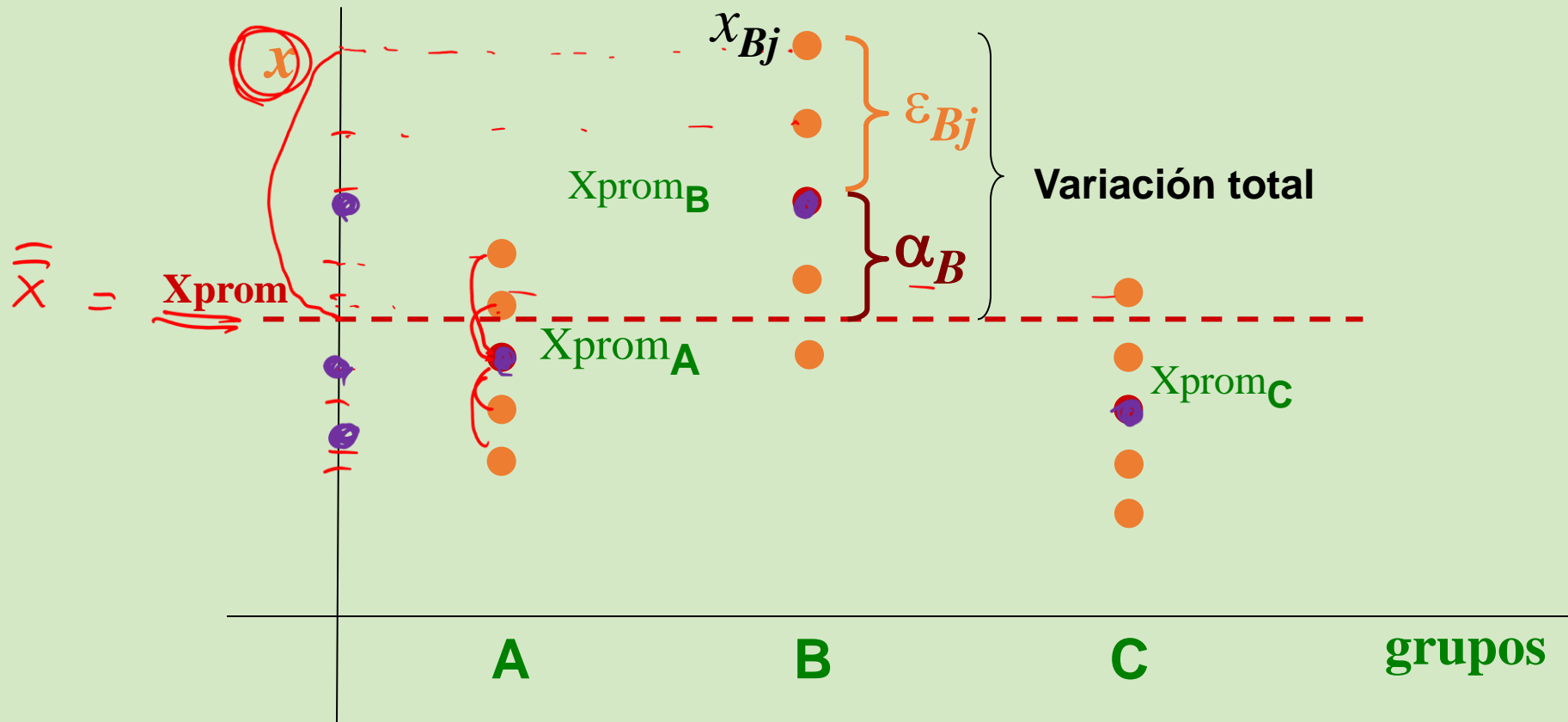
¿Las medias de cada grupo difieren o están “afectadas” por la variable cuanti??

- Número de grupos:
  - Si son dos grupos: utilizamos test de t
  - más de 2 grupos: comparaciones múltiples con ANOVA

# ANOVA (ANalysis Of Variance)

Modelo lineal de efectos fijos

$$x_{ij} = \mu + \alpha_i + \varepsilon_{ij}$$



$$\sum_{ij} (x_{ij} - \bar{\bar{x}})^2 = \sum_{ij} (\bar{x}_i - \bar{\bar{x}})^2 + \sum_{ij} (x_{ij} - \bar{x}_i)^2$$

*promedio de grupo i*
*promedio general*

**SCT TOTAL**

*dispersión  
s/ discriminar  
s/ considerar grupos*

SCF ENTRE  
grupos (a)

SCE DENTRO de  
grupos (residual)

$$CM_{entre} = \frac{SC_{entre}}{gl_{entre}} = \frac{\sum_{ij} (\bar{x}_{ij} - \bar{\bar{x}})^2}{a - 1}$$

Mira la dispersión de los  
promedios de cada grupo

$$CM_{dentro \text{ o } residual} = \frac{SC_{dentro}}{gl_{dentro}} = \frac{\sum_{ij} (x_{ij} - \bar{x}_i)^2}{n - a}$$

Mira la dispersión dentro  
de cada grupo

$$F = \text{CM\_entre} / \text{CM\_dentro}$$

tiene distribución  $\mathcal{F}$  con  $(a-1)$  y  $(n-a)$  grados de libertad

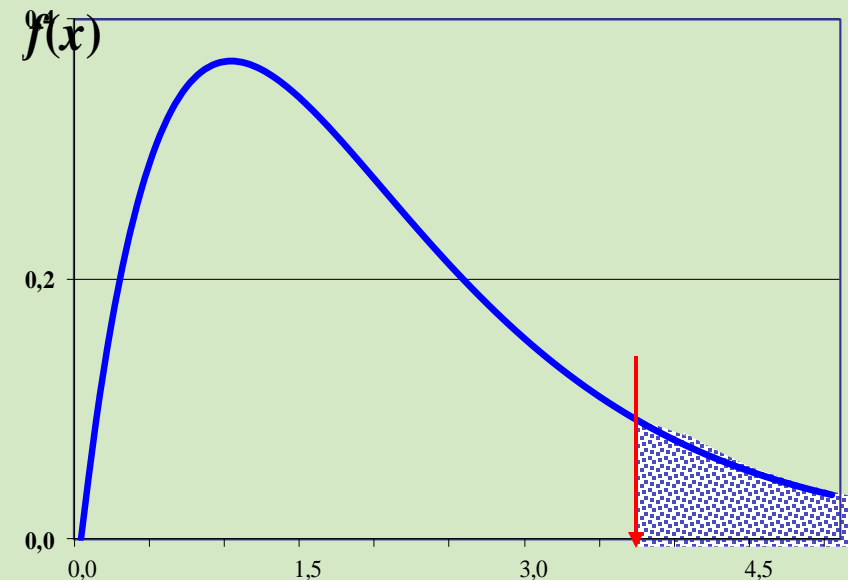
### Supuestos para la validez del test

- Normalidad
- Homocedasticidad (= varianzas)
- Independencia de las observaciones

## Tabla ANOVA

FUENTE DE VARIACION	SUMA DE CUADRADOS	GL	CUADRADOS MEDIOS	F
ENTRE GRUPOS	SCF entre	a-1	SC entre/(a - 1)	$\frac{CM \text{ entre}}{CM \text{ dentro}}$
DENTRO DE GRUPOS	SCE dentro	n-a	SC dentro/(n - a)	
TOTAL	SCT total	n-1		

El  $F_{\text{calculado}}$  se compara con  
el  $F_{\text{tabulado}}$  con (a-1) y (n-a) GL





Esto es  $n_i$  y a la variable cuanti

y  $G$ : variable cuali

Si  $G = \{a, b\}$  (2 valores posibles)  $\Rightarrow$  el Test t prueba:

$$H_0: \mu_Y^a = \mu_Y^b \quad H_1: \neq$$

Si  $G = \{a, b, c, \dots\}$  (+2 valores)  $\Rightarrow$  ANOVA prueba

$$H_0: \mu_Y^a = \mu_Y^b = \mu_Y^c = \dots$$

$H_1$ : "no son todos iguales"  
o sea hay alguna  $\neq$

Ej:  $Y = \text{fonaje}$  y  $G = \{\text{trat}_1, \text{trat}_2, \text{trat}_3\}$

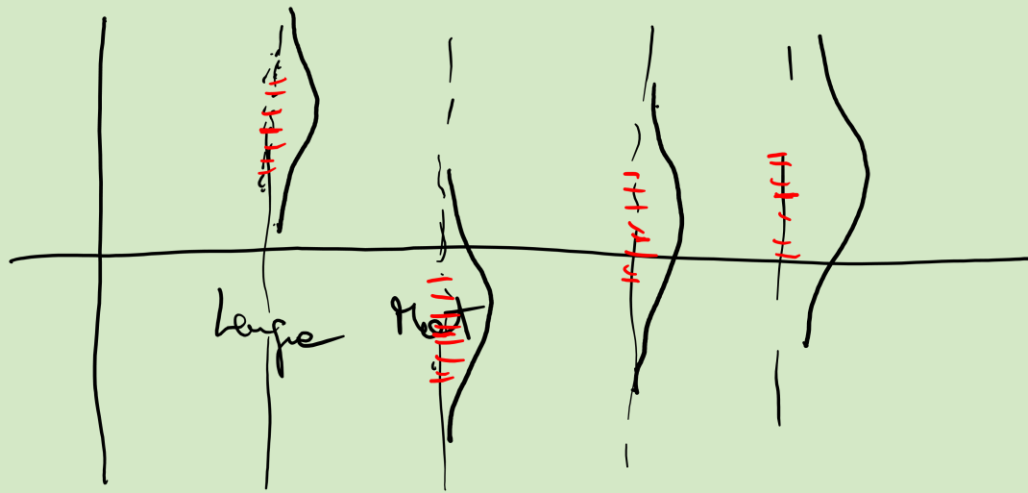
$$H_0: \mu_Y^1 = \mu_Y^2 = \mu_Y^3 \quad H_1: \text{alguna} \neq$$

En los datos puede que tenga

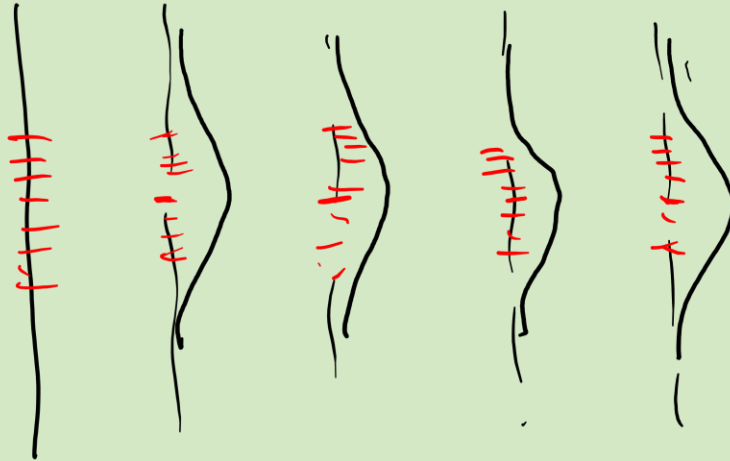
$$\text{mean}(Y) = \bar{X}_1 = 101.7$$

$$\bar{X}_2 = 99.98$$

$$\bar{X}_3 = 100.2$$



medias so  $\neq$



medias son =

# Ejemplo en R

Se quiere comparar las medias de notas de estudiantes según el área de interés que declararon.

Los datos corresponden a `estudiantes_completa.xlsx`

*POSTHOC*  
Comparaciones a posteriori: Esto es, si ANOVA indica  $\text{Rechazo } H_0 \Rightarrow$   
 $\Rightarrow$  hacemos estos pruebas p/decidir entre pares de medias si hay  $\neq$

- Tuckey: si hay muchas categorías de las cuales.
- Bonferroni: Se usa si son pocas las categorías de las cuales.

Se recomienda hacer varios test a posteriori, y decidir mirando todo.

# Tablas de contingencia

# Tablas de contingencia

- Analiza dos variables cualitativas, X e Y, buscando asociación entre ellas.
- Permite probar la hipótesis de independencia entre ellas, esto es:

$H_0$ : X e Y **independientes**

$H_1$ : dependientes, esto es, **relacionados**

# Tablas de contingencia

**Tabla de frecuencias observadas**

<b>X // Y</b>	<b>y1</b>	<b>y2</b>	<b>...</b>	<b>Frec. Marginal de X</b>
<b>x1</b>				
<b>x2</b>				
<b>...</b>			<b><math>n_{ij}</math></b>	
<b>Frec. Marginal de Y</b>				<b><math>n</math></b>

Se comparan las frecuencias observadas en cada casilla versus frecuencias esperadas bajo la hipótesis de independencia

# Ejemplo:

Se consulta sobre acuerdo con la construcción de Metrobus y nivel de educación.

Hay relación?

**Tabla de frecuencias observadas**

	Muy /bastante	Poco / nada	Total
Sin estudios / primaria	1130	1123	2253
Secundaria	860	599	1459
Universitarios	480	340	820
Total	2470	2062	4532

# Tablas de contingencia

frecuencias observadas (en %)

	Muy /bastante	Poco / nada	Total
Sin estudios / primaria	50.2%	49.8%	100%
Secundaria	58.9%	41.1%	100%
Universitarios	58.5%	41.5%	100%
Total	54.5%	45.5%	100%

¿Se podría hablar de **que existe asociación** entre nivel de estudios y acuerdo con Metrobus?

Hipótesis:

$H_0$ : hay indep entre 'acuerdo' y 'nivel de estudio'

$H_1$ : no indep (o no, hay alguna relación o asociación)

Para explorar si existe asociación se compara esta tabla con una tabla en la que no existiría asociación o también llamada **tabla de frecuencias esperadas** (suponiendo independencia)



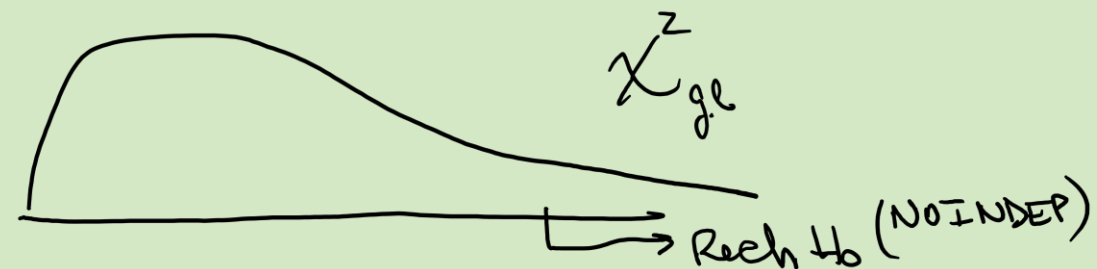
# Tablas de contingencia

Las hipótesis

- $H_0$ : independencia (no asociación)
- $H_1$ : existe asociación

Estadístico de prueba para comprobar la hipótesis:

$$\chi^2 = \sum \frac{(O - E)^2}{E} \sim \chi^2_{\underbrace{(r-1)(s-1)}_{d.f.}} \text{ si } H_0 \text{ cierta}$$



# Tablas de contingencia

Estadístico de prueba:

$$\chi^2 = \sum \frac{(O - E)^2}{E} \sim \chi^2_{(r-1)(s-1)} \text{ si } H_0 \text{ cierta}$$

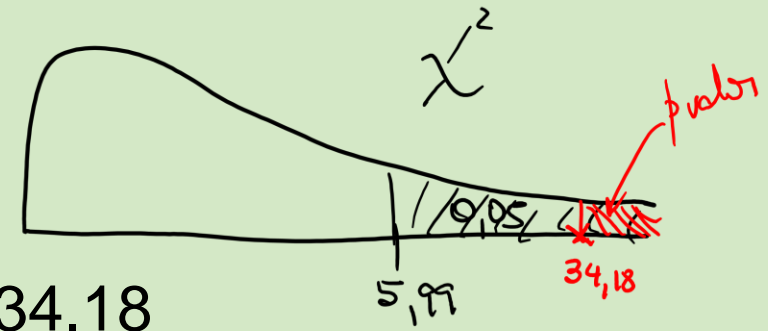
Con nuestros datos, se calculan las frec esperadas:

OBSERVADAS	Muy /bastante	Poco / nada	Total
Sin estudios / primaria	1130	1123	2253
Secundaria	860	599	1459
Universitarios	480	340	820
Total	2470	2062	4532

ESPERADAS	Muy /bastante	Poco / nada	Total
Sin estudios / primaria	1127.9	1025.1	2253
Secundaria	795.2	663.8	1459
Universitarios	446.9	373.1	820
Total	2470	2062	4532

Se calcula el valor del estadístico de prueba, que aquí da **34.18**

# Tablas de contingencia



- el valor obtenido en la prueba de chi cuadrado es 34.18
- es mayor que el valor de chi-cuadrado obtenido en la tabla para un nivel de significación de 0.05 con dos grados de libertad (5.99).
- Es decir, **34.18** pertenece a la Región de Rechazo,
- $\chi^2$  observado  $> 5.99$  por lo que rechazamos la  $H_0$  (hipótesis nula)  
*rechazo la independencia*
- **Esto permite decir que existe asociación entre nivel educativo y acuerdo con Metrobus,** con un nivel de significación del 5%.

# Tablas de contingencia

**¡Ojo!** Para poder realizar la prueba ji-cuadrado se debe cumplir un supuesto:

Como máximo un 20% de las celdas debe tener una frecuencia esperada menor a 5.

(no debe haber demasiadas celdas con pocos casos).

# Medidas de asociación

Observación:

- La prueba de contingencia no indica la dirección de la asociación
- No informa sobre el grado de asociación

## Coeficiente de Contingencia

$$C = \sqrt{\frac{\chi^2}{\chi^2 + n}}$$

## V de Cramer

$$V = \sqrt{\frac{\chi^2}{n \cdot t}}$$

- $t = \min\{\text{\#filas}-1, \text{\#col}-1\}$
- V de Cramer entre 0 y 1.
  - <0.05: asociación débil
  - 0.05-0.25 moderado
  - 0.25-1 fuerte

# Medidas de asociación

$$\chi^2 = 34.18$$
$$n = 4532$$

En los datos anteriores:

Coeficiente de Contingencia = **0.08652**

V de Cramer = **0.08684**

CONCLUSIÓN: La asociación detectada entre nivel de estudios y acuerdo, aunque significativa, **no es fuerte** (índices cercanos a 0).

## Ejemplo en R

Se quiere decidir si hay independencia entre el sexo del estudiante y el interés en un área en particular.

Los datos corresponden a `estudiantes_completa.xlsx`