

EDEM



Máster en Data Analytics

Estadística con Python IV
Miguel Rua del Barrio

Comparación de múltiples medias

En la unidad anterior hemos visto test para comparar hasta 2 medias de dos grupos distintos.

A partir de ahora vamos a trabajar el caso de cuando tenemos más de dos grupos que comparar.

Ejemplo: Tiempo medio de recuperación del paciente

- Según la enfermedad que padezcan (cáncer, bronquitis...)
- Por edad (mayor de 65, adulto, niño)

No podremos realizar estos análisis siempre que queramos, habrá unas condiciones de aplicabilidad:

- 1. Poblaciones normalmente distribuidas en cada grupo.**
- 2. Homocedasticidad de varianzas entre grupos.**
- 3. Muestra aleatoria e independientes.**

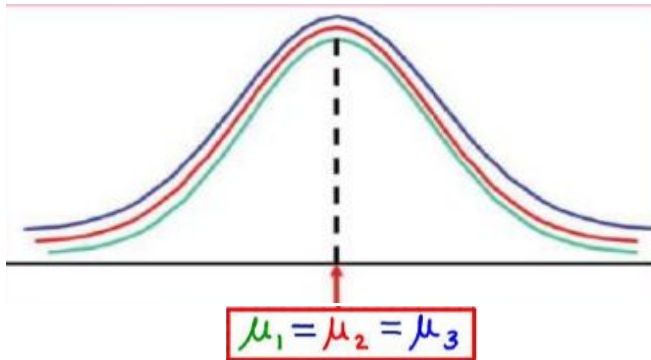
Nos podremos ayudar de las técnicas vistas anteriormente para comprobar las condiciones de aplicabilidad.

Comparación de múltiples medias

Hipótesis nula

- Todas las medias son iguales

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

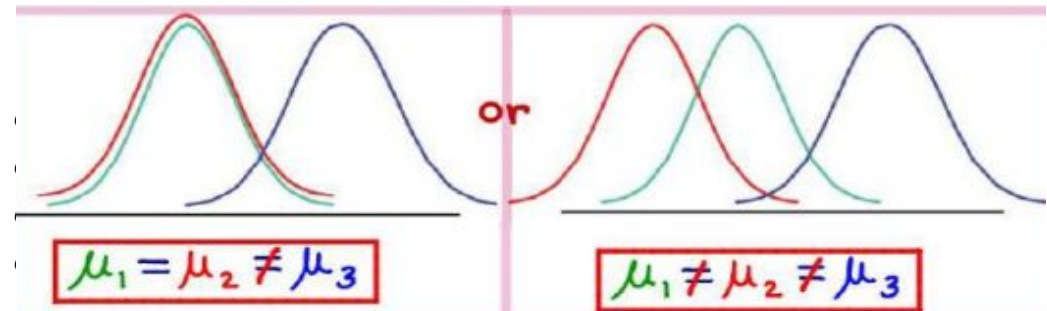


- No hay diferencia entre las medias de los grupos

Hipótesis alternativa

- Al menos una media es diferente

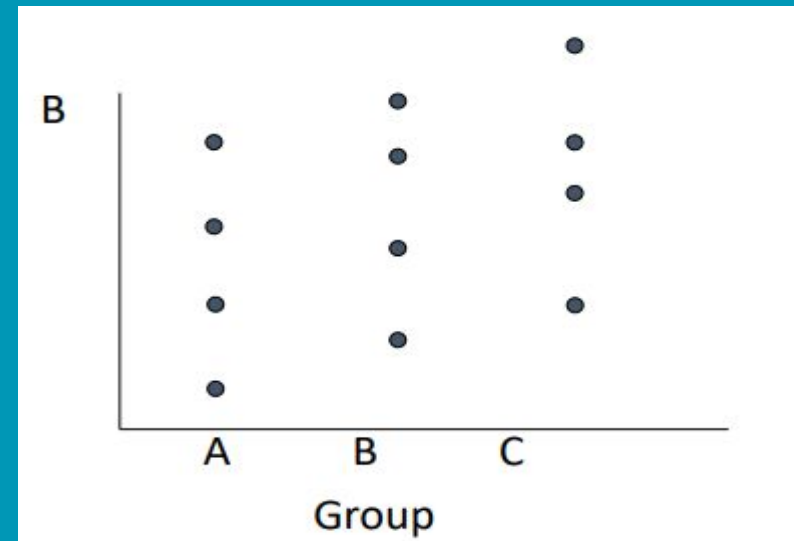
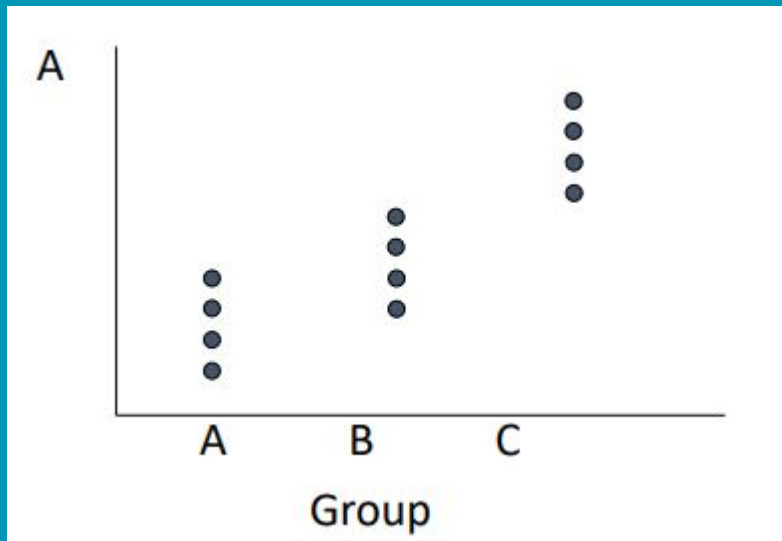
$$H_1: \mu_i \neq \mu_j \text{ for at least one } i, j \text{ pair}$$



- Algunos pares de medias pueden ser iguales

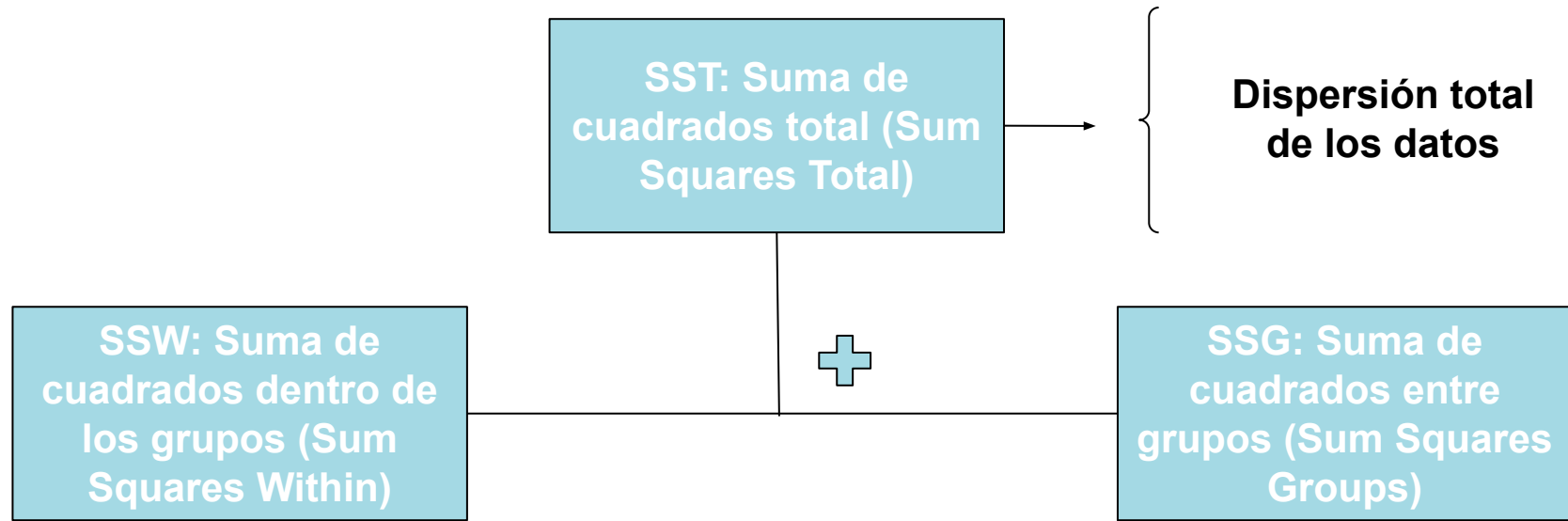
Variabilidad de los datos

- La variabilidad es un factor clave para comprobar la igualdad de medias.
- En cada caso, las medias parecen diferentes, pero una gran variabilidad en el grupo B hace que no sea tan significativa esta diferencia.

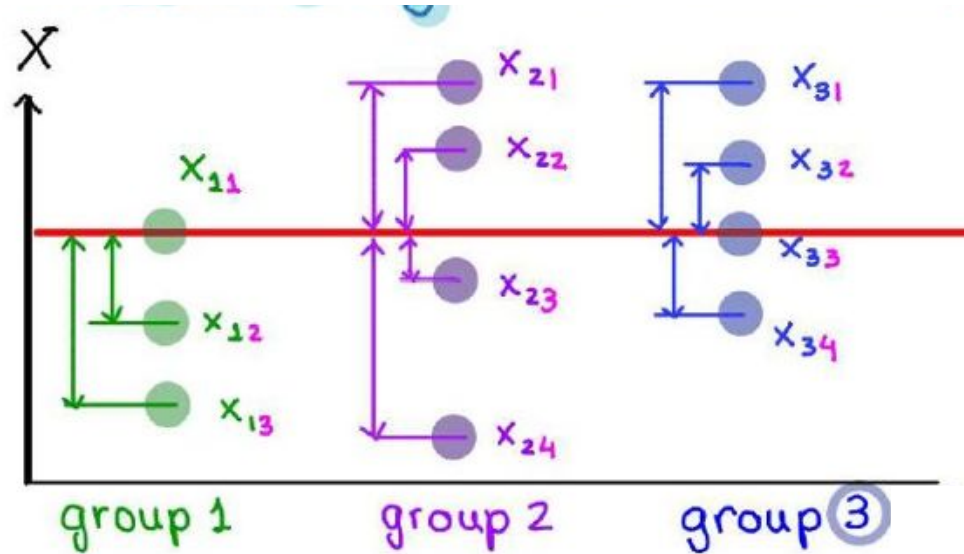


Comparación de múltiples medias

- Enfoque: Estudiar variabilidad de los datos para realizar el test



Variabilidad SST



\bar{x} = media total de los datos

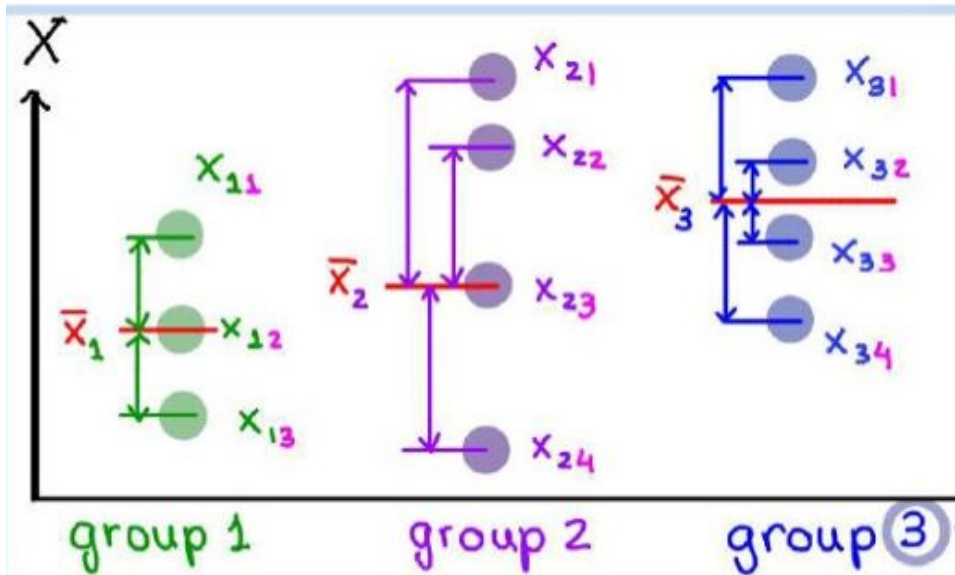
n_i = número de observaciones en el grupo i

k = número de grupos

$$SST = \left[\begin{array}{l} (x_{11} - \bar{x})^2 + (x_{12} - \bar{x})^2 + (x_{13} - \bar{x})^2 \\ + (x_{21} - \bar{x})^2 + \dots + (x_{24} - \bar{x})^2 \\ + (x_{31} - \bar{x})^2 + \dots + (x_{34} - \bar{x})^2 \end{array} \right] = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2$$

jth observation from group i

Variabilidad SSW



\bar{x}_i = media total de los datos

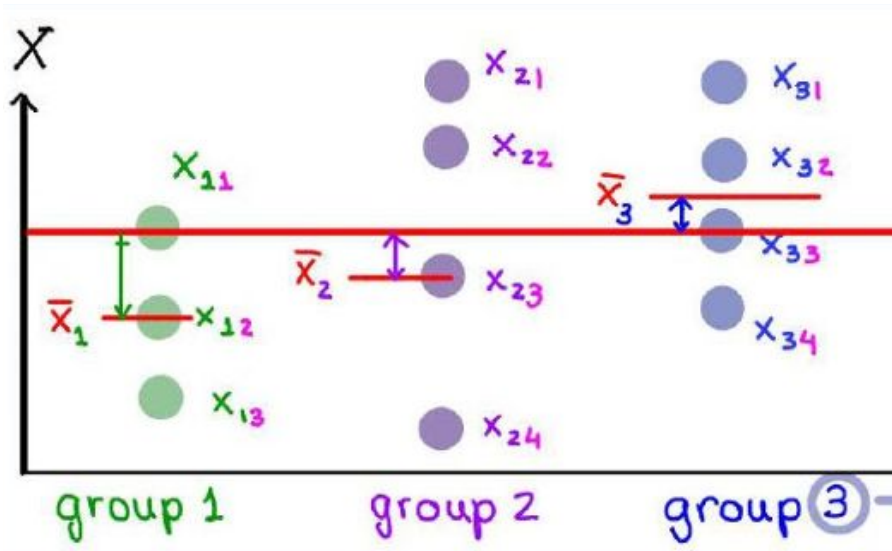
n_i = número de observaciones en el grupo i

k = número de grupos

$$SSW = \left[\begin{aligned} &(x_{11} - \bar{x}_1)^2 + (x_{12} - \bar{x}_1)^2 + (x_{13} - \bar{x}_1)^2 \\ &+ (x_{21} - \bar{x}_2)^2 + \dots + (x_{24} - \bar{x}_2)^2 \\ &+ (x_{31} - \bar{x}_3)^2 + \dots + (x_{34} - \bar{x}_3)^2 \end{aligned} \right] = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$$

jth observation from group i

Variabilidad SSG



\bar{x}_i = media de cada grupo

\bar{x} = media total de los datos

n_i = número de observaciones en el grupo i

k = número de grupos

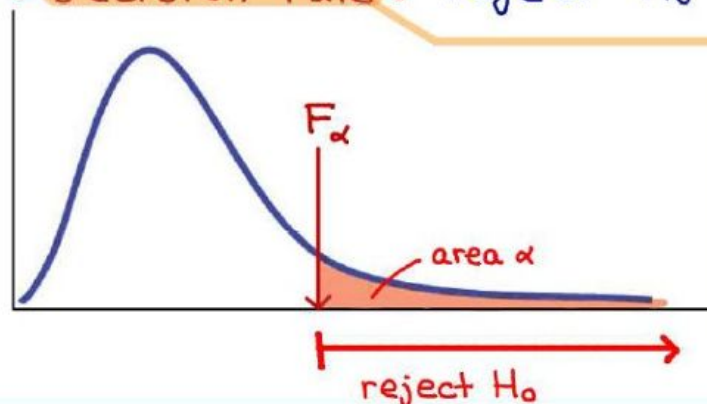
$$SSG = n_1 \cdot (\bar{x}_1 - \bar{x})^2 + n_2 \cdot (\bar{x}_2 - \bar{x})^2 + n_3 \cdot (\bar{x}_3 - \bar{x})^2 = \sum_{i=1}^k n_i \cdot (\bar{x}_i - \bar{x})^2$$

Test ANOVA

- Hipótesis $\begin{cases} H_0: \mu_1 = \mu_2 = \dots = \mu_k \\ H_1: \text{at least 2 } \mu_i\text{'s are } \neq \end{cases}$
- Estadístico del test $F = \frac{MSG}{MSW} = \frac{\text{between groups var}/(k-1)}{\text{within groups var}/(n-k)}$

$$\begin{cases} df_1 = k-1 & (\text{typically small}) \\ df_2 = n-k & (\text{" large}) \end{cases}$$

* Decision rule: reject H_0 if $F > F_{\text{critical}} = F_{k-1, n-k, \alpha}$



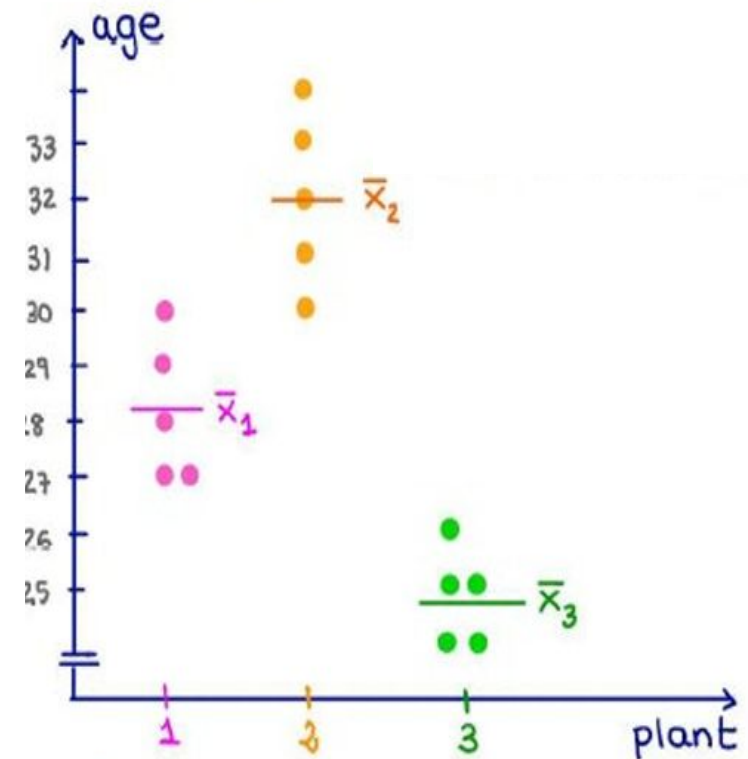
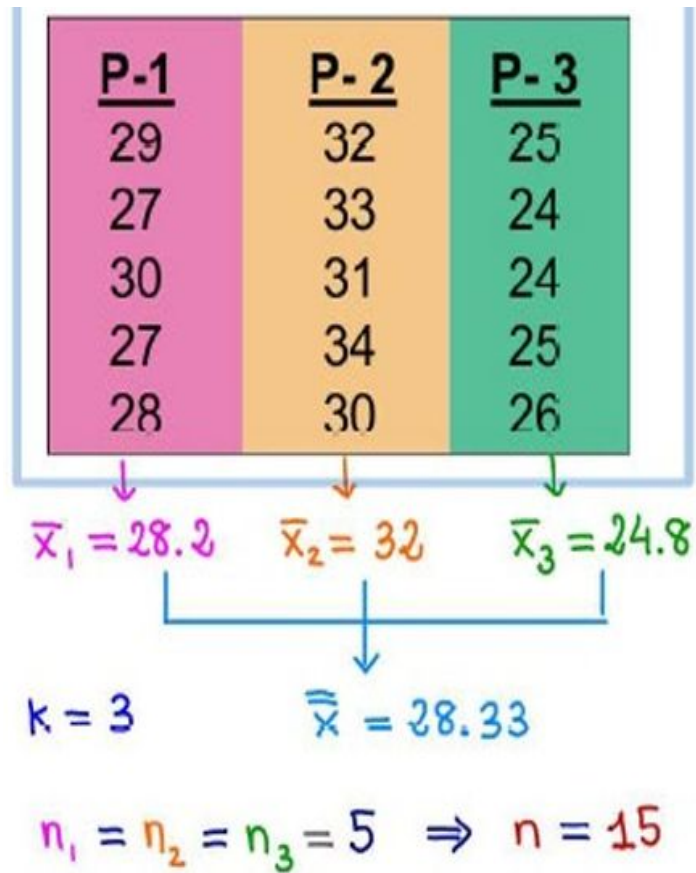
ANOVA de un factor

- Se suele utilizar la tabla ANOVA para representar el contraste

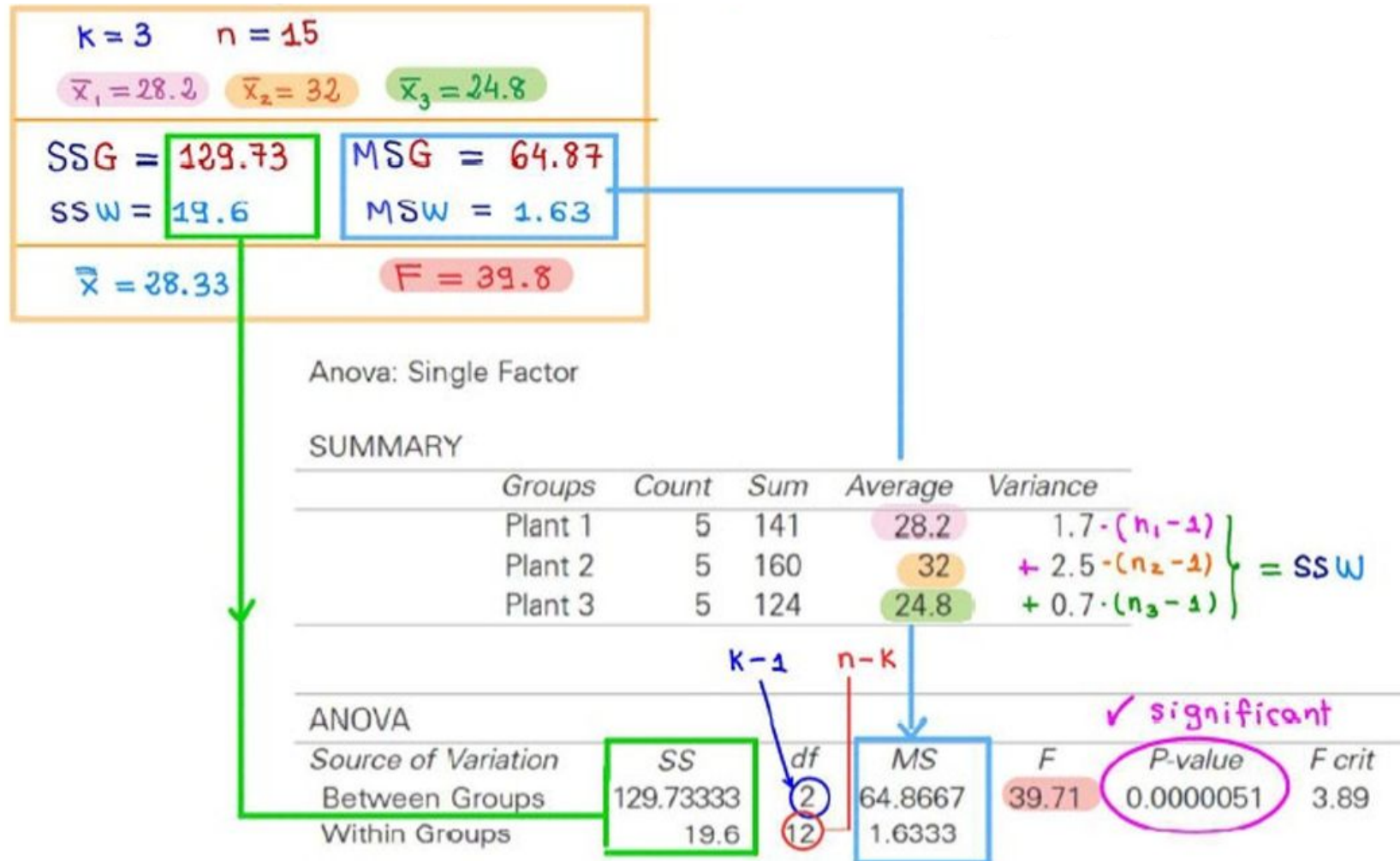
Source of variation	SS	df	MS (variance)	F-ratio
Between groups	SSG	$k - 1$	$MSG = \frac{SSG}{k - 1}$	$F = \frac{MSG}{MSW}$
Within groups	SSW	$n - k$	$MSW = \frac{SSW}{n - k}$	
Total	SST	$n - 1$		

ANOVA de un factor: ejemplo

- Se suele utilizar la tabla ANOVA para representar el contraste



ANOVA de un factor: ejemplo



ANOVA de un factor: ejemplo

* **Hypotheses**: $\left\{ \begin{array}{l} H_0 = \mu_1 = \mu_2 = \mu_3 \\ H_1 = \mu_i \text{ not all } = \end{array} \right\}$
 $\alpha = 0.05$

* **Test statistic**: $F = 39.71$

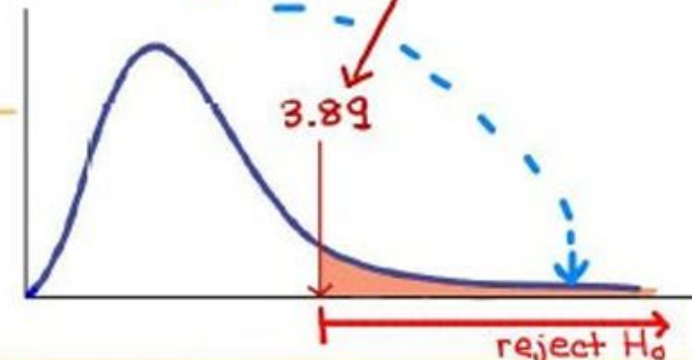
ANOVA				
Source of Variation	SS	df	MS	F
Between Groups	129.73333	2	64.8667	39.71
Within Groups	19.6	12	1.6333	

p-value for the
F-test

P-value
0.0000051

F crit
3.89

* **Decision rule**: reject H_0
if $F > F_{2, 12, 0.05} = 3.89$



Decisión: como $F \in$ a la región de rechazo
(o p valor $< \alpha$), rechazamos H_0

Conclusión: existen diferencias significativas en la
edad media de los tres grupos

Parte práctica de la sesión:

1. Realizar ANOVA con Python y test ad hoc

Durante este módulo hemos visto:

- **Nociones básicas de estadística como su nomenclatura (población, muestra..).**
- **Definición de las variables aleatorias y sus diferentes tipos.**
- **Descriptiva univariante y bivalente para cada tiempo de variable (gráficos y estadísticos)**
- **Distribuciones de probabilidad (definiciones, ejemplos)**
- **Contrastes de hipótesis (definición, concepto, ejemplos).**
- **ANOVA.**
- **Implementación de todos estos contenidos con Python.**
- **Realización de informes estadísticos con Google Colab.**

¡GRACIAS POR VUESTRA ATENCIÓN!



miguel.ruadelbarrio@ams-europe.com



[linkedin.com/in/miguel-rua-del-barrio-5214661b5](https://www.linkedin.com/in/miguel-rua-del-barrio-5214661b5)

Valoración DOCENCIA - Estadística
con Python - Miguel Rua

