

Biometría



Comparación de dos poblaciones

Comparando dos poblaciones

Frecuentemente el investigador está interesado en comparar dos poblaciones:

- La densidad de cierto pasto es mayor en zonas pastoreadas que no pastoreadas
- El descenso de peso con la dieta A es menor que con la B
- El porcentaje de luxados con la prótesis A es mayor que con la prótesis B
- La longitud corporal de los machos de cierta especie es más variable que la de las hembras

Inferencia basada en dos muestras

Para llevar a cabo esta comparación, el investigador necesita tomar **muestras**

Las muestras pueden ser:

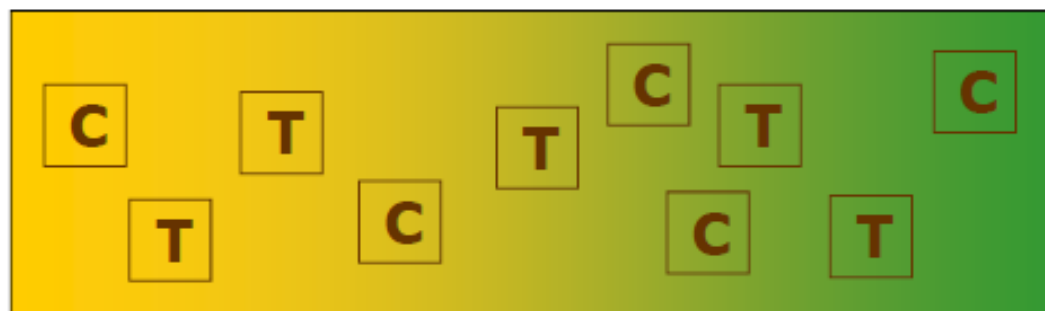
- Muestras independientes
 - 2 promedios
 - 2 proporciones
 - 2 varianzas o desvíos estándar
- Muestras dependientes o pareadas
 - diferencias

¿El pastoreo afecta la diversidad florística de los pastizales de la pampa de Achala?



Se desea estudiar el efecto del pastoreo y de su exclusión por 10 años sobre la composición y diversidad florística de un pastizal en las Sierras de Córdoba (2200 msnm)

Diseño 1: Se tomaron 10 parcelas de 20x20 m que se dividieron al azar en dos grupos: uno fue sometido a pastoreo mientras que en el otro fue excluido el ganado



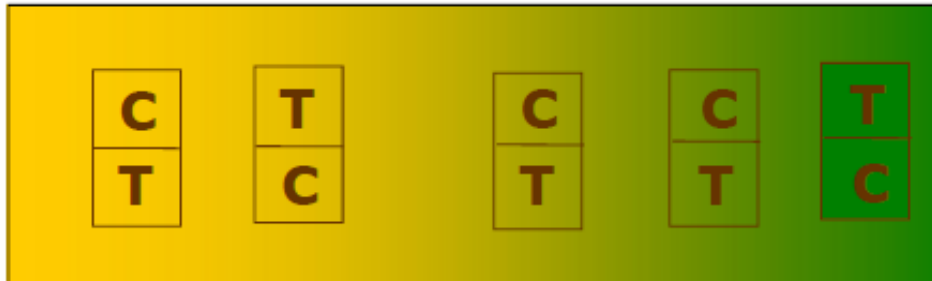
PASTOREADO	NO PASTOREADO

¿El pastoreo afecta la diversidad florística de los pastizales de la pampa de Achala?



Se desea estudiar el efecto del pastoreo y de su exclusión por 10 años sobre la composición y diversidad florística de un pastizal en las Sierras de Córdoba (2200 msnm)

- **Diseño 2:** Se tomaron 5 sectores. En cada uno se delimitaron dos parcelas adyacentes de 20 x 20 m c/u separadas por un alambrado, una pastoreada y otra excluida al ganado



Otro ejemplo



- Un investigador cree que los fumadores tienden a fumar más durante los períodos de stress.
- Para comprobarlo debe elegir entre dos metodologías:

Ejemplo



- Interroga a cada individuo con respecto a la cantidad de cigarrillos diarios fumados

m. indep

Sin stress Con stress

15	20
31	45
50	48
16	30
56	72



m. dep

Individuo Sin stress Con stress

1	15	→	20
2	31	→	45
3	50	→	48
4	16	→	30
5	56	→	72

Muestras independientes

- ✓ Los tratamientos son asignados al azar a las u.e.
- ✓ Cada observación en una muestra **no está relacionada** con ninguna observación en la otra muestra
- ✓ Cada individuo es observado **una sola vez**
- ✓ Las dos muestras pueden diferir en **varios** factores, no solo en el que interesa comparar
- ✓ Las dos muestras no necesariamente deben ser del mismo tamaño

Muestras dependientes o pareadas

- ✓ El investigador agrupa a las u.e. en pares y luego asigna al azar los tratamientos a las u.e. dentro de cada par
- ✓ Cada observación en una muestra está directamente **relacionada** con otra observación en la otra muestra
- ✓ Cada individuo es observado **dos veces**
- ✓ Las dos muestras difieren **solo** en el factor que interesa comparar
- ✓ Las dos muestras deben ser del mismo tamaño

Muestras dependientes vs independientes

- Se debe elegir entre dos preparaciones para el tratamiento de la dermatitis en codo.
- Se desea determinar si los carpinchos macho tienen un peso mayor al de las hembras
- Se desea determinar si una nueva droga para el tratamiento de úlceras es más efectiva que la que está actualmente en uso

Comparando dos proporciones

- En ciertos casos estamos interesados en comparar la proporción de “éxito” en dos poblaciones independientes.
- Para efectuar esta comparación se requiere

Una muestra aleatoria de tamaño n_1 extraída de la población 1 con parámetro π_1

Una muestra aleatoria de tamaño n_2 extraída de la población 2 con parámetro π_2

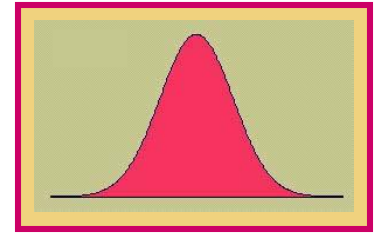
Comparando dos proporciones

Comparamos las dos proporciones haciendo inferencia sobre $\pi_1 - \pi_2$, la diferencia entre las dos proporciones poblacionales.

- Si las dos proporciones poblacionales son iguales, entonces $\pi_1 - \pi_2 = 0$.
- El mejor estimador de $\pi_1 - \pi_2$ es la diferencia entre las dos proporciones muestrales,

$$p_1 - p_2 = \frac{x_1}{n_1} - \frac{x_2}{n_2}$$

Distribución muestral de $p_1 - p_2$



1. La media de $p_1 - p_2$ es $\pi_1 - \pi_2$, la diferencia entre las proporciones poblacionales.
2. El desvío estándar (EE) de $p_1 - p_2$ es $\sqrt{\frac{p(1-p)}{n_1} + \frac{p(1-p)}{n_2}}$
donde $p = \frac{x_1 + x_2}{n_1 + n_2}$
3. Si el tamaño de las muestras es lo suficientemente grande, $pn > 5$ y $(1-p)n > 5$, entonces
$$z = \frac{(p_1 - p_2) - (\pi_1 - \pi_2)}{\sqrt{\frac{p(1-p)}{n_1} + \frac{p(1-p)}{n_2}}}$$
$$p_1 - p_2$$

sigue una distribución normal

¿Es efectiva la aspirina en la prevención de infartos?



- En 1982, 22000 hombres (todos médicos) de entre 40 y 84 años, sin antecedentes de cardiopatía o de accidente cerebrovascular, se sometieron a un estudio para evaluar la eficacia de la aspirina
- Se registró la presencia de infartos (incidencia) durante 5 años:

Grupo	Infarto	No infarto	n	Incidencia
Placebo	239	10795	11034	0,0217
Aspirina	139	10898	11037	0,0126

¿Es efectiva la aspirina en la prevención de infartos?



- H_0 :
- H_1 :
- CR:

Intervalo de confianza para la diferencia entre 2 proporciones

$$\hat{\theta} \pm VC EE_{(\hat{\theta})}$$

Si el tamaño de las muestras es lo suficientemente grande, $pn > 5$ y $(1-p)n > 5$,

Entonces $p_1 - p_2$ sigue una distribución normal

El desvío estándar (EE) de $\hat{p}_1 - \hat{p}_2$ es $\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$

$$\Delta p \pm Z_{1-\alpha/2} \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

Prevalencia de la diabetes mellitus no dependiente de la insulina en Lejona (Vizcaya)

TABLA 1

Características generales de los distintos grupos estudiados

	Normales (grupo N)		Diabéticos (grupo DM)	
	Varones	Mujeres	Varones	Mujeres
Número de individuos	370	347	15	16
PAS (mmHg)	125,7 ± 16	125,4 ± 21,1 ^a	148,6 ± 24	158,7 ± 29 ^a
PAD (mmHg)	71,3 ± 9,2	71,5 ± 11,8 ^a	82 ± 10,6	87,5 ± 14,4 ^a
Antecedentes familiares (%)	18,1	16,7 ^a	18,8	23,5 ^a

IMC: índice de masa corporal; PAS: presión arterial sistólica; PAD: presión arterial diastólica. Resultados expresados como $\bar{x} \pm DE$.

¿En pacientes diabéticos, existen diferencias en la variabilidad de la presión arterial diastólica entre hombres y mujeres?

Comparando la variabilidad de dos poblaciones: Prueba F

- $H_0: \sigma^2_1 = \sigma^2_2$

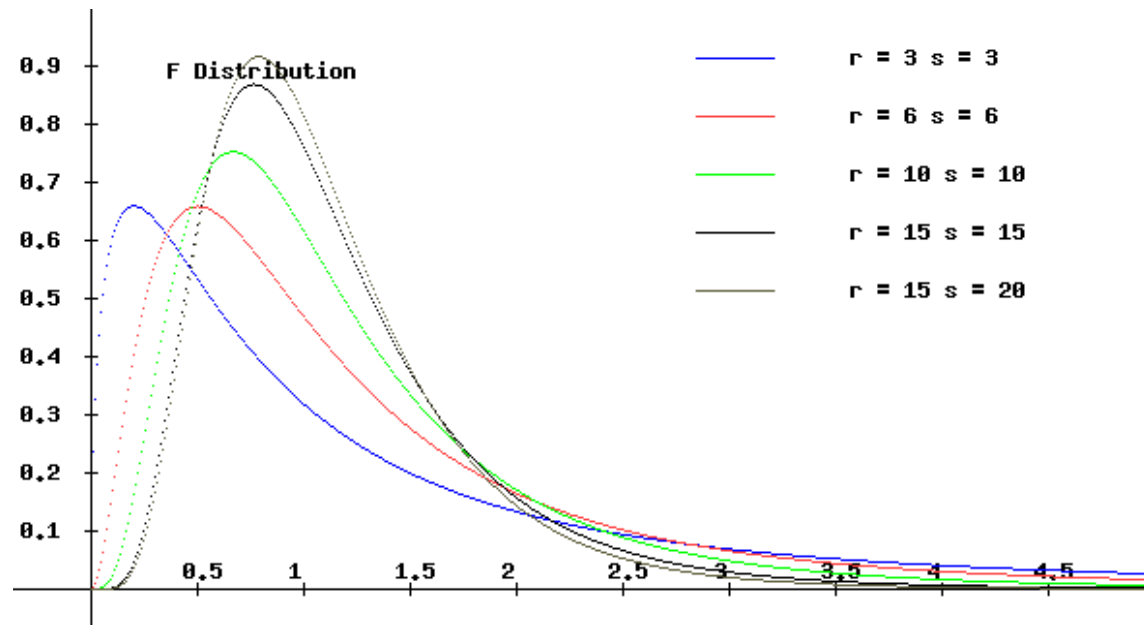
- $H_1: \sigma^2_1 \neq \sigma^2_2$

Distribución muestral de s^2_1 / s^2_2

1. El cociente s^2_1/s^2_2 estima el cociente σ^2_1 / σ^2_2
2. Si las poblaciones originales siguen una distribución **normal**, el cociente s^2_1/s^2_2 seguirá una distribución conocida como F de Fisher-Snedecor con GL_1 y GL_2 grados de libertad

Distribución F de Fisher-Snedecor

- es una familia de distribuciones
- Matemáticamente surge como el cociente de dos variables con distribución chi-cuadrado
- cada curva está caracterizada por dos GL
- la distribución es asimétrica positiva
- la variable solo toma valores positivos



¿En pacientes diabéticos, existen diferencias entre hombres y mujeres en los valores generales de presión arterial diastólica?

TABLA 1

Características generales de los distintos grupos estudiados

	Normales (grupo N)		Diabéticos (grupo DM)	
	Varones	Mujeres	Varones	Mujeres
Número de individuos	370	347	15	16
PAS (mmHg)	125,7 ± 16	125,4 ± 21,1 ^s	148,6 ± 24	158,7 ± 29 ^s
PAD (mmHg)	71,3 ± 9,2	71,5 ± 11,8 ^s	82 ± 10,6	87,5 ± 14,4 ^s
Antecedentes familiares (%)	18,1	16,7 ^s	18,8	23,5 ^s

IMC: índice de masa corporal; PAS: presión arterial sistólica; PAD: presión arterial diastólica. Resultados expresados como $\bar{x} \pm DE$;

$$F_{\text{obs}} = 14,4^2 / 10,6^2 = 207,36 / 112,36 = 1,85$$

$$F_{\text{critico}} = F_{15;14} = 2,62$$

FUNCIÓN DE DISTRIBUCIÓN FISHER-SNEDECOR (Continuación)

GL 1

1- α	GL 2	1	2	3	4	5	6	7	8	9	10	12	15	20	30	60	120	∞
0,9	9	3,36	3,01	2,81	2,69	2,61	2,55	2,51	2,47	2,44	2,42	2,38	2,34	2,30	2,25	2,21	2,18	2,16
0,95	9	5,12	4,26	3,86	3,63	3,48	3,37	3,29	3,23	3,18	3,14	3,07	3,01	2,94	2,86	2,79	2,75	2,71
0,975	9	7,21	5,71	5,08	4,72	4,48	4,32	4,20	4,10	4,03	3,96	3,87	3,77	3,67	3,56	3,45	3,39	3,33
0,99	9	10,56	8,02	6,99	6,42	6,06	5,80	5,61	5,47	5,35	5,26	5,11	4,96	4,81	4,65	4,48	4,40	4,31
0,995	9	13,61	10,11	8,72	7,96	7,47	7,13	6,88	6,69	6,54	6,42	6,23	6,03	5,83	5,62	5,41	5,30	5,19
0,9	10	3,29	2,92	2,73	2,61	2,52	2,46	2,41	2,38	2,35	2,32	2,28	2,24	2,20	2,16	2,11	2,08	2,06
0,95	10	4,96	4,10	3,71	3,48	3,33	3,22	3,14	3,07	3,02	2,98	2,91	2,85	2,77	2,70	2,62	2,58	2,54
0,975	10	6,94	5,46	4,83	4,47	4,24	4,07	3,95	3,85	3,78	3,72	3,62	3,52	3,42	3,31	3,20	3,14	3,08
0,99	10	10,04	7,56	6,55	5,99	5,64	5,39	5,20	5,06	4,94	4,85	4,71	4,56	4,41	4,25	4,08	4,00	3,91
0,995	10	12,83	9,43	8,08	7,34	6,87	6,54	6,30	6,12	5,97	5,85	5,66	5,47	5,27	5,07	4,86	4,75	4,64
0,9	12	3,18	2,81	2,61	2,48	2,39	2,33	2,28	2,24	2,21	2,19	2,15	2,10	2,06	2,01	1,96	1,93	1,90
0,95	12	4,75	3,89	3,49	3,26	3,11	3,00	2,91	2,85	2,80	2,75	2,69	2,62	2,54	2,47	2,38	2,34	2,30
0,975	12	6,55	5,10	4,47	4,12	3,89	3,73	3,61	3,51	3,44	3,37	3,28	3,18	3,07	2,96	2,85	2,79	2,72
0,99	12	9,33	6,93	5,95	5,41	5,06	4,82	4,64	4,50	4,39	4,30	4,16	4,01	3,86	3,70	3,54	3,45	3,36
0,995	12	11,75	8,51	7,23	6,52	6,07	5,76	5,52	5,35	5,20	5,09	4,91	4,72	4,53	4,33	4,12	4,01	3,90
0,9	15	3,07	2,70	2,49	2,36	2,27	2,21	2,16	2,12	2,09	2,06	2,02	1,97	1,92	1,87	1,82	1,79	1,76
0,95	15	4,54	3,68	3,29	3,06	2,90	2,79	2,71	2,64	2,59	2,54	2,48	2,40	2,33	2,25	2,16	2,11	2,07
0,975	15	6,20	4,77	4,15	3,80	3,58	3,41	3,29	3,20	3,12	3,06	2,96	2,86	2,76	2,64	2,52	2,46	2,40
0,99	15	8,68	6,36	5,42	4,89	4,56	4,32	4,14	4,00	3,89	3,80	3,67	3,52	3,37	3,21	3,05	2,96	2,87
0,995	15	10,80	7,70	6,48	5,80	5,37	5,07	4,85	4,67	4,54	4,42	4,25	4,07	3,88	3,69	3,48	3,37	3,26
0,9	20	2,97	2,59	2,38	2,25	2,16	2,09	2,04	2,00	1,96	1,94	1,89	1,84	1,79	1,74	1,68	1,64	1,61
0,95	20	4,35	3,49	3,10	2,87	2,71	2,60	2,51	2,45	2,39	2,35	2,28	2,20	2,12	2,04	1,95	1,90	1,84
0,975	20	5,87	4,46	3,86	3,51	3,29	3,13	3,01	2,91	2,84	2,77	2,68	2,57	2,46	2,35	2,22	2,16	2,09
0,99	20	8,10	5,85	4,94	4,43	4,10	3,87	3,70	3,56	3,46	3,37	3,23	3,09	2,94	2,78	2,61	2,52	2,42
0,995	20	9,94	6,99	5,82	5,17	4,76	4,47	4,26	4,09	3,96	3,85	3,68	3,50	3,32	3,12	2,92	2,81	2,69
0,9	30	2,88	2,49	2,28	2,14	2,05	1,98	1,93	1,88	1,85	1,82	1,77	1,72	1,67	1,61	1,54	1,50	1,46
0,95	30	4,17	3,32	2,92	2,69	2,53	2,42	2,33	2,27	2,21	2,16	2,09	2,01	1,93	1,84	1,74	1,68	1,62
0,975	30	5,57	4,18	3,59	3,25	3,03	2,87	2,75	2,65	2,57	2,51	2,41	2,31	2,20	2,07	1,94	1,87	1,79
0,99	30	7,56	5,39	4,51	4,02	3,70	3,47	3,30	3,17	3,07	2,98	2,84	2,70	2,55	2,39	2,21	2,11	2,01
0,995	30	9,18	6,35	5,24	4,62	4,23	3,95	3,74	3,58	3,45	3,34	3,18	3,01	2,82	2,63	2,42	2,30	2,18
0,9	60	2,79	2,39	2,18	2,04	1,95	1,87	1,82	1,77	1,74	1,71	1,66	1,60	1,54	1,48	1,40	1,35	1,29
0,95	60	4,00	3,15	2,76	2,53	2,37	2,25	2,17	2,10	2,04	1,99	1,92	1,84	1,75	1,65	1,53	1,47	1,39
0,975	60	5,29	3,93	3,34	3,01	2,79	2,63	2,51	2,41	2,33	2,27	2,17	2,06	1,94	1,82	1,67	1,58	1,48
0,99	60	7,08	4,98	4,13	3,65	3,34	3,12	2,95	2,82	2,72	2,63	2,50	2,35	2,20	2,03	1,84	1,73	1,60
0,995	60	8,49	5,79	4,73	4,14	3,76	3,49	3,29	3,13	3,01	2,90	2,74	2,57	2,39	2,19	1,96	1,83	1,69
0,9	120	2,75	2,35	2,13	1,99	1,90	1,82	1,77	1,72	1,68	1,65	1,60	1,55	1,48	1,41	1,32	1,26	1,19
0,95	120	3,92	3,07	2,68	2,45	2,29	2,18	2,09	2,02	1,96	1,91	1,83	1,75	1,66	1,55	1,43	1,35	1,25
0,975	120	5,15	3,80	3,23	2,89	2,67	2,52	2,39	2,30	2,22	2,16	2,05	1,94	1,82	1,69	1,53	1,43	1,31
0,99	120	6,85	4,79	3,95	3,48	3,17	2,96	2,79	2,66	2,56	2,47	2,34	2,19	2,03	1,86	1,66	1,53	1,38
0,995	120	8,18	5,54	4,50	3,92	3,55	3,28	3,09	2,93	2,81	2,71	2,54	2,37	2,19	1,98	1,75	1,61	1,43

Comparando dos promedios con muestras independientes

Comparamos dos promedios haciendo inferencia sobre $\mu_1 - \mu_2$, la diferencia entre los dos promedios poblacionales.

- Si los dos promedios poblacionales son iguales, entonces $\mu_1 - \mu_2 = 0$.
- El mejor estimador de $\mu_1 - \mu_2$ es la diferencia entre los dos promedios muestrales,

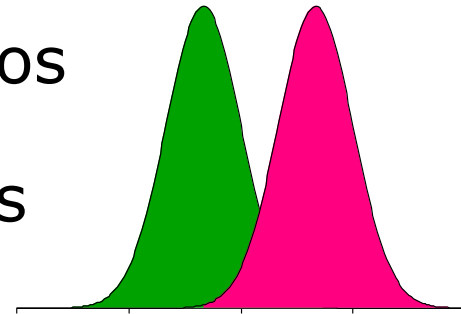
$$\bar{x}_1 - \bar{x}_2$$

Comparando dos promedios con muestras independientes

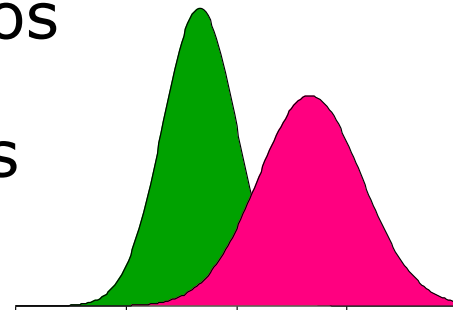
Hay tres situaciones posibles:

❑ que los desvíos poblacionales de las dos poblaciones sean **conocidos** (muy raro!)

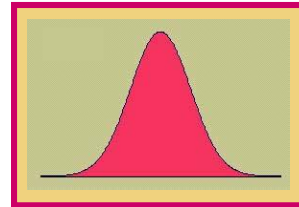
❑ que los desvíos poblacionales de las dos poblaciones sean **desconocidos**, pero que la variabilidad de las dos poblaciones que se comparan **no difiera**



❑ que los desvíos poblacionales de las dos poblaciones sean **desconocidos**, pero que la variabilidad de las dos poblaciones que se comparan **sea distinta**



Distribución muestral de $\bar{x}_1 - \bar{x}_2$ con desvíos poblacionales conocidos



1. La media de $\bar{x}_1 - \bar{x}_2$ es $\mu_1 - \mu_2$, la diferencia entre las medias poblacionales.
2. El desvío estándar (EE) de $\bar{x}_1 - \bar{x}_2$ es $\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$
3. Si la población original sigue una distribución normal o si el tamaño de ambas muestras es lo suficientemente grande, $Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$ sigue una distribución normal

¿Y si queremos comparar dos promedios y no conocemos los desvíos poblacionales?

- ❑ Se utiliza la distribución ***t* de Student**
- ❑ Para poder comparar los promedios es necesario determinar si las varianzas de las dos poblaciones son iguales o no.
- ❑ Por ello deben compararse previamente las varianzas de las dos poblaciones (**Prueba F**)

Comparando dos promedios con desvíos poblacionales desconocidos

- Si se concluye que las varianzas de las poblaciones no difieren se calcula una varianza amalgamada s_a^2 “promediando” las varianzas de las dos muestras

$$S_a^2 = \frac{S_1^2(n_1 - 1) + S_2^2(n_2 - 1)}{n_1 + n_2 - 2}$$

y los GL de la t resultan de sumar los GL de las dos muestras ($n_1 + n_2 - 2$)

- Si se concluye que las poblaciones poseen varianzas distintas no es correcto amalgamar las varianzas muestrales y se pierden GL.
- Los GL de la t resultan de una fórmula (Welsch)

Comparación de dos promedios con desvíos poblacionales desconocidos y supuestamente iguales

$$t_{muestral} = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{S_a^2}{n_1} + \frac{S_a^2}{n_2}}}$$

$$t_{crit} = t_{n_1+n_2-2}$$

Comparación de dos promedios con desvíos poblacionales desconocidos y supuestamente distintos

$$t_{muestral} = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} \quad t_{crit} = t_{GL} \quad GL(Welsh) = \frac{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}\right)^2}{\left(\frac{S_1^2}{n_1}\right)^2 \frac{1}{n_1-1} + \left(\frac{S_2^2}{n_2}\right)^2 \frac{1}{n_2-1}}$$

Volviendo al ejemplo:

TABLA 1

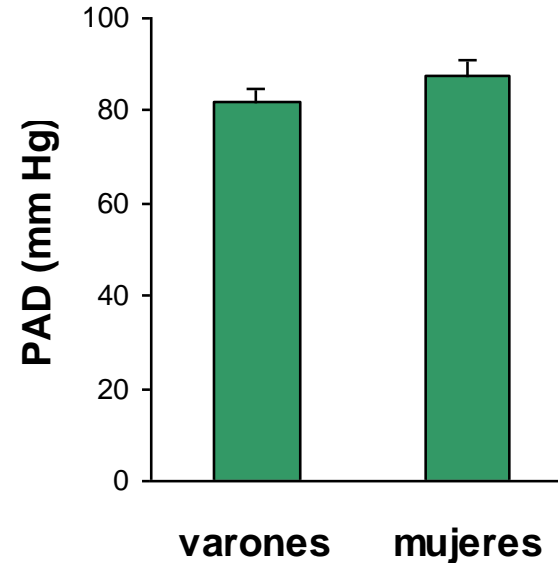
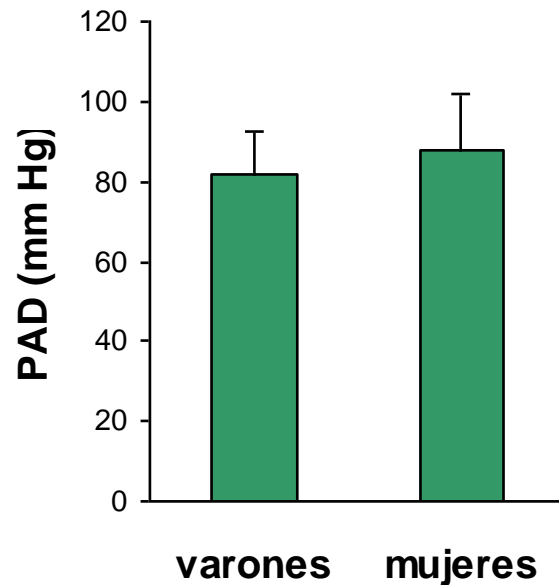
Características generales de los distintos grupos estudiados

	Normales (grupo N)		Diabéticos (grupo DM)	
	Varones	Mujeres	Varones	Mujeres
Número de individuos	370	347	15	16
PAS (mmHg)	125,7 ± 16	125,4 ± 21,1 ^a	148,6 ± 24	158,7 ± 29 ^a
PAD (mmHg)	71,3 ± 9,2	71,5 ± 11,8 ^a	82 ± 10,6	87,5 ± 14,4 ^a
Antecedentes familiares (%)	18,1	16,7 ^a	18,8	23,5 ^a

IMC: índice de masa corporal; PAS: presión arterial sistólica; PAD: presión arterial diastólica. Resultados expresados como $\bar{x} \pm DE$.

¿En pacientes diabéticos, existen diferencias entre hombres y mujeres en los valores generales de presión arterial diastólica?

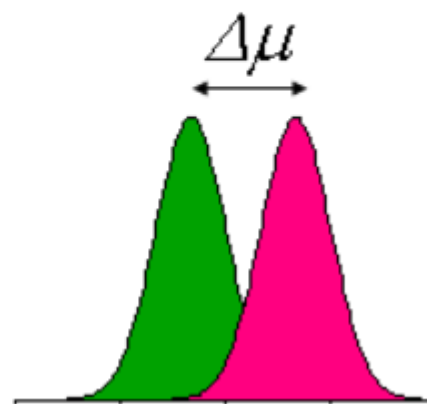
¿Cómo se informan los resultados?



Media, desvío estándar: da idea de la dispersión de la variable
Media, EE: da idea de la precisión en la estimación del parámetro
Magnitud del efecto d: las mujeres poseen en promedio 5.5 mm Hg más de PAS que los hombres, es decir un 7% más ($p < 0.05$)

Estimación de la diferencia entre dos medias

- ❑ Se desea estimar la diferencia promedio entre dos poblaciones
- ❑ Se dispone de dos muestras **independientes**
- ❑ Se desconocen las varianzas poblacionales de ambos grupos, **pero se supone que son iguales**
- ❑ Se supone que ambas poblaciones poseen distribución normal o bien el tamaño de ambas muestras debe ser grande



Intervalo de confianza para la diferencia entre dos medias

$$\hat{\theta} \pm VC EE_{(\hat{\theta})}$$

$$\Delta \bar{x} \pm t_{n_1+n_2-2; 1-\alpha/2} s_a \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$



Comparando dos muestras dependientes



- Un investigador cree que los fumadores tienden a fumar más durante los períodos de estrés.
- Encuesta a un grupo de fumadores en condiciones normales y al mismo grupo cuando está bajo estrés
- Para efectuar esta comparación se requiere:

Una muestra aleatoria de tamaño n de las diferencias d_i entre las dos situaciones extraída de la población con parámetro μ_d

Ejemplo

cantidad de cigarrillos diarios fumados

Individuo	Sin estrés	Con estrés	d
1	15	20	
2	31	45	
3	50	48	
4	16	30	
5	56	72	

- La media de \bar{x}_d es μ_d
- El desvío estándar (EE) de \bar{x}_d es s_d / \sqrt{n}
- Como el desvío poblacional es desconocido, $t = \frac{\bar{x}_d - \mu_d}{s_d / \sqrt{n}}$ la distribución muestral es **t de Student**

Muestras dependientes

Dif= con estrés - sin estrés

- $H_0: \mu_d \leq 0$

- $H_1: \mu_d > 0$

- CR: p - valor $< \alpha$

$$\alpha = 0,05$$

Individuo	Dif
1	5
2	14
3	- 2
4	14
5	16

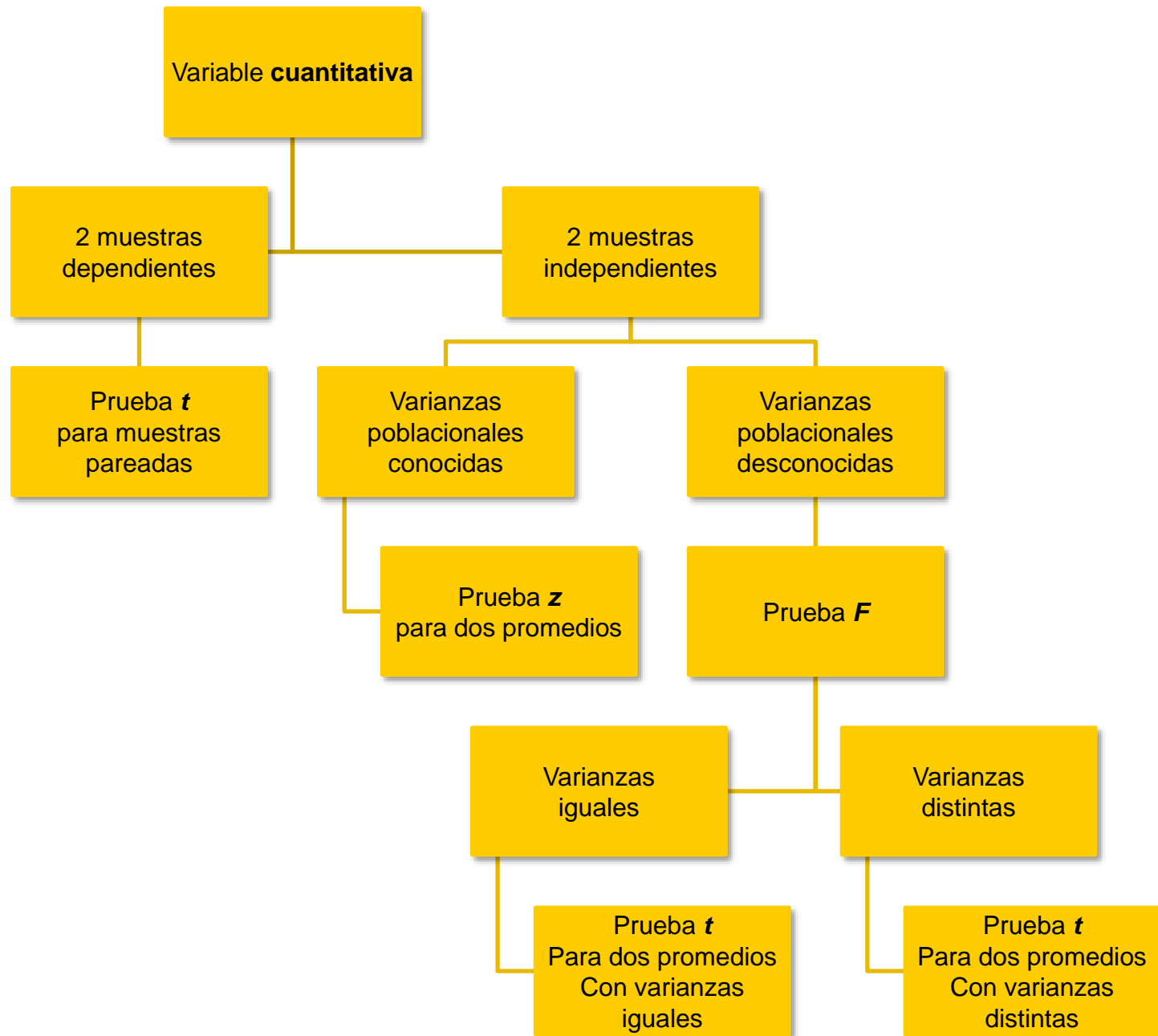
$$\bar{x}_d = 9,40$$

$$S_d = 7,67$$

$$t = \frac{\bar{x}_d - \mu_d}{s_d / \sqrt{n}} = 2,74 \quad t_{gl} = 4$$

P-valor = 0.026

En resumen:



Supuestos

Para que las conclusiones sean válidas, se deben verificar los supuestos de la prueba.

Para PH para 2 medias con desvíos poblacionales desconocidos:

- muestras aleatorias y observaciones independientes
- distribución normal de la variable en ambas poblaciones o tamaño de ambas muestras suficientemente grandes

Para PH para una proporción:

- muestras aleatorias y observaciones independientes
- tamaños de ambas muestras suficientemente grandes; $pn > 5$ y $qn > 5$

Pruebas no paramétricas

- Para todas las pruebas vistas existen pruebas equivalentes no paramétricas, es decir que no exigen que la variable siga cierta distribución de probabilidades.
- Se las denomina por eso pruebas libres de distribución
- En general no trabajan con promedios sino con estadísticos de posición
- Prueba de t para dos muestras independientes:
Prueba de Mann Whitney
- Prueba de t para dos muestras dependientes:
Prueba de Wilcoxon