

Lectura 5

Inferencia con medias muestrales

Prueba z

Como ya adelantamos, la prueba Z es adecuada para inferir acerca de las medias con una o dos muestras, aunque aquí solo veremos el primer caso. Para poder usarla, debemos **verificar el cumplimiento** de algunas condiciones, muchas de las cuales están asociadas al modelo normal que conocimos en el capítulo anterior:

- La muestra debe tener al menos 30 observaciones. Si la muestra tiene menos de 30 observaciones, se debe conocer la varianza de la población.
- Las observaciones deben ser independientes, es decir que la elección de una observación para la muestra no influye en la selección de las otras.
- La población de donde se obtuvo la muestra sigue aproximadamente una distribución normal.

Esta prueba resulta adecuada si queremos **asegurar o descartar** que la media de la población tiene un cierto **valor hipotético**

Ejemplo de tabla

Obs.	Utilidad [M\$]	Obs.	Utilidad [M\$]	Obs.	Utilidad [M\$]	Obs.	Utilidad [M\$]
1	19,33	6	22,22	11	22,55	16	29,68
2	29,37	7	31,26	12	20,69	17	29,27
3	29,14	8	26,92	13	24,68	18	26,72
4	32,10	9	31,40	14	28,74	19	27,08
5	25,04	10	17,66	15	26,85	20	20,62

Script 5.1: prueba Z para una muestra.

```
1 library(TeachingDemos)
2 library(ggpubr)
3
4 # Ingresar los datos.
5 muestra <- c(19.33, 29.37, 29.14, 32.10, 25.04, 22.22, 31.26, 26.92,
6             31.40, 17.66, 22.55, 20.69, 24.68, 28.74, 26.85, 29.68,
7             29.27, 26.72, 27.08, 20.62)
8
9 # Establecer los datos conocidos.
10 desv_est <- 2.32
11 n <- length(muestra)
12 valor_nulo <- 20
13
14 # Crear gráfico Q-Q para verificar la distribución de la muestra,
15 datos <- data.frame(muestra)
16
17 g <- ggqqplot(datos, x = "muestra", color = "SteelBlue")
18 print(g)
19
20 # Verificar distribución muestral usando la prueba de normalidad
21 # de Shapiro-Wilk.
22 normalidad <- shapiro.test(muestra)
23 print(normalidad)
24
25 # Fijar un nivel de significación.
26 alfa <- 0.01
27
28 # Calcular la media de la muestra.
29
30 cat("\tPrueba Z para una muestra\n\n")
31 media <- mean(muestra)
32 cat("Media =", media, "M$\n")
33
34 # Calcular el estadístico de prueba.
35 Z <- (media - valor_nulo) / desv_est
36 cat("Z =", Z, "\n")
37
38 # Calcular el valor p.
39 p <- 2 * pnorm(Z, lower.tail = FALSE)
40 cat("p =", p, "\n")
41
42 # Hacer la prueba Z con R.
43 prueba <- z.test(media, mu = valor_nulo, alternative = "two.sided",
44                  stdev = desv_est, conf.level = 1-alfa)
45 print(prueba)
```

Prueba t

Aunque la prueba t no opera bajo el control de supuesto de normalidad, aún así requiere verificar algunas condiciones para poder usarla:

- 1) Las observaciones son independientes entre sí
- 2) Las observaciones provienen de una distribución cercana a la normal

->No limitan el tamaño de la muestra para que sea mayor a 30, es decir permite su uso para muestras pequeñas

Ejemplo de tablas

Obs.	t [ms]	Obs.	t [ms]	Obs.	t [ms]
1	411,5538	6	388,6731	11	418,1169
2	393,2753	7	430,0382	12	408,4110
3	445,8905	8	469,4734	13	463,3733
4	411,4022	9	409,5844	14	407,0908
5	498,8969	10	442,0800	15	516,5222

Script 5.2: prueba t para una muestra.

```
1 library(ggpubr)
2
3 # Cargar los datos.
4 tiempo <- c(411.5538, 393.2753, 445.8905, 411.4022, 498.8969,
5            388.6731, 430.0382, 469.4734, 409.5844, 442.0800,
6            418.1169, 408.4110, 463.3733, 407.0908, 516.5222)
7
8 # Establecer los datos conocidos.
9 n <- length(tiempo)
10 grados_libertad <- n - 1
11 valor_nulo <- 500
12
13
14 # Verificar si la distribución se acerca a la normal.
15 g <- ggqqplot(data = data.frame(tiempo),
16              x = "tiempo",
17              color = "steelblue",
18              xlab = "Teórico",
19              ylab = "Muestra",
20              title = "Gráfico Q-Q muestra v/s distr. normal")
21
22 print(g)
23
24 # Fijar un nivel de significación.
25 alfa <- 0.025
26
27 # Calcular el estadístico de prueba.
28 cat("\tPrueba t para una muestra\n\n")
29 media <- mean(tiempo)
30 cat("Media =", media, "M$\n")
31 desv_est <- sd(tiempo)
32 error <- desv_est / sqrt(n)
33 t <- (media - valor_nulo) / error
34 cat("t =", t, "\n")
--
```

```

--
36 # Calcular el valor p.
37 p <- pt(t, df = grados_libertad, lower.tail = TRUE)
38 cat("p =", p, "\n")
39
40 # Construir el intervalo de confianza.
41 t_critico <- qt(alfa, df = grados_libertad, lower.tail = FALSE)
42 superior <- media + t_critico * error
43 cat("Intervalo de confianza = (-Inf, ", superior, "]\n", sep = "")
44
45 # Aplicar la prueba t de Student con la función de R.
46 prueba <- t.test(tiempo,
47                  alternative = "less",
48                  mu = valor_nulo,
49                  conf.level = 1 - alfa)
50
51 print(prueba)

```

Prueba t para dos muestras pareadas

Ejemplo de tablas

instancia	t_A [ms]	t_B [ms]	dif [ms]	instancia	t_A [ms]	t_B [ms]	dif [ms]
1	436,5736	408,5142	28,0594	19	438,5959	458,2536	-19,6577
2	470,7937	450,1075	20,6862	20	439,7409	474,9863	-35,2454
3	445,8354	490,2311	-44,3957	21	464,5916	496,0153	-31,4237
4	470,9810	513,6910	-42,7100	22	467,9926	485,8112	-17,8186
5	485,9394	467,6467	18,2927	23	415,3252	457,4253	-42,1001
6	464,6145	484,1897	-19,5752	24	495,4094	483,3700	12,0394
7	466,2139	465,9334	0,2805	25	493,7082	510,7131	-17,0049
8	468,9065	502,6670	-33,7605	26	433,1082	467,5739	-34,4657
9	473,8778	444,9693	28,9085	27	445,7433	482,5621	-36,8188
10	413,0639	456,3341	-43,2702	28	515,2049	453,5986	61,6063
11	496,8705	501,1443	-4,2738	29	441,9420	385,9391	56,0029
12	450,6578	471,7833	-21,1255	30	472,1396	548,7884	-76,6488
13	502,9759	441,1206	61,8553	31	451,2234	467,2533	-16,0299
14	465,6358	544,1575	-78,5217	32	476,5149	494,7049	-18,1900
15	437,6397	447,8844	-10,2447	33	440,7918	451,9716	-11,1798
16	458,8806	432,4108	26,4698	34	460,1070	522,3699	-62,2629
17	503,1435	477,1712	25,9723	35	450,1008	444,1270	5,9738
18	430,0524	482,4828	-52,4304				

```

1 # Cargar los datos.
2 instancia <- seq(1, 35, 1)
3
4 t_A <- c(436.5736, 470.7937, 445.8354, 470.9810, 485.9394,
5         464.6145, 466.2139, 468.9065, 473.8778, 413.0639,
6         496.8705, 450.6578, 502.9759, 465.6358, 437.6397,
7         458.8806, 503.1435, 430.0524, 438.5959, 439.7409,
8         464.5916, 467.9926, 415.3252, 495.4094, 493.7082,
9         433.1082, 445.7433, 515.2049, 441.9420, 472.1396,
10        451.2234, 476.5149, 440.7918, 460.1070, 450.1008)
11
12 t_B <- c(408.5142, 450.1075, 490.2311, 513.6910, 467.6467,
13         484.1897, 465.9334, 502.6670, 444.9693, 456.3341,
14         501.1443, 471.7833, 441.1206, 544.1575, 447.8844,
15         432.4108, 477.1712, 482.4828, 458.2536, 474.9863,
16         496.0153, 485.8112, 457.4253, 483.3700, 510.7131,
17         467.5739, 482.5621, 453.5986, 385.9391, 548.7884,
18         467.2533, 494.7049, 451.9716, 522.3699, 444.1270)
19
20 diferencia <- t_A - t_B
21
22 # Verificar si la distribución se acerca a la normal.
23 normalidad <- shapiro.test(diferencia)
24 print(normalidad)
25
26 # Fijar un nivel de significación.
27 alfa <- 0.05
28
29 # Aplicar la prueba t de Student a la diferencia de medias.
30 prueba_1 <- t.test(diferencia,
31                   alternative = "two.sided",
32                   mu = valor_nulo,
33                   conf.level = 1 - alfa)
34
35 print(prueba_1)
36
37 # Otra alternativa puede ser aplicar la prueba t de Student
38 # para dos muestras pareadas.
39 prueba_2 <- t.test(x = t_A,
40                   y = t_B,

```

```

41         paired = TRUE,
42         alternative = "two.sided",
43         mu = valor_nulo,
44         conf.level = 1 - alfa)
45
46 print(prueba_2)

```

Prueba t para muestras independientes

En este caso, la prueba t se usa para comparar las medias de dos poblaciones en que las observaciones con que se cuenta no tienen relación con ninguna de las otras observaciones, ni influyen en su selección, ni en la misma ni en la otra muestra. En este caso la inferencia se hace sobre la diferencia de las medias: $\mu_1 - \mu_2 = d_0$, donde d_0 es un valor hipotético fijo para la diferencia. Usualmente se usa $d_0 = 0$, en cuyo caso las muestras podrían provenir de dos poblaciones distintas con igual media, o desde la misma población. Para ello, la prueba usa como estimador puntual la diferencia de las medias muestrales ($\bar{x}_1 - \bar{x}_2$). Así, el estadístico T en este caso toma la forma de la ecuación 5.3.

$$T = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{SE_{(\bar{x}_1 - \bar{x}_2)}} \quad (5.3)$$

Al usar la distribución t Student para la diferencia de medias, se deben cumplir los siguientes requisitos:

- 1) Cada una de las muestras cumple con las condiciones para usar la distribución t
- 2) Las muestras son independientes entre sí

Ejemplo de tablas

Anticuerpos [mg/ml]	
Vacuna A	Vacuna B
6,04	5,32
19,84	3,31
8,62	5,68
13,02	5,73
12,20	4,86
14,78	5,68
4,53	2,93
26,67	5,48
3,14	6,10
19,14	2,56
10,86	7,52
13,13	7,41
6,34	4,02
11,16	
7,62	

Tabla 5.4: Concentración de anticuerpos de los pacientes vacunados.

Script 5.4: prueba t para dos muestras independientes.

```
1 library(ggpubr)
2
3 # Cargar los datos.
4 vacuna_A <- c(6.04, 19.84, 8.62, 13.02, 12.20, 14.78, 4.53, 26.67,
5              3.14, 19.14, 10.86, 13.13, 6.34, 11.16, 7.62)
6
7 vacuna_B <- c(5.32, 3.31, 5.68, 5.73, 4.86, 5.68, 2.93, 5.48, 6.10,
8              2.56, 7.52, 7.41, 4.02)
9
10 # Verificar si las muestras se distribuyen de manera cercana
11 # a la normal.
12 normalidad_A <- shapiro.test(vacuna_A)
13 print(normalidad_A)
14 normalidad_B <- shapiro.test(vacuna_B)
15 print(normalidad_B)
16
17 # Fijar un nivel de significación.
18 alfa <- 0.01
19
20 # Aplicar la prueba t para dos muestras independientes.
21 prueba <- t.test(x = vacuna_A,
22                 y = vacuna_B,
23                 paired = FALSE,
24                 alternative = "greater",
25                 mu = 0,
26                 conf.level = 1 - alfa)
27
28 print(prueba)
29
30 # Calcular la diferencia entre las medias.
31 media_A <- mean(vacuna_A)
32 media_B <- mean(vacuna_B)
33 diferencia <- media_A - media_B
34 cat("Diferencia de las medias =", diferencia, "[mg/ml]\n")
```