

Estadística

Práctica 9

Contrastes de hipótesis paramétricos

Contenido

1. Introducción	3
2. Contraste de hipótesis para la media poblacional	5
3. Ejemplo: Contraste de hipótesis sobre la media de calificaciones de una población de estudiantes (muestra grande)	7
4. Contraste de hipótesis para proporciones poblacionales	15
5. Ejemplo: Contraste de hipótesis sobre la proporción de estudiantes que tienen móvil Android... ..	17
6. Contraste de hipótesis sobre la diferencia de medias	20
7. Ejemplo: Contraste de hipótesis sobre la diferencia de medias de calificaciones de estudiantes del turno de mañana y del turno de tarde (muestras grandes e independientes)	21
8. Contraste de hipótesis sobre la diferencia de proporciones	24
9. Ejemplo: Contraste de hipótesis sobre la diferencia de proporciones de estudiantes que tienen móvil Android en el turno de mañana y en el turno de tarde (muestras grandes e independientes) ..	25
10. Ejercicios propuestos	28

1. Introducción

Con esta práctica se utiliza R y RStudio para determinar si se puede aceptar una hipótesis, con un nivel de confianza determinado, para el valor de un parámetro de una población (media) a partir de una muestra de la misma.

En estadística hay que diferenciar entre conceptos “parámetros” y “estadísticos”. El término “parámetro” suele utilizarse en referencia a una población mientras que el término “estadístico” suele utilizarse en referencia a una muestra extraída de dicha población. En la siguiente tabla se muestra la nomenclatura utilizada habitualmente.

PARÁMETROS (POBLACIONALES)	ESTADÍSTICOS (MUESTRALES)
Media poblacional (μ)	Media muestral (\bar{x})
Desviación típica o estándar poblacional (σ)	Desviación típica o estándar muestral (s)
Varianza poblacional (σ^2)	Varianza muestral (s^2)
Tamaño población (N)	Tamaño muestra (n)
Proporción poblacional (p)	Proporción muestral (\hat{p})

En un contraste de hipótesis:

$$H_0: \text{Hipótesis nula} \text{ vs } H_A: \text{Hipótesis alternativa}$$

Se manejan varios conceptos:

- **Nivel de significación del contraste (α):** Es la probabilidad de rechazar la hipótesis H_0 cuando dicha hipótesis es cierta. Se denomina error tipo I.
- **Nivel de confianza del contraste ($1-\alpha$):** Es la probabilidad de no rechazar la hipótesis H_0 cuando dicha hipótesis es cierta.
- **Potencia del contraste.** Es la probabilidad de rechazar la hipótesis H_0 cuando dicha hipótesis es falsa.
- **Debilidad del contraste (1-potencia del contraste).** Es la probabilidad de no rechazar la hipótesis H_0 cuando dicha hipótesis es falsa. Se denomina error tipo II.

Se resumen en la siguiente tabla.

	Decisión: Rechazar H_0	Decisión: No rechazar H_0	Probabilidad de acertar
Realidad: Hipótesis H_0 CIERTA	FALLO Error tipo I	ACIERTO	($1 - \alpha$) es la probabilidad de acertar cuando la hipótesis es cierta
Realidad: Hipótesis H_0 FALSA	ACIERTO	FALLO Error tipo II	($1 - \text{potencia}$) es la probabilidad de acertar cuando la hipótesis es falsa

Un contraste de hipótesis puede ser

- Paramétrico: cuando los datos de la población tienen una distribución de probabilidad Normal
- No Paramétrico: Cuando no hay evidencias de que los datos tengan una distribución Normal

Por tanto, lo primero que hay que hacer antes de resolver un contraste es comprobar si la población tiene una distribución Normal, para decidir si hay que hacer un contraste paramétrico o no paramétrico. Si no se tiene acceso a los datos de toda la población habrá que hacer la comprobación con la muestra disponible.

Si superan la prueba de normalidad se pueden aplicar los contrastes no paramétricos tratados en esta práctica.

Para realizar la prueba de normalidad, se suele usar la función `shapiro.test()` de R cuando la muestra es pequeña¹ (prueba de Shapiro-Wilk), o la función `ks.test()` cuando la muestra es grande (prueba de Kolmogorov-Smirnov), con el siguiente formato:

```
> shapiro.test(muestra)

> ks.test(muestra, pnorm, mean(muestra), sd(muestra))
```

Si el resultado obtiene un p-valor < 0.05, con un nivel de confianza del 95% no puede afirmarse que los datos tengan una distribución Normal.

Por ejemplo:

```
> x = c(4.4, 3, 0.6, 1.7, 0.4, 1.0, 1.4, 0.9, 4.5, 0.2)
> shapiro.test(x)
      Shapiro-Wilk normality test
data:  x
W = 0.84364, p-value = 0.04878
```

En este caso el p-valor es menor que 0.05, por lo que los datos de la muestra no son Normales y hay que aplicar contrastes de hipótesis no paramétricos.

En esta práctica nos limitaremos a los contrastes paramétricos. Existen dos métodos para realizar un contraste de hipótesis paramétrico:

1. Comprobando la región de no rechazo
2. Calculando el p-valor

¹ No hay unanimidad respecto a lo que se considera muestra pequeña o grande. En general, para estas pruebas se suele suponer el límite en 50 datos. Aunque hay [autores](#) que defienden que el test de Shapiro-Wilk sólo debería realizarse para una muestra de como mínimo 30 datos.

2. Contraste de hipótesis para la media poblacional

Se calcula el valor del estadístico de contraste (z_0 si es muestra grande, o t_0 si es muestra pequeña), después se calculan los límites de la región de aceptación, y si el valor del estadístico de contraste está:

- Dentro de la región, entonces no se puede rechazar la hipótesis H_0 , por lo que se acepta y se rechaza H_A .
- Fuera de la región, entonces se rechaza la hipótesis H_0 y se acepta H_A .

Otra opción es calcular el p-valor y si el p-valor es:

- Mayor que α , entonces no se puede rechazar H_0 .
- Menor que α , entonces se rechaza H_0 .

Varianza conocida	Contraste de hipótesis	Región de aceptación y código R para calcularla	p-valor y código R para calcularlo
No $n \geq 30$	$H_0: \mu = \mu_0$ vs $H_A: \mu < \mu_0$	$(-z_\alpha, +\infty)$ > z.alfa=qnorm(1-alfa, 0, 1)	$P\{Z \leq Z_0\}$ > p.valor=pnorm(z0, 0, 1)
	$H_0: \mu = \mu_0$ vs $H_A: \mu > \mu_0$	$(-\infty, z_\alpha)$ > z.alfa=qnorm(1-alfa, 0, 1)	$P\{Z \geq Z_0\}$ > p.valor=1-pnorm(z0, 0, 1)
	$H_0: \mu = \mu_0$ vs $H_A: \mu \neq \mu_0$	$(-z_{\alpha/2}, z_{\alpha/2})$ > z.alfa.medios=qnorm(1-alfa/2, 0, 1)	$2P\{Z \geq Z_0 \}$ > p.valor=2*(1-pnorm(abs(z0), 0, 1))
No $n < 30$	$H_0: \mu = \mu_0$ vs $H_A: \mu < \mu_0$	$(-t_{\alpha, n-1}, +\infty)$ > t.alfa=qt(1-alfa, n-1)	$P\{t \leq t_0\}$ > p.valor=pt(t0, n-1)
	$H_0: \mu = \mu_0$ vs $H_A: \mu > \mu_0$	$(-\infty, t_{\alpha, n-1})$ > t.alfa=qt(1-alfa, n-1)	$P\{t \geq t_0\}$ > p.valor=1-pt(t0, n-1)
	$H_0: \mu = \mu_0$ vs $H_A: \mu \neq \mu_0$	$(-t_{\alpha/2, n-1}, t_{\alpha/2, n-1})$ > t.alfa.medios=qt(1-alfa/2, n-1)	$2P\{t \geq t_0 \}$ > p.valor=2*(1-pt(abs(t0), n-1))

Cálculo del estadístico de contraste

Estadístico de contraste	Código R
$Z_o = \frac{\bar{x} - \mu_o}{\frac{s}{\sqrt{n}}}$	<pre>x=muestra de X m=mean(x) s=sd(x) n=length(x) mo=valor a comprobar en la hipótesis zo=(m-mo) / (s/sqrt(n))</pre>
$t_o = \frac{\bar{x} - \mu_o}{s/\sqrt{n}}$	<pre>x=muestra de X m=mean(x) s=sd(x) n=length(x) mo=valor a comprobar en la hipótesis to=(m-mo) / (s/sqrt(n))</pre>

En lugar de calcular las fórmulas, se pueden usar las funciones z.test y t.test de R.

Varianza conocida	Contraste de hipótesis	Código R
No $n \geq 30$	$H_0: \mu = \mu_o$ vs $H_A: \mu < \mu_o$	<pre>install.packages("BSDA") library(BSDA) > z.test(x, alternative="less", mu=mo, sigma.x=s)</pre>
	$H_0: \mu = \mu_o$ vs $H_A: \mu > \mu_o$	<pre>> z.test(x, alternative="greater", mu=mo, sigma.x=s)</pre>
	$H_0: \mu = \mu_o$ vs $H_A: \mu \neq \mu_o$	<pre>> z.test(x, alternative="two.sided", mu=mo, sigma.x=s)</pre>
No $n < 30$	$H_0: \mu = \mu_o$ vs $H_A: \mu < \mu_o$	<pre>> t.test(x, alternative="less", mu=mo)</pre>
	$H_0: \mu = \mu_o$ vs $H_A: \mu > \mu_o$	<pre>> t.test(x, alternative="greater", mu=mo)</pre>
	$H_0: \mu = \mu_o$ vs $H_A: \mu \neq \mu_o$	<pre>> t.test(x, alternative="two.sided", mu=mo)</pre>

3. Ejemplo: Contraste de hipótesis sobre la media de calificaciones de una población de estudiantes (muestra grande)

ENUNCIADO

Mediante una encuesta, se saben las calificaciones de acceso a la universidad de 74 estudiantes de un curso de la asignatura Estadística del Grado en Ingeniería en Sistemas de Información de la Universidad de Alcalá, de una población de 108 matriculados.

- a) Comprobar si las calificaciones tienen una distribución Normal, para poder aplicar contrastes paramétricos
- b) ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de toda la población es diferente a 8.6?
 1. Resolverlo comprobando la región de aceptación.
 2. Resolverlo calculando el p-valor aplicando fórmulas
 3. Resolverlo usando la función `z.test()` del paquete BSDA de R
- c) ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de toda la población de alumnos matriculados en la asignatura es inferior a 8.6?
 1. Resolverlo comprobando la región de aceptación.
 2. Resolverlo calculando el p-valor aplicando fórmulas
 3. Resolverlo usando la función `z.test()` del paquete BSDA de R
- d) ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de toda la población de alumnos matriculados en la asignatura es superior a 8.6?
 1. Resolverlo comprobando la región de aceptación
 2. Resolverlo calculando el p-valor aplicando fórmulas
 3. Resolverlo usando la función `z.test()` del paquete BSDA de R

SOLUCIÓN

Primero hay que leer los datos de las notas, disponibles en el archivo [encuesta.csv](#):

```
> encuesta = read.csv2("encuesta.csv")
> (nota=encuesta$NOTA)
[1] 8.50 7.10 8.63 8.62 8.20 8.70 7.21 7.63 8.40 8.30
[11] 8.20 9.10 9.79 10.11 8.02 7.31 7.50 8.71 8.34 9.21
[21] 7.80 10.30 7.99 6.90 7.80 10.00 8.59 7.00 8.05 10.80
[31] 7.99 8.55 7.34 6.75 9.56 7.42 6.94 7.21 7.68 10.28
[41] 7.86 10.26 7.27 5.80 7.30 7.14 8.60 7.50 8.00 7.54
[51] 7.29 7.83 6.75 9.81 6.80 6.44 6.65 7.80 10.27 7.60
[61] 7.87 7.00 7.08 7.48 8.07 5.82 6.50 9.90 7.50 6.50
[71] 9.46 8.00 7.80 7.65
```

a) Comprobar si las calificaciones tienen una distribución Normal, para poder aplicar contrastes paramétricos

Como el tamaño de la muestra es grande, se realiza una prueba de normalidad con el test de Kolmogorov-Smirnov.

```
> ks.test(nota, pnorm, mean(nota), sd(nota))  
  
Asymptotic one-sample Kolmogorov-Smirnov test  
  
data: nota  
D = 0.11842, p-value = 0.2505  
alternative hypothesis: two-sided
```

Se obtiene un p-valor mayor que 0.05, por lo que se puede aceptar que la población es Normal.

b) ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de toda la población es diferente a 8.6?

No se conoce la varianza poblacional y es una muestra grande ($74 > 30$).

Primero se crean las variables alfa, n, s, m y mo.

```
> (alfa=1-0.95)  
[1] 0.05  
> (n=length(nota))  
[1] 74  
> (s=sd(nota))  
[1] 1.13154  
> (m=mean(nota))  
[1] 8.022568  
> (mo=8.6)  
[1] 8.6
```

Se formula el siguiente contraste de hipótesis:

$$H_0: \mu = 8.6 \text{ vs } H_A: \mu \neq 8.6$$

b.1) Resolverlo comprobando la región de aceptación.

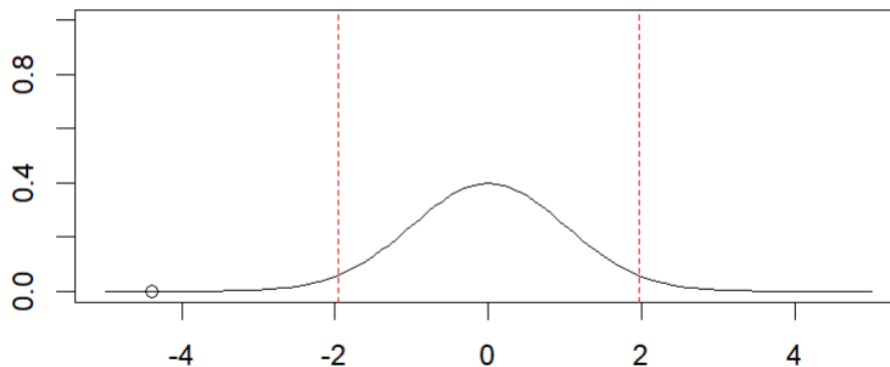
La región de aceptación es $(-z_{\alpha/2}, z_{\alpha/2})$.

```
> (zo=(m-mo)/(s/sqrt(n)))  
[1] -4.389822  
> (z.alfa.medios=qnorm(1-alfa/2, 0, 1))  
[1] 1.959964  
> (-z.alfa.medios<=zo)&(zo<=z.alfa.medios)  
[1] FALSE
```


Como z_0 (-4.39) no está en la región de aceptación $(-1.96, 1.96)$, entonces se rechaza la hipótesis nula (H_0) y se acepta la alternativa (H_A), por lo que la respuesta es **Sí se puede afirmar con una confianza del 95% que la media poblacional es diferente a 8.6.**

Podemos hacer la comprobación visualmente.

```
> curve(dnorm(x,0,1),-5,5, add = TRUE)
> abline(v=-z.alfa.medios, col="red", lty=2)
> abline(v=z.alfa.medios, col="red", lty=2)
```



En el diagrama, el área bajo la función entre las dos líneas es $1 - \alpha$ (0.95), el área desde la primera línea hacia la izquierda es $\alpha/2$ (0.025), y el área desde la segunda línea hacia la derecha es $\alpha/2$ (0.025). El área total es 1.

La región de aceptación está limitada por las dos líneas rojas, y z_0 es el punto dibujado. Se comprueba visualmente que z_0 está fuera de la región.

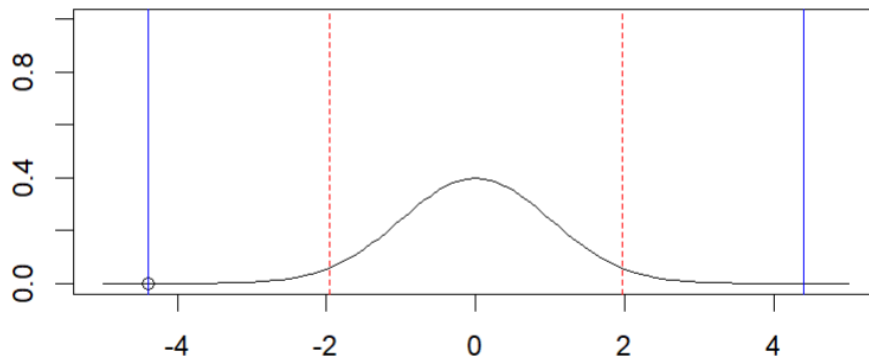
b.2) Resolverlo calculando el p-valor aplicando fórmulas

```
> (p.valor=2*(1-pnorm(abs(z0),0,1)))
[1] 0.00001134432
```

Como el p-valor (0.00001) es menor que α (0.05), como ocurrió en el apartado anterior, se rechaza la hipótesis nula y se acepta la alternativa de que la media poblacional es diferente a 8.6.

En el mismo diagrama del apartado anterior podemos dibujar una línea vertical azul continua en el punto z_0 y otra en $-z_0$.

```
> abline(v=z0, col="blue")
> abline(v=-z0, col="blue")
```



el área desde la primera línea roja hacia la izquierda es $\alpha/2$ (0.025), y el área desde la segunda línea hacia la derecha es $\alpha/2$ (0.025). El área desde la primera línea azul continua hacia la izquierda es $p\text{-valor}/2$ y el área desde la segunda línea azul continua hacia la derecha es $p\text{-valor}/2$, y la suma de ambas es $p\text{-valor}$. Se observa visualmente que $p\text{-valor} < \alpha$.

b.3) Resolverlo usando la función `z.test()` del paquete BSDA de R

```
> z.test(nota, alternative="two.sided", mu=mo, sigma.x=s, conf.level=0.95)
One-sample z-Test
data: nota
z = -4.3898, p-value = 0.00001134
alternative hypothesis: true mean is not equal to 8.6
95 percent confidence interval:
 7.764756 8.280379
sample estimates:
mean of x
 8.022568
```

El $p\text{-valor}$ es el mismo que en apartado anterior. Como el $p\text{-valor}$ es menor que α (0.05), se rechaza la hipótesis nula y se acepta la alternativa de que la media poblacional es diferente a 8.6.

El resultado de la función `z.test()` incluye más información:

- Indica que el valor del estadístico de contraste z_0 (llamado z en el resultado de la función) es -4.3898, que coincide con los apartados anteriores.
- Indica que la hipótesis alternativa es que la media fuera diferente a 8.6, como así se indicó al llamar a la función con el parámetro `alternative="two.sided"`
- Indica que el nivel de confianza es del 95%, como así se indicó al llamar a la función con el parámetro `conf.level=0.95`.
- Indica que el intervalo de confianza para la media poblacional con un nivel del 95% es, redondeando a dos decimales: (7.76, 8.28). Puede comprobarse que 8.6 está fuera del intervalo.
- Indica que la media muestral es 8.022568.

c) ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de toda la población de alumnos matriculados en la asignatura es inferior a 8.6?

Se formula el siguiente contraste de hipótesis:

$$H_0: \mu = 8.6 \text{ vs } H_A: \mu < 8.6$$

c.1) Resolverlo comprobando la región de aceptación

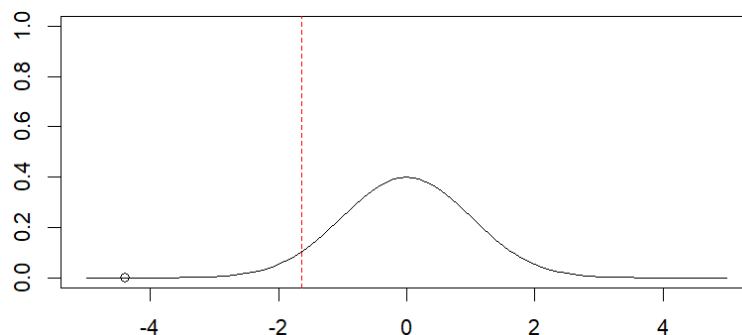
La región de aceptación es $(-z_{\alpha}, +\infty)$.

```
> (zo=(m-mo)/(s/sqrt(n)))
[1] -4.389822
> (z.alfa=qnorm(1-alfa,0,1))
[1] 1.644854
> (zo>-z.alfa) #Como es FALSE se acepta la alternativa. Respuesta SI
[1] FALSE
```

Como z_0 (-4.39) está fuera de la región de aceptación $(-1.64, +\infty)$, entonces se rechaza la hipótesis nula (H_0), y se acepta la hipótesis alternativa (H_A), por lo que la respuesta es **SÍ se puede afirmar con una confianza del 95% que la media poblacional es menor a 8.6.**

Podemos hacer la comprobación visualmente, dibujando la curva de la función de densidad de probabilidad normal $Z: N(0,1)$ para la variable estadístico de contraste, señalando un punto para el valor calculado para el estadístico (z_0), y una línea vertical roja con trazo discontinuo, en el extremo inferior de la región de aceptación ($-z.alfa$), que no tiene límite superior, pues sería $-\infty$. Puede comprobarse visualmente que el punto z_0 está fuera de la región de aceptación.

```
> plot(zo,0, xlim = c(-5,5), ylim=c(0,1))
> curve(dnorm(x,0,1),-5,5, add = TRUE)
> abline(v=-z.alfa, col="red", lty=2)
```



En el diagrama, el área bajo la función desde la línea roja hacia la izquierda es α (0.05), mientras que el área hacia la derecha es $1-\alpha$ (0.95). La región de aceptación es desde la línea roja hasta $+\infty$, y z_0 es el punto dibujado. Se comprueba visualmente que z_0 está fuera de la región.

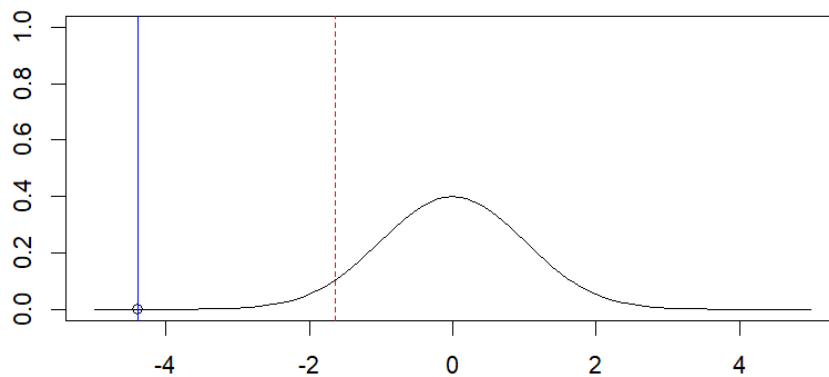
c.2) Resolverlo calculando el p-valor aplicando fórmulas

```
> (p.valor=pnorm(zo,0,1))
[1] 0.000005672161
```

Como el p-valor (0.00000567) es menor que alfa (0.05), se rechaza la hipótesis H_0 , como en el apartado anterior.

En el mismo diagrama del apartado anterior podemos dibujar una línea vertical azul continua en el punto z_0 .

```
> abline(v=zo, col="blue")
```



En el diagrama, el área bajo la función desde la línea roja discontinua hacia la izquierda es alfa (0.05), mientras que el área desde la línea azul continua hacia la izquierda es el p-valor (0.00000563). Se observa visualmente que $p\text{-valor} < \alpha$.

c.3) Resolverlo usando la función `z.test()` del paquete BSDA de R

```
> z.test(nota, alternative="less", mu=mo, sigma.x=s, conf.level=0.95)
One-sample z-Test
data: nota
z = -4.3898, p-value = 0.000005672
alternative hypothesis: true mean is less than 8.6
95 percent confidence interval:
 NA 8.23893
sample estimates:
mean of x
 8.022568
```

El p-valor es el mismo que en el apartado anterior.

d) ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de toda la población es superior a 8.6?

Se formula el siguiente contraste de hipótesis:

$$H_0: \mu = 8.6 \text{ vs } H_A: \mu > 8.6$$

d.1) Resolverlo comprobando la región de aceptación

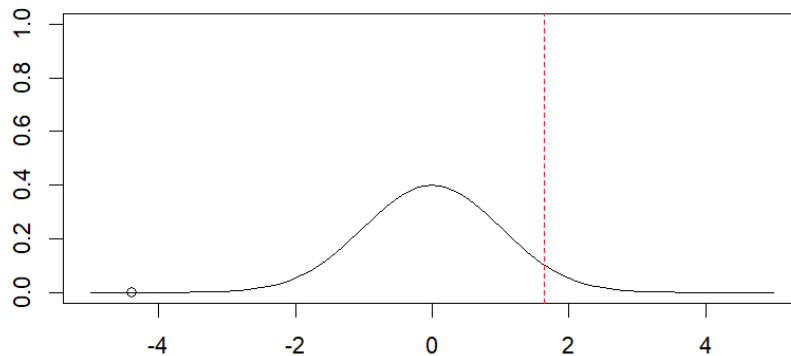
La región de aceptación es $(-\infty, z_\alpha)$.

```
> (zo=(m-mo)/(s/sqrt(n)))  
[1] -4.389822  
> (z.alfa=qnorm(1-alfa,0,1))  
[1] 1.644854  
> (zo<z.alfa)  
[1] TRUE
```

Como z_0 (-4.39) está en la región de aceptación $(-\infty, 1.64)$, entonces no se rechaza la hipótesis nula, y se rechaza la hipótesis alternativa, por lo que la respuesta es que **NO se puede afirmar con una confianza del 95% que la media poblacional sea superior a 8.6.**

Podemos hacer la comprobación visualmente.

```
> plot(zo,0, xlim = c(-5,5), ylim=c(0,1))  
> curve(dnorm(x,0,1),-5,5, add = TRUE)  
> abline(v=z.alfa, col="red", lty=2)
```



En el diagrama, el área bajo la función desde la línea hacia la derecha es alfa (0.05), mientras que el área hacia la izquierda es 1-alfa (0.95). La región de aceptación es desde $-\infty$ hasta la línea roja, y z_0 es el punto dibujado. Se comprueba visualmente que z_0 está dentro de la región

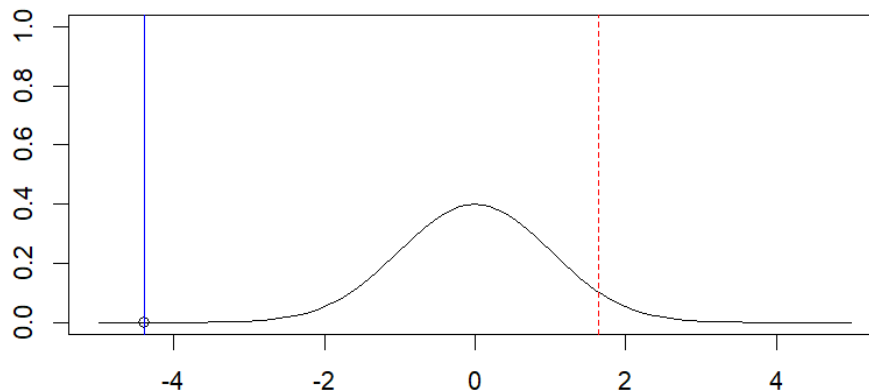
d.2) Resolverlo calculando el p-valor aplicando fórmulas

```
> (p.valor=1-pnorm(zo,0,1))
[1] 0.9999943
```

Como el p-valor (0.9999943) es mayor que alfa (0.05), no se rechaza la hipótesis nula, y se rechaza la hipótesis alternativa de que la media poblacional es mayor que 8.6.

En el mismo diagrama del apartado anterior podemos dibujar una línea vertical azul continua en el punto z_0 .

```
> abline(v=zo, col="blue")
```



En el diagrama, el área bajo la función desde la línea roja discontinua hacia la derecha es alfa (0.05), mientras que el área desde la línea azul continua hacia la derecha es p-valor (0.999). Se observa visualmente que $p\text{-valor} > \alpha$.

d.3) Resolverlo usando la función `z.test()` del paquete BSDA de R

```
> z.test(nota, alternative="greater", mu=mo, sigma.x=s, conf.level=0.95)
One-sample z-Test
data: nota
z = -4.3898, p-value = 1
alternative hypothesis: true mean is greater than 8.6
95 percent confidence interval:
 8.23893      NA
sample estimates:
mean of x
 8.022568
```

El p-valor es el mismo que en apartado anterior (0.9999943), pero está redondeado a 1.

4. Contraste de hipótesis para proporciones poblacionales

Se calcula el valor del estadístico de contraste (z_o), después se calculan los límites de la región de aceptación, y si el valor del estadístico de contraste está:

- Dentro de la región, entonces no se puede rechazar la hipótesis H_0 , por lo que se acepta y se rechaza H_A .
- Fuera de la región, entonces se rechaza la hipótesis H_0 y se acepta H_A .

Otra opción es calcular el p-valor y si el p-valor es:

- Mayor que α , entonces no se puede rechazar H_0 .
- Menor que α , entonces se rechaza H_0 .

Contraste de hipótesis	Región de aceptación y código R para calcularla	p-valor y código R para calcularlo
$H_0: p = p_o$ vs $H_A: p < p_o$	$(-z_\alpha, +\infty)$ > <code>z.alfa=qnorm(1-alfa,0,1)</code>	$P\{Z \leq Z_o\}$ > <code>p.valor=pnorm(z_o,0,1)</code>
$H_0: p = p_o$ vs $H_A: p > p_o$	$(-\infty, z_\alpha)$ > <code>z.alfa=qnorm(1-alfa,0,1)</code>	$P\{Z \geq Z_o\}$ > <code>p.valor=1-pnorm(z_o,0,1)</code>
$H_0: p = p_o$ vs $H_A: p \neq p_o$	$(-z_{\alpha/2}, z_{\alpha/2})$ > <code>z.alfa.medios=qnorm(1-alfa/2,0,1)</code>	$2P\{Z \geq Z_o \}$ > <code>p.valor=2*(1-pnorm(abs(z_o),0,1))</code>

Cálculo del estadístico de contraste.

Estadístico de contraste	Código R
$Z_o = \frac{\hat{p} - p_o}{\sqrt{\frac{p_o(1-p_o)}{n}}}$	<code>x=muestra de X</code> <code>n=length(x)</code> <code>p=proporción muestral</code> <code>po=valor a comprobar en la hipótesis</code> > <code>zo=(p-po)/sqrt(po*(1-po)/n)</code>

En lugar de calcular la fórmula, se puede usar la función `prop.test` de R.

Contraste de hipótesis	Código R
$H_0: p = p_o$ vs $H_A: p < p_o$	<code>> prop.test(ng, n, p=po, alternative="less", conf.level=0.95, correct=FALSE)</code>
$H_0: p = p_o$ vs $H_A: p > p_o$	<code>> prop.test(ng, n, p=po, alternative="greater", conf.level=0.95, correct=FALSE)</code>
$H_0: p = p_o$ vs $H_A: p \neq p_o$	<code>> prop.test(ng, n, p=po, alternative="two.sided", conf.level=0.95, correct=FALSE)</code>

Donde n es el tamaño de la muestra, y n_g es el tamaño del grupo de datos de la muestra para el que se calcula la proporción, es decir la proporción muestral sería n_g/n .

5. Ejemplo: Contraste de hipótesis sobre la proporción de estudiantes que tienen móvil Android

ENUNCIADO

Mediante una encuesta, se sabe que 50 de una muestra de 74 estudiantes de un curso de la asignatura Estadística del Grado en Ingeniería en Sistemas de Información de la Universidad de Alcalá tienen teléfonos móviles con sistema operativo Android. La población es de 108 estudiantes matriculados. Responder a la siguiente pregunta:

¿Se puede afirmar con una confianza del 95% que más de la mitad de los alumnos matriculados en la asignatura tienen teléfono Android?

- Resolverlo comprobando la región de aceptación
- Resolverlo calculando el p-valor aplicando fórmulas
- Resolverlo usando la función `prop.test()`

SOLUCIÓN

Se formula el siguiente contraste de hipótesis, ya que la mitad es el 50%, es decir una proporción de 0.5:

$$H_0: p = 0.5 \text{ vs } H_A: p > 0.5$$

a) Resolverlo comprobando la región de aceptación

La región de aceptación es $(-\infty, z_\alpha)$.

```
> so=encuesta$SO
> (n=length(so))
[1] 74
> so.android=so[so=="Android"]
> (ng=length(so.android))
[1] 50
> (p=ng/n)
[1] 0.6756757
> po=0.5
> alfa=0.05

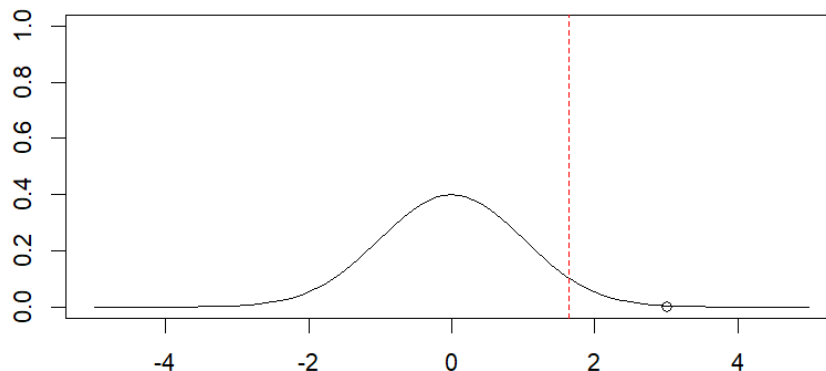
> (zo=(p-po)/sqrt(po*(1-po)/n))
[1] 3.022439

> (z.alfa=qnorm(1-alfa,0,1))
[1] 1.644854
> (zo<z.alfa)
[1] FALSE
```

Como z_0 (3.02) está fuera de la región de aceptación $(-\infty, 1.64)$, entonces se rechaza la hipótesis nula (H_0), y se acepta la hipótesis alternativa (H_A), por lo que la respuesta es **SÍ se puede afirmar con una confianza del 95% que la proporción poblacional de alumnos con Android es superior al 50%**.

Podemos hacer la comprobación visualmente.

```
> plot(z0,0, xlim = c(-5,5), ylim=c(0,1))  
> curve(dnorm(x,0,1),-5,5, add = TRUE)  
> abline(v=z.alfa, col="red", lty=2)
```



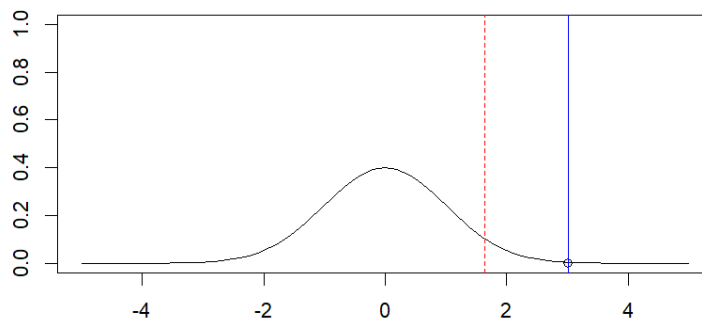
En el diagrama, el área bajo la función desde la línea hacia la izquierda es $1-\alpha$ (0.95), mientras que el área hacia la derecha es α (0.055). La región de aceptación es desde $-\infty$ hasta la línea roja, y z_0 es el punto dibujado. Se comprueba visualmente que z_0 está fuera de la región.

b) Resolverlo calculando el p-valor aplicando fórmulas

```
> (p.valor=1-pnorm(z0,0,1))  
[1] 0.001253735
```

Como el p-valor (0.001) es menor que α (0.05), se rechaza la hipótesis H_0 , como en el apartado anterior.

En el mismo diagrama del apartado anterior podemos dibujar una línea vertical azul continua en el punto z_0 .



En el diagrama, el área bajo la función desde la línea roja discontinua hacia la derecha es alfa (0.05), mientras que el área desde la línea azul continua hacia la derecha es el p-valor (0.001). Se observa visualmente que p-valor < alfa.

c) Resolverlo usando la función `prop.test()`

```
> prop.test(ng, n, p=po, alternative="greater", conf.level=0.95, correct=FALSE)
  1-sample proportions test without continuity correction
data:  ng out of n, null probability po
X-squared = 9.1351, df = 1, p-value = 0.001254
alternative hypothesis: true p is greater than 0.5
95 percent confidence interval:
 0.5813442 1.0000000
sample estimates:
      p
0.6756757
```

El p-valor es el mismo que en el apartado anterior.

6. Contraste de hipótesis sobre la diferencia de medias

Por simplicidad se utilizará sólo el método del p-valor usando las funciones `z.test` y `t.test` de R.

Si x e y son dos vectores con las muestras de las poblaciones X e Y independientes, entonces se utilizarán las funciones indicadas en la siguiente tabla.

Varianzas conocidas	Contraste de hipótesis	Función R
No $n, m \geq 30$	$H_0: \mu_X = \mu_Y$ vs $H_A: \mu_X < \mu_Y$	<code>z.test(x, y, alternative="less", sigma.x=sd(x), sigma.y=sd(y), conf.level=valor)</code>
	$H_0: \mu_X = \mu_Y$ vs $H_A: \mu_X > \mu_Y$	<code>z.test(x, y, alternative="greater", sigma.x=sd(x), sigma.y=sd(y), conf.level=valor)</code>
	$H_0: \mu_X = \mu_Y$ vs $H_A: \mu_X \neq \mu_Y$	<code>z.test(x, y, alternative="two.sided", sigma.x=sd(x), sigma.y=sd(y), conf.level=valor)</code>
No $n, m < 30$ $\sigma_X \neq \sigma_Y$	$H_0: \mu_X = \mu_Y$ vs $H_A: \mu_X < \mu_Y$	<code>t.test(x, y, alternative="less", conf.level=valor, var.equal=FALSE)</code>
	$H_0: \mu_X = \mu_Y$ vs $H_A: \mu_X > \mu_Y$	<code>t.test(x, y, alternative="greater", conf.level=valor, var.equal=FALSE)</code>
	$H_0: \mu_X = \mu_Y$ vs $H_A: \mu_X \neq \mu_Y$	<code>t.test(x, y, alternative="two.sided", conf.level=valor, var.equal=FALSE)</code>
No $n, m < 30$ $\sigma_X = \sigma_Y$	$H_0: \mu_X = \mu_Y$ vs $H_A: \mu_X < \mu_Y$	<code>t.test(x, y, alternative="less", conf.level=valor, var.equal=TRUE)</code>
	$H_0: \mu_X = \mu_Y$ vs $H_A: \mu_X > \mu_Y$	<code>t.test(x, y, alternative="greater", conf.level=valor, var.equal=TRUE)</code>
	$H_0: \mu_X = \mu_Y$ vs $H_A: \mu_X \neq \mu_Y$	<code>t.test(x, y, alternative="two.sided", conf.level=valor, var.equal=TRUE)</code>

7. Ejemplo: Contraste de hipótesis sobre la diferencia de medias de calificaciones de estudiantes del turno de mañana y del turno de tarde (muestras grandes e independientes)

ENUNCIADO

Mediante una encuesta, se saben las calificaciones de acceso a la universidad de 42 estudiantes del turno de mañana y de 32 estudiantes del turno de tarde de un curso de la asignatura Estadística del Grado en Ingeniería en Sistemas de Información de la Universidad de Alcalá, de una población de 52 matriculados en el turno de mañana y 56 en el turno de tarde. Si se supone que las calificaciones tienen una distribución Normal, responder a la siguiente pregunta:

- ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de los alumnos de mañana es diferente a la media de notas de acceso de los alumnos de tarde?
- ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de los alumnos de mañana es mayor que la media de notas de acceso de los alumnos de tarde?
- ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de los alumnos de mañana es menor que la media de notas de acceso de los alumnos de tarde?

SOLUCIÓN

Primero hay que crear los vectores con los datos de las dos muestras.

```
> encuesta=read.csv2("encuesta.csv")  
  
> encuesta.mañana=encuesta[(encuesta$GRUPO=="A1")|(encuesta$GRUPO=="A2"),]  
> encuesta.tarde=encuesta[(encuesta$GRUPO=="B1")|(encuesta$GRUPO=="B2"),]  
> notaM=encuesta.mañana$NOTA  
> notaT=encuesta.tarde$NOTA
```

a) ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de los alumnos de mañana es diferente a la media de notas de acceso de los alumnos de tarde?

Se formula el contraste:

$$H_0: \mu_{Mañana} = \mu_{Tarde} \text{ vs } H_A: \mu_{Mañana} \neq \mu_{Tarde}$$

```
> z.test(notaM, notaT, alternative="two.sided", sigma.x=sd(notaM),  
sigma.y=sd(notaT), conf.level=0.95)
```

```
Two-sample z-Test  
data: notaM and notaT  
z = 2.9889, p-value = 0.0028  
alternative hypothesis: true difference in means is not equal to 0  
95 percent confidence interval:  
 0.2597101 1.2491589  
sample estimates:  
mean of x mean of y  
 8.348810  7.594375
```

Como el p-valor (0.0028) es menor que alfa (0.05), se rechaza la Hipótesis H_0 y se acepta H_A , por lo que **Sí se puede afirmar con una confianza del 95% que la media de notas de acceso de los alumnos de mañana es diferente a la media de notas de acceso de los alumnos de tarde.**

b) ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de los alumnos de mañana es mayor que la media de notas de acceso de los alumnos de tarde?

Se formula el contraste:

$$H_0: \mu_{Mañana} = \mu_{Tarde} \text{ vs } H_A: \mu_{Mañana} > \mu_{Tarde}$$

```
> z.test(notaM, notaT, alternative="greater", sigma.x=sd(notaM),  
sigma.y=sd(notaT), conf.level=0.95)
```

```
Two-sample z-Test  
data: notaM and notaT  
z = 2.9889, p-value = 0.0014  
alternative hypothesis: true difference in means is greater than 0  
95 percent confidence interval:  
 0.3392487          NA  
sample estimates:  
mean of x mean of y  
 8.348810  7.594375
```

Como el p-valor (0.0014) es menor que alfa (0.05), se rechaza la Hipótesis H_0 y se acepta H_A , por lo que **Sí se puede afirmar con una confianza del 95% que la media de notas de acceso de los alumnos de mañana es mayor que la media de notas de acceso de los alumnos de tarde.**

c) ¿Se puede afirmar con una confianza del 95% que la media de notas de acceso de los alumnos de mañana es menor que la media de notas de acceso de los alumnos de tarde?

Se formula el contraste:

$$H_0: \mu_{\text{Mañana}} = \mu_{\text{Tarde}} \text{ VS } H_A: \mu_{\text{Mañana}} < \mu_{\text{Tarde}}$$

```
> z.test(notaM, notaT, alternative="less", sigma.x=sd(notaM), sigma.y=sd(notaT),  
conf.level=0.95)
```

```
Two-sample z-Test  
data: notaM and notaT  
z = 2.9889, p-value = 0.9986  
alternative hypothesis: true difference in means is less than 0  
95 percent confidence interval:  
NA 1.16962  
sample estimates:  
mean of x mean of y  
8.348810 7.594375
```

Como el p-valor (0.9986) es mayor que alfa (0.05), no se puede rechazar la Hipótesis H_0 y, por tanto, se rechaza H_A , por lo que **NO se puede afirmar con una confianza del 95% que la media de notas de acceso de los alumnos de mañana es menor que la media de notas de acceso de los alumnos de tarde.**

8. Contraste de hipótesis sobre la diferencia de proporciones

Por simplicidad se utilizará sólo el método del p-valor usando la función `prop.test` de R.

Si x e y son dos vectores con las muestras de las poblaciones X e Y independientes, entonces se utilizarán las funciones indicadas en la siguiente tabla.

Contraste de hipótesis	Función R
$H_0: p_1 = p_2$ vs $H_A: p_1 < p_2$	<code>prop.test(c (ng1,ng2), c (n1,n2), alternative="less", conf.level=valor, correct=FALSE)</code>
$H_0: p_1 = p_2$ vs $H_A: p_1 > p_2$	<code>prop.test(c (ng1,ng2), c (n1,n2), alternative="greater", conf.level=valor, correct=FALSE)</code>
$H_0: p_1 = p_2$ vs $H_A: p_1 \neq p_2$	<code>prop.test(c (ng1,ng2), c (n1,n2), alternative="two.sided", conf.level=valor, correct=FALSE)</code>

n_1 y n_2 son los tamaños de las muestras, y ng_1 y ng_2 son los tamaños de los grupos de datos de las muestras para los que se calculan las proporciones, es decir, las proporciones muestrales serían $p_1=ng_1/n_1$ y $p_2=ng_2/n_2$.

9. Ejemplo: Contraste de hipótesis sobre la diferencia de proporciones de estudiantes que tienen móvil Android en el turno de mañana y en el turno de tarde (muestras grandes e independientes)

ENUNCIADO

Mediante una encuesta se sabe que, en una muestra de 42 estudiantes del turno de mañana de la asignatura Estadística, 27 tienen teléfonos móviles con sistema operativo Android, y que, en una muestra de 32 estudiantes del turno de tarde, 23 tienen también teléfono Android. Responder a las siguientes preguntas:

- ¿Se puede afirmar con una confianza del 95% que la proporción de alumnos con Android en el turno de mañana es diferente a la proporción de alumnos con Android en el turno de tarde?
- ¿Se puede afirmar con una confianza del 95% que la proporción de alumnos con Android en el turno de mañana es mayor que la proporción de alumnos con Android en el turno de tarde?
- ¿Se puede afirmar con una confianza del 95% que la proporción de alumnos con Android en el turno de mañana es menor que la proporción de alumnos con Android en el turno de tarde?

SOLUCIÓN

Primero creamos variables con el número de alumnos total y los que tiene Android en cada turno:

```
> soM=encuesta.mañana$SO
> soT=encuesta.tarde$SO

> (nM=length(soM))
[1] 42
> (nT=length(soT))
[1] 32

> (nMA=length(soM[soM=="Android"]))
[1] 27
> (nTA=length(soT[soT=="Android"]))
[1] 23
```

- a) ¿Se puede afirmar con una confianza del 95% que la proporción de alumnos con Android en el turno de mañana es diferente a la proporción de alumnos con Android en el turno de tarde?

Se formula el contraste:

$$H_0: p_{Mañana} = p_{Tarde} \text{ vs } H_A: p_{Mañana} \neq p_{Tarde}$$

```
> prop.test(c(nMA,nTA), c(nM,nT), alternative="two.sided", conf.level=0.95,  
correct=FALSE)
```

```
2-sample test for equality of proportions without continuity correction  
data: c(nMA, nTA) out of c(nM, nT)  
X-squared = 0.47737, df = 1, p-value = 0.4896  
alternative hypothesis: two.sided  
95 percent confidence interval:  
-0.2886515 0.1368658  
sample estimates:  
prop 1 prop 2  
0.6428571 0.7187500
```

Como el p-valor (0.4896) es mayor que alfa (0.05), no se puede rechazar la Hipótesis H_0 y, por tanto se rechaza H_A , por lo que **NO se puede afirmar con una confianza del 95% que la proporción de alumnos con Android en el turno de mañana es diferente a la proporción de alumnos con Android en el turno de tarde.**

b) ¿Se puede afirmar con una confianza del 95% que la proporción de alumnos con Android en el turno de mañana es mayor que la proporción de alumnos con Android en el turno de tarde?

Se formula el contraste:

$$H_0: p_{\text{Mañana}} = p_{\text{Tarde}} \text{ vs } H_A: p_{\text{Mañana}} > p_{\text{Tarde}}$$

```
> prop.test(c(nMA,nTA), c(nM,nT), alternative="greater", conf.level=0.95,  
correct=FALSE)
```

```
2-sample test for equality of proportions without continuity correction  
data: c(nMA, nTA) out of c(nM, nT)  
X-squared = 0.47737, df = 1, p-value = 0.7552  
alternative hypothesis: greater  
95 percent confidence interval:  
-0.2544456 1.0000000  
sample estimates:  
prop 1 prop 2  
0.6428571 0.7187500
```

Como el p-valor (0.7552) es mayor que alfa (0.05), no se puede rechazar la Hipótesis H_0 y, por tanto se rechaza H_A , por lo que **NO se puede afirmar con una confianza del 95% que la proporción de alumnos con Android en el turno de mañana es mayor que la proporción de alumnos con Android en el turno de tarde.**

c) ¿Se puede afirmar con una confianza del 95% que la proporción de alumnos con Android en el turno de mañana es menor que la proporción de alumnos con Android en el turno de tarde?

Se formula el contraste:

$$H_0: p_{\text{Mañana}} = p_{\text{Tarde}} \text{ vs } H_A: p_{\text{Mañana}} < p_{\text{Tarde}}$$

```
> prop.test(c(nMA,nTA), c(nM,nT), alternative="less", conf.level=0.95,  
correct=FALSE)
```

```
2-sample test for equality of proportions without continuity correction  
data: c(nMA, nTA) out of c(nM, nT)  
X-squared = 0.47737, df = 1, p-value = 0.2448  
alternative hypothesis: less  
95 percent confidence interval:  
-1.0000000 0.1026599  
sample estimates:  
prop 1 prop 2  
0.6428571 0.7187500
```

Como el p-valor (0.2448) es mayor que alfa (0.05), no se puede rechazar la Hipótesis H_0 y, por tanto se rechaza H_A , por lo que **NO se puede afirmar con una confianza del 95% que la proporción de alumnos con Android en el turno de mañana es menor que la proporción de alumnos con Android en el turno de tarde.**

10. Ejercicios propuestos

- 1) Mediante una encuesta, se sabe el tiempo en minutos del viaje a la Escuela Politécnica de 74 estudiantes de un curso de la asignatura Estadística del Grado en Ingeniería en Sistemas de Información de la Universidad de Alcalá, de una población de 108 matriculados. Responder a las siguientes preguntas:
 - a. ¿Se puede afirmar con una confianza del 95% que la media del tiempo de viaje de toda la población es inferior o igual a 50 minutos?
 1. Resolverlo comprobando la región de no rechazo.
 2. Resolverlo calculando el p-valor aplicando fórmulas
 3. Resolverlo calculando el p-valor usando la función `z.test()` del paquete BSDA de R.