

Ejemplo de solución ejercicio práctico N°6

ANOVA para muestras correlacionadas

Enunciado

Un equipo de investigación del área de interacción humano-información está estudiando si el área temática y el nivel de dificultad del problema de información influyen en el tiempo (en segundos) que le toma a una persona en formular una consulta de búsqueda para resolver dicho problema.

Para ello, han reclutado a un grupo voluntario de participantes, asignados aleatoriamente a distintos grupos. Cada participante debe resolver tres problemas de información con diferentes niveles de dificultad: baja, media y alta. A su vez, cada grupo debe resolver problemas relacionados a una temática diferente. Los datos recolectados contemplan las siguientes variables:

Variable	Descripción
id	Identificador único de cada participante.
area	Área temática de los problemas que cada participante debe responder. Variable categórica con los niveles Arquitectura, Biología, Computación, Economía, Física, Leyes, Literatura, Matemáticas, Música, Pedagogía, Psicología, Química.
dificultad	Nivel de dificultad del problema resuelto. Variable categórica con los niveles Baja, Media y Alta.
tiempo	Tiempo, en segundos, que toma a cada participante formular la consulta

En este momento, equipo de investigación busca determinar si existen diferencias en el tiempo que tardan las personas en formular consultas para problemas con diferentes niveles de dificultad en el área de matemáticas.

En esta pregunta se pide inferir acerca de medias de una variable numérica (tiempo) medida en condiciones distintas (niveles de dificultad) para una misma área de conocimiento, lo que correlaciona las mediciones. Luego se requiere usar un **procedimiento ANOVA para muestras correlacionadas**. Las hipótesis serían:

H₀: no hay diferencia en los tiempos requeridos para formular una consulta asociada a un problema de información en el área de las matemáticas al considerar niveles de dificultad bajo (μ_B), medio (μ_M) y alto (μ_A), es decir $\mu_{(B-M)} = \mu_{(B-A)} = \mu_{(M-A)} = 0$.

H_A: hay diferencia en los tiempos requeridos para formular consultas asociadas a problemas de información en el área de las matemáticas con diferentes niveles de dificultad, es decir $\exists i, j \in \{B, M, A\} : \mu_{(i-j)} \neq 0$.

Obtengamos la muestra de datos (desde el archivo disponible para el ejercicio práctico anterior) que debemos utilizar.

```
src_dir <- "~/Downloads"
src_basename <- "EP06 Datos.csv"
src_file <- file.path(src_dir, src_basename)
datos <- read.csv(file = src_file, stringsAsFactors = TRUE)

# Seleccionamos datos de interés
datos_largos <- datos %>%
  filter(area == "Matemáticas") %>%
  select(id, dificultad, tiempo) %>%
  droplevels()
datos_largos[["id"]] <- factor(datos_largos[["id"]])
datos_largos[["dificultad"]] <- factor(datos_largos[["dificultad"]],
  levels = c("Baja", "Media", "Alta"))
```

Procedemos a verificar las condiciones para asegurar que podemos aplicar el procedimiento para muestras correlacionadas con validez.

La variable dependiente corresponde a tiempo que, como vimos, se mide en una escala continua de intervalos iguales.

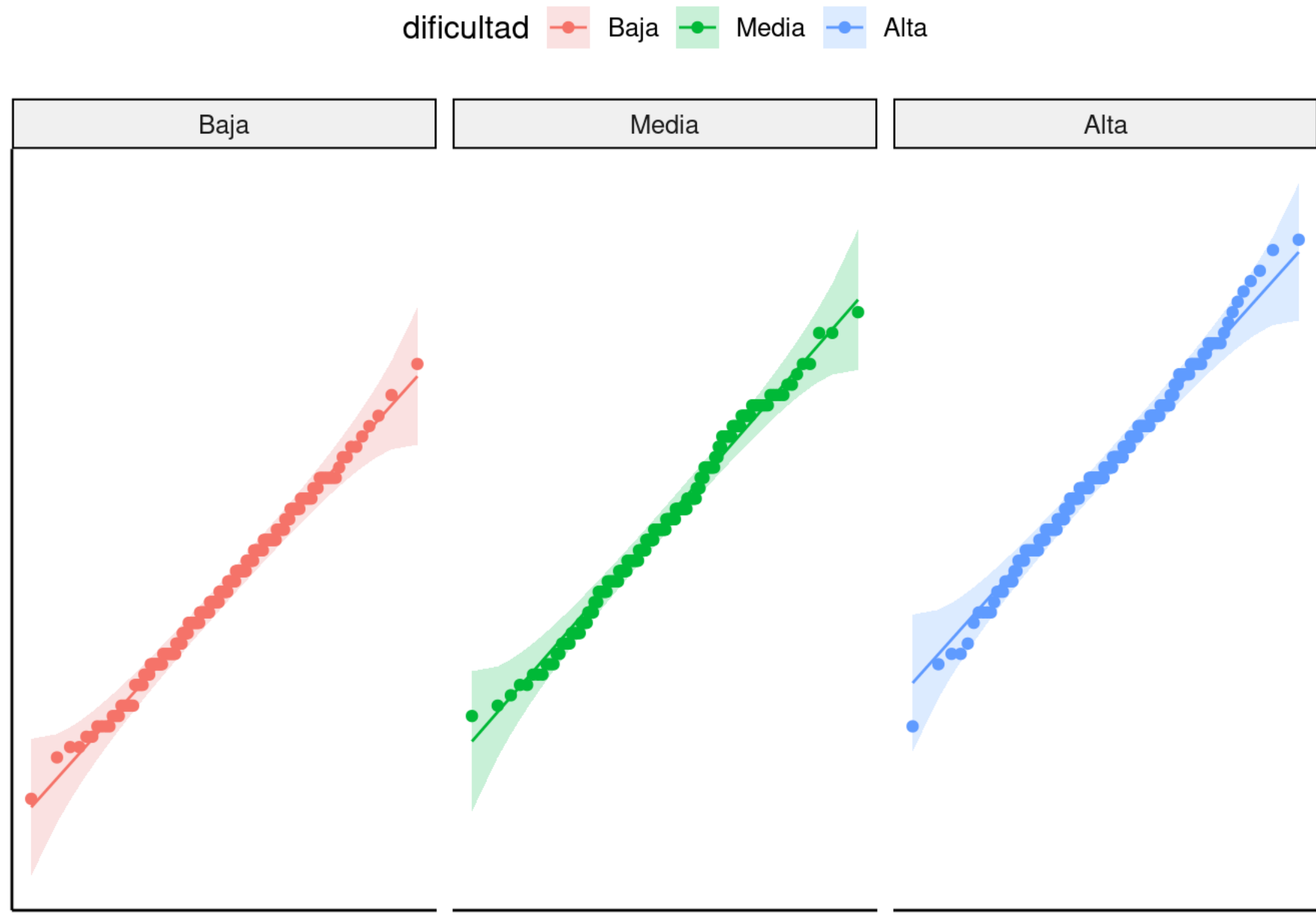
Por otro lado, los tríos de observaciones son independientes entre sí, pues provienen de personas diferentes que fueron elegidos de manera aleatoria.

Revisemos ahora la condición de normalidad por medio de un gráfico Q-Q.

```
library(ggpubr)

g <- ggqqplot(datos_largos,
  x = "tiempo",
  y = "dificultad",
  color = "dificultad")

g <- g + facet_wrap(~ dificultad)
g <- g + rremove("x.ticks") + rremove("x.text")
g <- g + rremove("y.ticks") + rremove("y.text")
g <- g + rremove("axis.title")
print(g)
```



El gráfico sugiere que los datos siguen una distribución cercana a la normal, puesto que se encuentran dentro de la región aceptable del gráfico Q-Q y no se observan patrones no aleatorios, aunque se observa cierta desviación en el extremo superior de las preguntas con dificultad alta.

En cuanto a la condición de esfericidad, se posterga su discusión hasta ver el resultado de la prueba de Mauchly efectuada por `ezAnova()`.

Así, vamos a proceder con el procedimiento ANOVA para muestras correlacionadas considerando un nivel de significación de 0,025 para compensar la posible desviación de la normalidad observada en el gráfico Q-Q.

```
library(ez)

alfa <- 0.025

omnibus <- ezANOVA(
  data = datos_largos,
  dv = tiempo,
  within = dificultad,
  wid = id,
  type = 3
)

print(omnibus)
```

```
$ANOVA
      Effect DFn DFd      F      p p<. $\alpha$       ges
2 dificultad   2 398 114.8477 4.213847e-40 * 0.2827608

$`Mauchly's Test for Sphericity`
      Effect      W      p p<. $\alpha$ 
2 dificultad 0.9849307 0.2224141

$`Sphericity Corrections`
      Effect      GGe      p[GG] p[GG]<. $\alpha$       HFe      p[HF] p[HF]<. $\alpha$ 
2 dificultad 0.9851544 1.505931e-39 * 0.9949307 6.509552e-40 *
```

Podemos ver que la prueba de esfericidad de Mauchly resulta no significativa con 97,5% de confianza ($W = 0,985; p = 0,222$), por lo que se falla en rechazar la hipótesis nula de esta prueba. Así, debemos concluir que **no hay suficiente evidencia estadística para descartar que se cumple la condición de esfericidad** en estos datos.

Interpretemos este resultado ómnibus.

El procedimiento ANOVA correlacionado resultó significativo ($F(2, 398) = 114,848; p < 0,001$). En consecuencia, con 97,5% de confianza, rechazamos la hipótesis nula en favor de la hipótesis alternativa y concluimos que hay diferencias en el tiempo requerido para formular consultas asociadas a un problema de información de en el área de las matemáticas dificultad con diferentes niveles de complejidad (baja, media y alta).

Puesto que el procedimiento ómnibus encuentra diferencias estadísticamente significativas, es necesario realizar un procedimiento post-hoc. Puesto que no requerimos hacer contrastes adicionales, usaremos la prueba HSD de Tukey (haciendo uso de un modelo mixto y de la estimación medias marginales, implementadas en R en los paquetes `nlme` y `emmeans` respectivamente).

```
library(emmeans)
library(nlme)

mixto <- lme(tiempo ~ dificultad, data = datos, random = ~1 | id)
medias <- emmeans(mixto, "dificultad")
post_hoc <- pairs(medias, adjust = "tukey")
conf_int <- confint(post_hoc, level = 1 - alfa)

print(post_hoc)
```

```
contrast estimate SE df t.ratio p.value
Alta - Baja      6.244 0.238 4398  26.232 <.0001
Alta - Media     5.338 0.238 4398  22.428 <.0001
Baja - Media     -0.905 0.238 4398   -3.804 0.0004
```

Degrees-of-freedom method: containment
P value adjustment: tukey method for comparing a family of 3 estimates

```
print(conf_int)
```

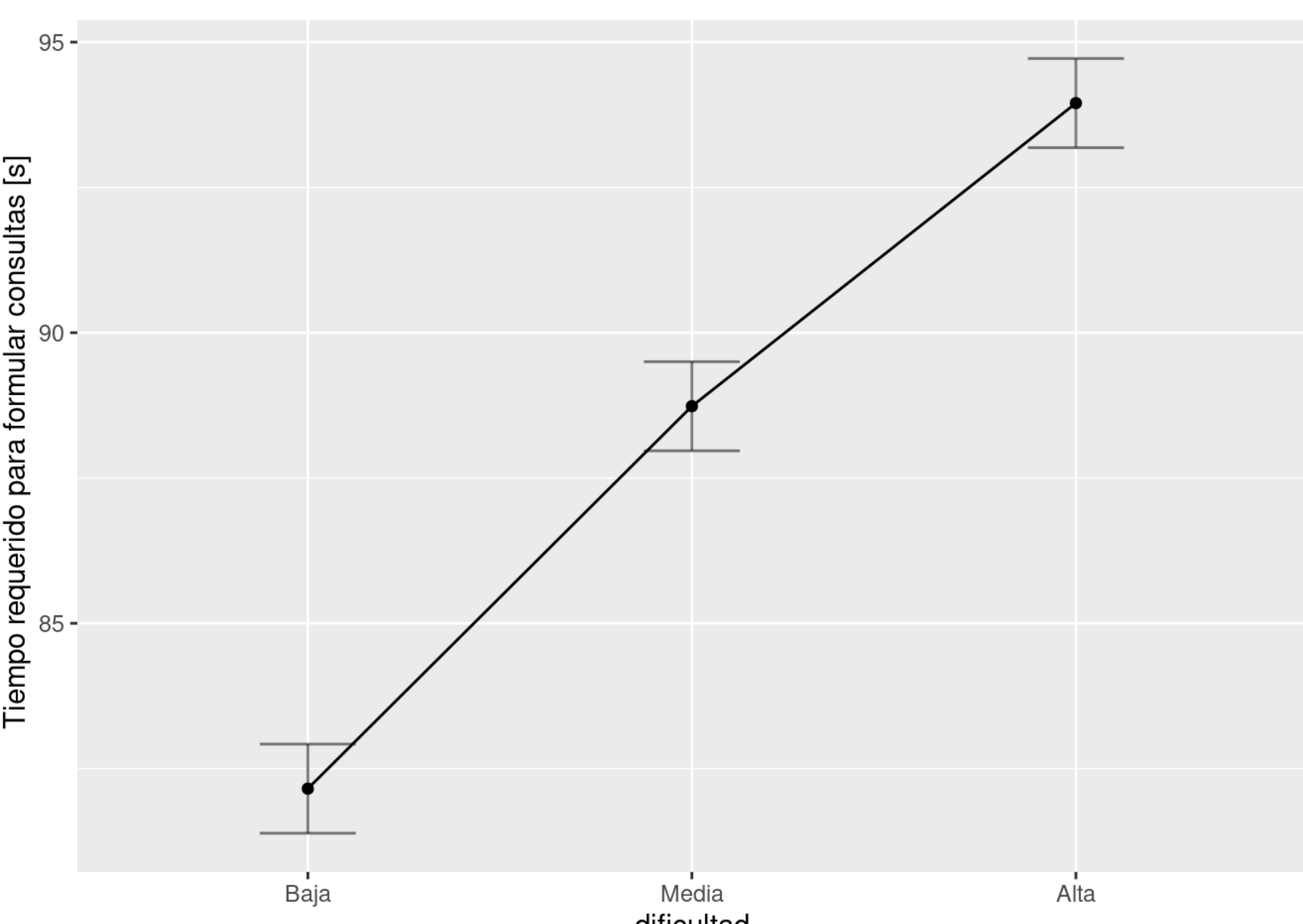
```
contrast estimate SE df lower.CL upper.CL
Alta - Baja      6.244 0.238 4398      5.62      6.864
Alta - Media     5.338 0.238 4398      4.72      5.958
Baja - Media     -0.905 0.238 4398     -1.53     -0.285
```

Degrees-of-freedom method: containment
Confidence level used: 0.975
Conf-level adjustment: tukey method for comparing a family of 3 estimates

Veamos si estos resultados coincide con el efecto (que tiene la variable independiente `dificultad` en la variable dependiente `tiempo`) encontrado en el procedimiento ANOVA para muestras correlacionadas.

```
ezp <- ezPlot(data = datos_largos, dv = tiempo, wid = id,
  within = dificultad, x = dificultad,
  y_lab = "Tiempo requerido para formular consultas [s]"
)

print(ezp)
```



Vemos que el gráfico del efecto coincide bien con los resultados de la prueba post-hoc. Redactemos la conclusión.

El análisis post-hoc usando el método de la diferencia honestamente significativa de Tukey indica que en, el área de las matemáticas, el tiempo requerido por una persona para formular consultas aumenta con el nivel de dificultad del problema de información (Alta-Baja: 97,5% CI=[5.620, 6.864], $t(4.398) = 26,232, p < 0,001$; Alta-Media: 97,5% CI=[4.720, 5.958], $t(4.398) = 22.428, p < 0,001$; Media-Baja: 97,5% CI=[0.285, 0.905], $t(4.398) = 3.804, p < 0,001$).