

UNIVERSIDAD DE EL SALVADOR.
FACULTAD MULTIDISCIPLINARIA DE OCCIDENTE.
DEPARTAMENTO DE MATEMATICAS.



PRACTICA 7. SEMANA 9

CARRERA:
LICENCIATURA EN ESTADISTICA.

ASIGNATURA:
ANALISIS ESTADISTICO CON EL PAQUETE R

DOCENTE:
JAIME ISAAC PEÑA

PRESENTADO POR:
NELSON DE JESUS MAGAÑA GODINEZ

FECHA:
20 SEPTIEMBRE DE 2022

1 DISEÑO ESTADÍSTICO DE EXPERIMENTOS

1.1 Introducción al Diseño Estadístico de Experimentos

En la práctica 6 hemos descrito métodos de inferencias sobre la media y la varianza de una población y de dos poblaciones. En esta práctica 7 ampliamos dichos métodos a más de dos poblaciones e introducimos algunos aspectos elementales del Diseño Estadístico de Experimentos y del Análisis de la Varianza.

El diseño estadístico de experimentos incluye un conjunto de técnicas de análisis y un método de construcción de modelos estadísticos que, conjuntamente, permiten llevar a cabo el proceso completo de planificar un experimento para obtener datos apropiados, que puedan ser analizados con métodos estadísticos, con objeto de obtener conclusiones válidas y objetivas.

El análisis de la varianza o abreviadamente ANOVA (del inglés analysis of variance) es un procedimiento estadístico que permite dividir la variabilidad observada en componentes independientes que pueden atribuirse a diferentes causas de interés. Es una técnica estadística para comparar más de dos grupos, es decir un método para comparar más de dos tratamientos y la variable de estudio o variable respuesta es numérica.

En esta práctica presentamos el Diseño Completamente Aleatorio con efectos fijos y con efectos aleatorios, el Diseño en Bloques Completos Aleatorizados, Diseño en Bloques Incompletos Balanceados (BIB), el Diseño en Cuadrados Latinos, el Diseño en Cuadrados Greco-Latinos, el Diseño en Cuadrados de Jouden, el Diseño Bifactorial de efectos fijos y el Diseño Trifactorial de efectos fijos.

1.1.1 Diseño Completamente Aleatorio con efectos fijos (Diseño unifactorial de efectos fijos)

El primer diseño que presentamos es el diseño completamente aleatorio de efectos fijos y la técnica estadística es el análisis de la varianza de una vía o un factor. La descripción del diseño así como la terminología subyacente la vamos a introducir mediante el siguiente supuesto práctico.

Supuesto práctico 1

La contaminación es uno de los problemas ambientales más importantes que

afectan a nuestro mundo. En las grandes ciudades, la contaminación del aire se debe a los escapes de gases de los motores de explosión, a los aparatos domésticos de la calefacción, a las industrias, . . . El aire contaminado nos afecta en nuestro vivir diario, manifestándose de diferentes formas en nuestro organismo. Con objeto de comprobar la contaminación del aire en una determinada ciudad, se ha realizado un estudio en el que se han analizado las concentraciones de monóxido de carbono (CO) durante cinco días de la semana (lunes, martes, miércoles, jueves y viernes).

En el ejemplo disponemos de una colección de 40 unidades experimentales y queremos estudiar el efecto de las concentraciones de monóxido de carbono en 5 días distintos. Es decir, estamos interesados en contrastar el efecto de un solo factor, que se presenta con cinco niveles, sobre la variable respuesta.

Nos interesa saber si las concentraciones medias de monóxido de carbono son iguales en los cinco días de la semana, para ello realizamos el siguiente contraste de hipótesis:

$$\mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5 = \mu$$

$$\mu_i \neq \mu_j \text{ para algún } i \neq j$$

Es decir, contrastamos que no hay diferencia en las medias de los cinco tratamientos frente a la alternativa de que al menos una media difiere de otra.

Variable respuesta: Concentración de CO.

Factor: Día de la semana que tiene cinco niveles. Es un factor de efectos fijos ya que viene decidido qué niveles concretos se van a utilizar (5 días de la semana).

Modelo equilibrado: Los niveles de los factores tienen el mismo número de elementos (8 elementos).

Tamaño del experimento: Número total de observaciones, en este caso 40 unidades experimentales. El problema planteado se modeliza a través de un diseño unifactorial totalmente aleatorizado de efectos fijos equilibrado.

Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos realizarlo directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R.

En este caso lo hacemos en un archivo de texto:

```
contaminacion<-read.table("Datos.txt", header = TRUE)
```

Se puede realizar de dos formas:

Transformar la variable referente a los niveles del factor fijo como factor

```
contaminacion$Dia<-factor(contaminacion$Dia)
contaminacion$Dia

## [1] Lunes      Lunes      Lunes      Lunes      Lunes      Lunes      Lunes
## [8] Lunes      Martes     Martes     Martes     Martes     Martes     Martes
## [15] Martes     Martes     Miercoles  Miercoles  Miercoles  Miercoles  Miercoles
## [22] Miercoles  Miercoles  Miercoles  Jueves     Jueves     Jueves     Jueves
## [29] Jueves     Jueves     Jueves     Jueves     Viernes    Viernes    Viernes
## [36] Viernes    Viernes    Viernes    Viernes    Viernes
## Levels: Jueves Lunes Martes Miercoles Viernes
```

Para calcular la tabla ANOVA primero hacemos uso de la función "aov" de la siguiente forma:

```
mod <- aov(Concentracion ~ Dia, data = contaminacion)
mod

## Call:
## aov(formula = Concentracion ~ Dia, data = contaminacion)
##
## Terms:
##              Dia Residuals
## Sum of Squares 119416.4 219710.4
## Deg. of Freedom      4      35
##
## Residual standard error: 79.23029
## Estimated effects may be unbalanced
```

se puede mostrar un resumen de los resultados con la funcion "summary"

```
summary(mod)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Dia              4 119416    29854   4.756 0.0036 **
## Residuals       35 219710     6277
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Si el valor de F es mayor que uno quiere decir que hay un efecto positivo del factor día. Se observa que el P-valor (Sig.) tiene un valor de 0.003524, que es menor que el nivel de significación 0.05. Por lo tanto, hemos comprobado estadísticamente que estos cinco grupos son distintos. Es decir, existen diferencias significativas en las concentraciones medias de monóxido de carbono entre los cinco días de la semana. Por lo tanto no se puede rechazar la hipótesis alternativa que dice que al menos dos grupos son diferentes, pero ¿Cuáles son esos grupos? ¿Los cinco grupos son distintos o sólo alguno de ellos? Pregunta que

resolveremos más adelante mediante los contrastes de comparaciones múltiples.

2. En la expresión del comando "aov" indicar el factor

```
mod1 <- aov(Concentracion ~ factor(Dia), data = contaminacion)
summary(mod1)

##              Df Sum Sq Mean Sq F value Pr(>F)
## factor(Dia)   4 119416   29854   4.756 0.0036 **
## Residuals    35 219710    6277
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

También se puede utilizar el comando "anova" y no es necesario el comando "summary"

```
mod2 <- anova(lm(Concentracion ~ factor(Dia), data=contaminacion))
mod2

## Analysis of Variance Table
##
## Response: Concentracion
##              Df Sum Sq Mean Sq F value    Pr(>F)
## factor(Dia)   4 119416 29854.1   4.7558 0.003598 **
## Residuals    35 219710   6277.4
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Los datos pueden venir dados en diferentes formatos:

1. Caso en el que los datos se muestran de forma que se analiza la contaminación con cada uno de los días de la semana (de lunes a viernes). Como se muestra a continuación

```
contaminacion1 <- read.table("Supuesto1.txt", header = TRUE)
contaminacion1

##   Lunes Martes Miercoles Jueves Viernes
## 1   420    450      355    321     238
## 2   390    390      462    254     255
## 3   480    430      286    412     366
## 4   430    521      238    368     389
## 5   440    320      344    340     198
## 6   324    360      423    258     256
## 7   450    342      123    433     248
## 8   460    423      196    489     324
```

En primer lugar apilaremos las columnas, para ello utilizamos el comando "stack" de la siguiente forma

```
tats <- stack(contaminacion1)
tats
```

##	values	ind
## 1	420	Lunes
## 2	390	Lunes
## 3	480	Lunes
## 4	430	Lunes
## 5	440	Lunes
## 6	324	Lunes
## 7	450	Lunes
## 8	460	Lunes
## 9	450	Martes
## 10	390	Martes
## 11	430	Martes
## 12	521	Martes
## 13	320	Martes
## 14	360	Martes
## 15	342	Martes
## 16	423	Martes
## 17	355	Miercoles
## 18	462	Miercoles
## 19	286	Miercoles
## 20	238	Miercoles
## 21	344	Miercoles
## 22	423	Miercoles
## 23	123	Miercoles
## 24	196	Miercoles
## 25	321	Jueves
## 26	254	Jueves
## 27	412	Jueves
## 28	368	Jueves
## 29	340	Jueves
## 30	258	Jueves
## 31	433	Jueves
## 32	489	Jueves
## 33	238	Viernes
## 34	255	Viernes
## 35	366	Viernes
## 36	389	Viernes
## 37	198	Viernes
## 38	256	Viernes
## 39	248	Viernes

```
## 40      324   Viernes
```

Nos muestra dos columnas:

- La primera columna: **values** nos muestra los valores de la variable respuesta. En este caso la contaminación
- La segunda columna: **ind** nos muestra los diferentes tratamientos

Podemos realizar el Análisis de la varianza utilizando el comando **anova**

```
anova(lm(values ~ ind, data = tats))

## Analysis of Variance Table
##
## Response: values
##           Df Sum Sq Mean Sq F value    Pr(>F)
## ind           4 119484  29871.1     4.775 0.003518 **
## Residuals    35 218949   6255.7
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Los datos vienen dados de la siguiente forma:

Lunes: 420, 390, 480, 430, 440, 324, 450, 460

Martes: 450, 390, 430, 521, 320, 360, 342, 423

Miércoles: 355, 462, 286, 238, 344, 423, 123, 196

Jueves: 321, 254, 412, 368, 340, 258, 433, 489

Viernes: 238, 255, 366, 389, 198, 256, 248, 324

Se crean cinco vectores, cada uno de ellos representando la contaminación con un tratamiento.

```
Lunes= c(420, 390, 480, 430, 440, 324, 450, 460)
Martes=c(450, 390, 430, 521, 320, 360, 342, 423)
Miercoles<-c(355, 462, 286, 238, 344, 423, 123, 196)
Jueves = c(321, 254, 412, 368, 340, 258, 433, 489)
Viernes<-c(238, 255, 366, 389, 198, 256, 248, 324)
```

Acontinuación creamos un data.frame para poder resolver el ANOVA

```
datos <- data.frame(Lunes, Martes, Miercoles, Jueves, Viernes)
datos
```

##	Lunes	Martes	Miercoles	Jueves	Viernes
## 1	420	450	355	321	238
## 2	390	390	462	254	255
## 3	480	430	286	412	366
## 4	430	521	238	368	389
## 5	440	320	344	340	198
## 6	324	360	423	258	256
## 7	450	342	123	433	248
## 8	460	423	196	489	324

De esta forma hemos creado una nueva base de datos que hemos llamado "datos". Para resolver el ANOVA tenemos primero que apilar las columnas con el comando "stack"

```
datos1 <- stack(datos)
datos1
```

##	values	ind
## 1	420	Lunes
## 2	390	Lunes
## 3	480	Lunes
## 4	430	Lunes
## 5	440	Lunes
## 6	324	Lunes
## 7	450	Lunes
## 8	460	Lunes
## 9	450	Martes
## 10	390	Martes
## 11	430	Martes
## 12	521	Martes
## 13	320	Martes
## 14	360	Martes
## 15	342	Martes
## 16	423	Martes
## 17	355	Miercoles
## 18	462	Miercoles
## 19	286	Miercoles
## 20	238	Miercoles
## 21	344	Miercoles
## 22	423	Miercoles
## 23	123	Miercoles
## 24	196	Miercoles
## 25	321	Jueves


```
## 26    254    Jueves
## 27    412    Jueves
## 28    368    Jueves
## 29    340    Jueves
## 30    258    Jueves
## 31    433    Jueves
## 32    489    Jueves
## 33    238    Viernes
## 34    255    Viernes
## 35    366    Viernes
## 36    389    Viernes
## 37    198    Viernes
## 38    256    Viernes
## 39    248    Viernes
## 40    324    Viernes
```

Recordemos el anova del caso anterior

```
anova(lm(values~ind, data = datos1))

## Analysis of Variance Table
##
## Response: values
##           Df Sum Sq Mean Sq F value    Pr(>F)
## ind         4 119484  29871.1    4.775 0.003518 **
## Residuals   35 218949   6255.7
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

3. Los datos se muestran en un solo vector que tiene todos los datos de la contaminación tanto si se ha medido el lunes, el martes, el miércoles, el jueves o el viernes

```
contaminacion2 <- c(Lunes, Martes, Miercoles, Jueves, Viernes)
contaminacion2

## [1] 420 390 480 430 440 324 450 460 450 390 430 521 320 360 342 423 355 462 286
## [20] 238 344 423 123 196 321 254 412 368 340 258 433 489 238 255 366 389 198 256
## [39] 248 324
```

Este vector está formado por los 40 datos que podemos comprobarlo con el comando **length**

```
length(contaminacion2)

## [1] 40
```

Para realizar el ANOVA, ya tenemos los datos de la variable respuesta y a continuación tenemos que crear el factor tratamiento, para ello vamos a utilizar la función generador de niveles, `gl`, y le decimos que nos genere 5 niveles que son los cinco tratamientos, cada uno repetido 8 veces con un total de 40 datos y para identificar que nivel es cada uno, creamos las etiquetas Lunes, Martes, Miercoles, Jueves y Viernes.

```
trat <- gl(5, 8, 40, labels = c("Lunes", "Martes",
                                "Miercoles", "Jueves",
                                "Viernes"))

trat

## [1] Lunes    Lunes    Lunes    Lunes    Lunes    Lunes    Lunes
## [8] Lunes    Martes   Martes   Martes   Martes   Martes   Martes
## [15] Martes   Martes   Miercoles Miercoles Miercoles Miercoles Miercoles
## [22] Miercoles Miercoles Miercoles Jueves    Jueves    Jueves    Jueves
## [29] Jueves    Jueves    Jueves    Jueves    Viernes   Viernes   Viernes
## [36] Viernes   Viernes   Viernes   Viernes   Viernes
## Levels: Lunes Martes Miercoles Jueves Viernes
```

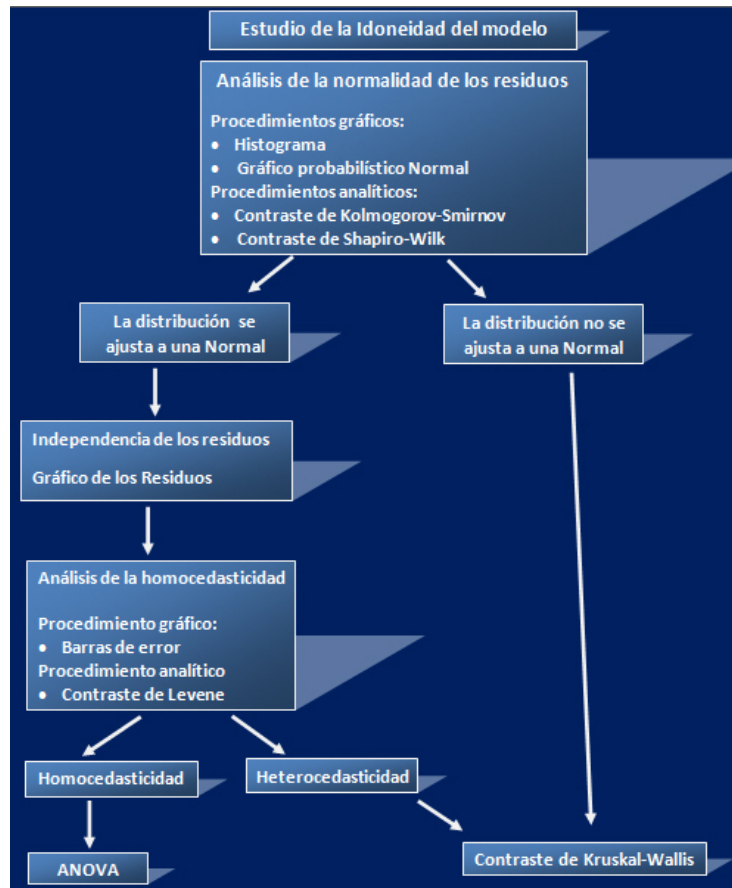
```
anova(lm(contaminacion2 ~ trat))

## Analysis of Variance Table
##
## Response: contaminacion2
##           Df Sum Sq Mean Sq F value    Pr(>F)
## trat         4 119484  29871.1    4.775 0.003518 **
## Residuals   35  218949   6255.7
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El modelo que hemos propuesto hay que validarlo, para ello hay que comprobar si se verifican las hipótesis básicas del modelo, es decir, si las perturbaciones son variables aleatorias independientes con distribución normal de media 0 y varianza constante (homocedasticidad).

1.2 Estudio de la Idoneidad del modelo

Como hemos dicho anteriormente, validar el modelo propuesto consiste en estudiar si las hipótesis básicas del modelo están o no en contradicción con los datos observados. Es decir si se satisfacen los supuestos del modelo: Normalidad, Independencia, Homocedasticidad. Para ello utilizamos procedimientos gráficos y analíticos.



1.3 Hipótesis de normalidad

En primer lugar, analizamos la normalidad de las concentraciones y continuamos con el análisis de la normalidad de los residuos.

Para analizar la normalidad de las concentraciones utilizamos el test de Shapiro-Wilks

```

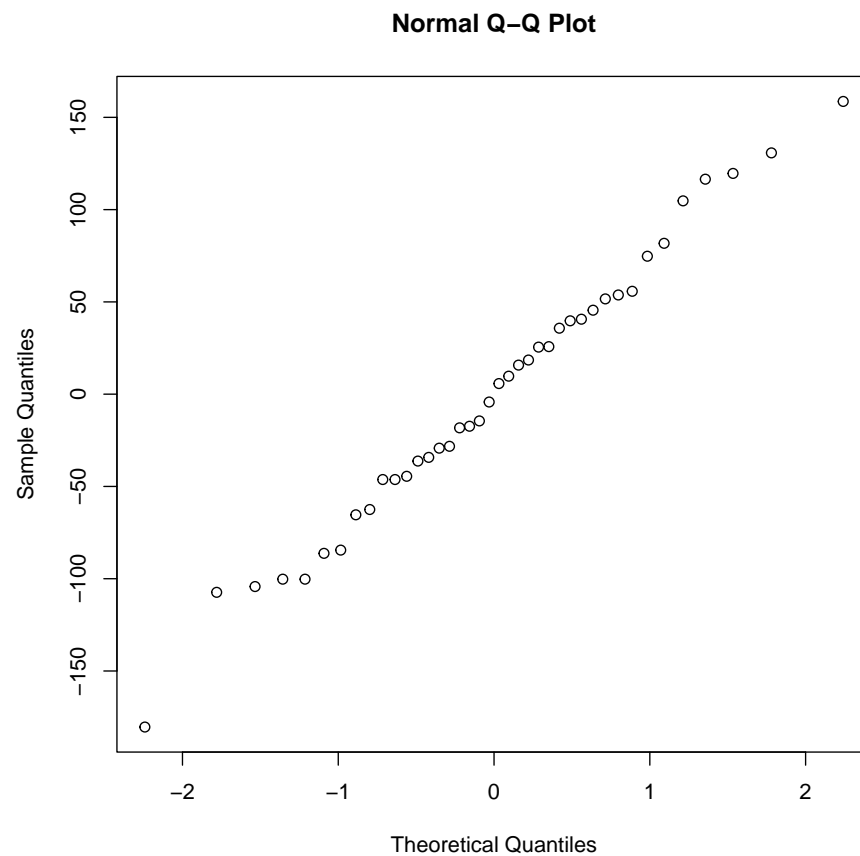
shapiro.test(mod$residuals)

##
##  Shapiro-Wilk normality test
##
## data:  mod$residuals
## W = 0.98937, p-value = 0.966
  
```

Observamos el contraste de Shapiro-Wilk que es adecuado cuando las muestras son pequeñas ($n \leq 50$) y es una alternativa más potente que el test de Kolmogorov-Smirnov. El p-valor es mayor que el nivel de significación del 5%, concluyendo que las muestras de las concentraciones se distribuyen de forma normal en cada día de la semana.

Podemos verlo también gráficamente con la orden "qqnorm"

```
qqnorm(mod$residuals)
```



Podemos apreciar en este gráfico que los puntos aparecen próximos a la línea diagonal. Esta gráfica no muestra una desviación marcada de la normalidad.

1.4 Hipótesis de homocedasticidad

Para comprobar la hipótesis de igualdad entre las varianzas del factor utilizamos el Test de Barlett.

```
bartlett.test(contaminacion$Concentracion, contaminacion$Dia)

##
## Bartlett test of homogeneity of variances
##
## data:  contaminacion$Concentracion and contaminacion$Dia
## Bartlett's K-squared = 5.5055, df = 4, p-value = 0.2392
```

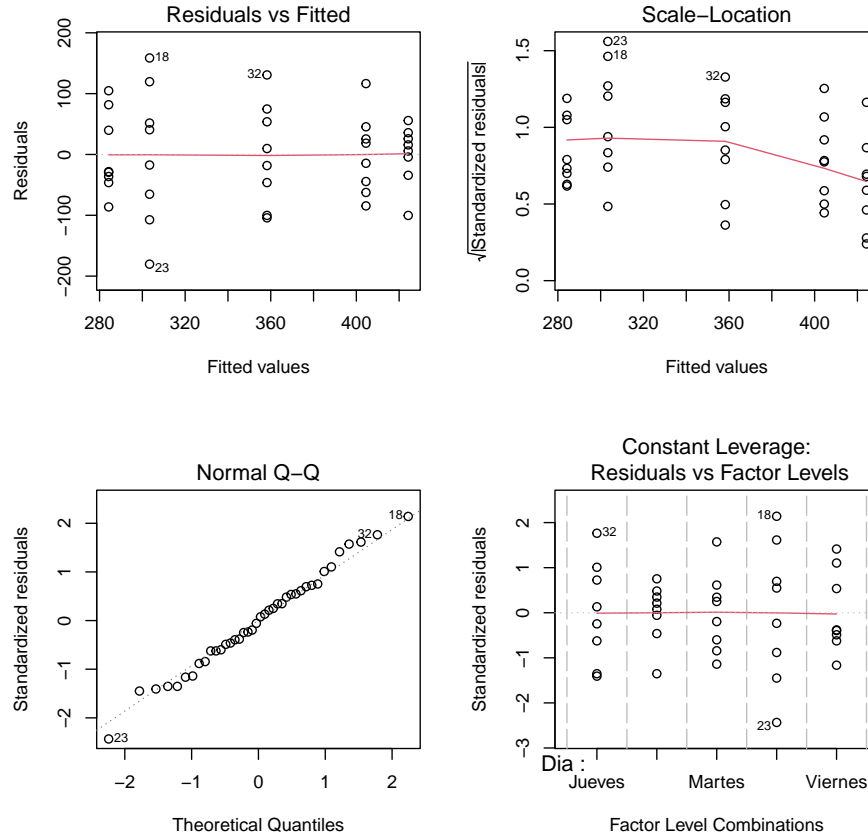
El p-valor es del 0.2402 que al ser mayor del nivel significación usual del 5% no podemos rechazar la hipótesis de igualdad de varianzas, es decir, se acepta la igualdad de varianzas en el factor.

1.5 Hipótesis de independencia

Para comprobar que se satisface el supuesto de independencia entre los residuos analizamos el gráfico de los residuos frente a los valores pronosticados o predichos por el modelo. El empleo de este gráfico es útil puesto que la presencia de alguna tendencia en el mismo puede ser indicio de una violación de dicha hipótesis. En R obtenemos varios gráficos a la vez que están incluidos en la estimación del modelo.

Para verlos de forma correcta hacemos uso de las siguientes órdenes:

```
layout(matrix(c(1,2,3,4),2,2)) # para que salgan en la misma pantalla
plot(mod)
```



En la Figura 5 se muestran cuatro gráficos, en el primero de ellos que se representan los residuos en el eje de ordenadas y los valores pronosticados en el eje de abscisas. No observamos, en dicho gráfico, ninguna tendencia sistemática que haga sospechar del incumplimiento de la suposición de independencia.

Anteriormente, hemos comprobado estadísticamente que estos cinco grupos son distintos. Es decir no se puede rechazar la hipótesis alternativa que dice que al menos dos grupos son diferentes, pero ¿Cuáles son esos grupos? ¿Los cinco grupos son distintos o sólo alguno de ellos? Pregunta que resolveremos más adelante mediante los contrastes de comparaciones múltiples.

1.6 Comparaciones múltiples

Para saber entre que parejas de días las diferencias entre concentraciones medias de CO son significativas aplicamos la prueba Post-hoc de Tukey

```
mod.tukey <- TukeyHSD(mod, ordered = TRUE)
mod.tukey

##    Tukey multiple comparisons of means
##      95% family-wise confidence level
##      factor levels have been ordered
##
## Fit: aov(formula = Concentracion ~ Dia, data = contaminacion)
##
## $Dia
##              diff          lwr          upr          p adj
## Miercoles-Viernes 19.125 -94.770937 133.0209 0.9884567
## Jueves-Viernes    74.000 -39.895937 187.8959 0.3528788
## Martes-Viernes    120.250  6.354063 234.1459 0.0342270
## Lunes-Viernes     140.000 26.104063 253.8959 0.0096790
## Jueves-Miercoles   54.875 -59.020937 168.7709 0.6408983
## Martes-Miercoles  101.125 -12.770937 215.0209 0.1019544
## Lunes-Miercoles   120.875  6.979063 234.7709 0.0329483
## Martes-Jueves      46.250 -67.645937 160.1459 0.7694367
## Lunes-Jueves       66.000 -47.895937 179.8959 0.4672516
## Lunes-Martes      19.750 -94.145937 133.6459 0.9869747
```

Esta salida nos muestra los intervalos de confianza simultáneos construidos por el método de Tukey. En la tabla se muestra un resumen de las comparaciones de cada tratamiento con los restantes. Es decir, aparecen comparadas dos a dos las cinco medias de los tratamientos.

En esta tabla, las columnas:

diff: muestra las medias de cada par

p adj: muestra los p-valores de los contrastes, que permiten conocer si la diferencia entre cada pareja de medias es significativa al nivel de significación considerado (en este caso 0.05)

lwr y upr: proporcionan los intervalos de confianza al 95% para cada diferencia.

Así por ejemplo, si comparamos la concentración media de CO del Lunes con el Martes, tenemos una diferencia entre ambas medias de 19.750, un p-valor (Sig.) de 0.9868896 no significativo puesto que la concentración de CO no difiere significativamente el lunes del martes y un intervalo de confianza con un límite inferior negativo y un límite superior positivo y por lo tanto contiene al cero de lo que también deducimos que no hay diferencias significativas entre los dos

grupos que se comparan o que ambos grupos son homogéneos.

En cambio si observamos el grupo formado por el Lunes y el Miércoles, vemos que ambos extremos del intervalo son del mismo signo y el p-valor es significativo deduciendo que si hay diferencias significativas entre ambos. Las otras comparaciones se interpretan de forma análoga.

Por lo tanto la tabla se interpreta observando los valores de p adj menores que el 5%, o si el intervalo de confianza contiene al cero.

Concluimos que se detectan diferencias significativas en las concentraciones de CO entre lunes y miércoles; lunes y viernes; martes y viernes.

Supuesto práctico 2

Los medios de cultivo bacteriológico en los laboratorios de los hospitales proceden de diversos fabricantes. Se sospecha que la calidad de estos medios de cultivo varía de un fabricante a otro. Para comprobar esta teoría, se hace una lista de fabricantes de un medio de cultivo concreto, se seleccionan aleatoriamente los nombres de cinco de los que aparecen en la lista y se comparan las muestras de los instrumentos procedentes de éstos. La comprobación se realiza colocando sobre una placa dos dosis, en gotas, de una suspensión medida de un microorganismo clásico, *Escherichia coli*, dejando al cultivo crecer durante veinticuatro horas, y determinando después el número de colonias (en millares) del microorganismo que aparecen al final del período. Se quiere comprobar si la calidad del instrumental difiere entre fabricantes.

```
bacterias <- read.table("supuesto2.txt", header = TRUE)
bacterias
```

##	Calidad	Fabricante
## 1	120	1
## 2	240	2
## 3	240	3
## 4	300	4
## 5	300	5
## 6	240	1
## 7	360	2
## 8	270	3
## 9	240	4
## 10	360	5
## 11	300	1
## 12	180	2
## 13	300	3
## 14	300	4
## 15	240	5
## 16	360	1


```
## 17      180      2
## 18      360      3
## 19      360      4
## 20      360      5
## 21      240      1
## 22      300      2
## 23      360      3
## 24      360      4
## 25      360      5
## 26      180      1
## 27      240      2
## 28      300      3
## 29      360      4
## 30      360      5
## 31      144      1
## 32      360      2
## 33      360      3
## 34      360      4
## 35      360      5
## 36      300      1
## 37      360      2
## 38      360      3
## 39      360      4
## 40      300      5
## 41      240      1
## 42      360      2
## 43      300      3
## 44      300      4
## 45      360      5
```

Para calcular la tabla ANOVA primero hacemos uso de la función 'aov' de la siguiente forma:

```
mod <- aov(Calidad ~ Fabricante, data = bacterias)
```

donde:

Calidad = nombre de la columna de las observaciones.

Fabricante = nombre de la columna en la que están representados los tratamientos.

data = data.frame en el que están guardados los datos.

```
mod

## Call:
## aov(formula = Calidad ~ Fabricante, data = bacterias)
```

```
##
## Terms:
##               Fabricante Residuals
## Sum of Squares      49561.6  152073.6
## Deg. of Freedom         1        43
##
## Residual standard error: 59.46928
## Estimated effects may be unbalanced
```

y posteriormente mostramos un resumen de los resultados con la función "summary" (verdadera tabla ANOVA):

```
summary(mod)      # TABLA ANOVA

##              Df Sum Sq Mean Sq F value    Pr(>F)
## Fabricante    1  49562   49562    14.01 0.000534 ***
## Residuals    43 152074    3537
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Esta tabla muestra los resultados del contraste planteado. El valor del estadístico de contraste es igual a 3.976 que deja a la derecha un p-valor de 0.00827, así que la respuesta dependerá del nivel de significación que se fije. Si fijamos un nivel de significación de 0.05 se concluye que hay evidencia suficiente para afirmar la existencia de alguna variabilidad entre la calidad del material de los diferentes fabricantes. Si fijamos un nivel de significación de 0.001, no podemos hacer tal afirmación.

En el modelo de efectos aleatorios no se necesitan llevar a cabo más contrastes incluso aunque la hipótesis nula sea rechazada. Es decir, en el caso de rechazar H_0 no hay que realizar comparaciones múltiples para comprobar que medias son distintas, ya que el propósito del experimento es hacer un planteamiento general relativo a las poblaciones de las que se extraen las muestras.

Diseño en Bloques Aleatorizados En los diseños estudiados anteriormente hemos supuesto que existe bastante homogeneidad entre las unidades experimentales. Pero puede suceder que dichas unidades experimentales sean heterogéneas y contribuyan a la variabilidad observada en la variable respuesta. Si en esta situación se utiliza un diseño completamente aleatorizado, no sabremos si la diferencia entre dos unidades experimentales sometidas a distintos tratamientos se debe a una diferencia real entre los efectos de los tratamientos o a la heterogeneidad de dichas unidades. Como resultado, el error experimental reflejará esta variabilidad. En esta situación se debe sustraer del error experimental la variabilidad producida por las unidades experimentales y para ello el experimentador puede formar bloques de manera que las unidades experimentales de cada bloque sean lo más homogéneas posible y los bloques entre sí sean

heterogéneos.

En el diseño en bloques Aleatorizados, primero se clasifican las unidades experimentales en grupos homogéneos, llamados bloques, y los tratamientos son entonces asignados aleatoriamente dentro de los bloques. Esta estrategia de diseño mejora efectivamente la precisión en las comparaciones al reducir la variabilidad residual.

Distinguimos dos tipos de diseños en bloques aleatorizados:

Los diseños en bloques completos aleatorizados (Todos los tratamientos se prueban en cada bloque exactamente vez). Los diseños por bloques incompletos aleatorizados (Todos los tratamientos no están representados en cada bloque, y aquellos que sí están en uno en particular se ensayan en él una sola vez).

Diseño en Bloques Completos Aleatorizados

En esta sección presentamos el diseño en Bloques Completos Aleatorizados. La palabra bloque se refiere al hecho de que se ha agrupado a las unidades experimentales en función de alguna variable extraña; aleatorizado se refiere al hecho de que los tratamientos se asignan aleatoriamente dentro de los bloques; completo implica que se utiliza cada tratamiento exactamente una vez dentro de cada bloque y el término efectos fijos se aplica a bloques y tratamientos. Es decir, se supone que ni los bloques ni los tratamientos se eligen aleatoriamente. Además una caracterización de este diseño es que los efectos bloque y tratamiento son aditivos; es decir no hay interacción entre los bloques y los tratamientos.

La descripción del diseño así como la terminología subyacente la vamos a introducir mediante el siguiente supuesto práctico.

Supuesto práctico 3

Abeto blanco, Abeto del Pirineo, es un árbol de gran belleza por la elegancia de sus formas y el exquisito perfume balsámico que destilan sus hojas y cortezas. Destilando hojas y madera se obtiene aceite de trementina muy utilizado en medicina contra torceduras y contusiones. En estos últimos años se ha observado que la producción de semillas ha descendido y con objeto de conseguir buenas producciones se proponen tres tratamientos. Se observa que árboles diferentes tienen distintas características naturales de reproducción, este efecto de las diferencias entre los árboles se debe de controlar y este control se realiza mediante bloques. En el experimento se utilizan 10 abetos, dentro de cada abeto se seleccionan tres ramas semejantes. Cada rama recibe exactamente uno de los tres tratamientos que son asignados aleatoriamente. Constituyendo cada árbol un bloque completo. Los datos obtenidos se presentan en la siguiente

tabla donde se muestra el número de semillas producidas por rama.

- Son diez Abetos en los que se aplican cuatro tratamientos distintos
- No hay ningún otro factor que pueda afectar de forma significativa a los resultados
- Los tratamientos se asignan en orden aleatorio a cada abeto
- El número de semillas observadas se muestra en la Figura 8.

1. El experimentador forma bloques de manera que las unidades experimentales de cada bloque sean lo más homogéneas posible.
2. Los bloques entre sí han de ser heterogéneos
3. Variable o factor bloque: Variable cuyo efecto sobre la variable respuesta no es directamente de interés, pero que se introduce en el experimento para obtener comparaciones homogéneas.
4. Se reduce la variabilidad residual

Distinguimos dos tipos de diseños en bloques aleatorizados:

- Los diseños en bloques completos aleatorizados (Todos los tratamientos se prueban en cada bloque exactamente vez). Los diseños por bloques incompletos aleatorizados (Todos los tratamientos no están representados en cada bloque, y aquellos que sí están en uno en particular se ensayan en él una sola vez).
- En este caso se trata de un diseño en bloques completos aleatorizados. El objetivo del estudio es comparar los tres tratamientos, por lo que se trata de un factor con tres niveles. Sin embargo, al realizar la medición sobre los distintos abetos, es posible que estos influyan sobre el número de semillas observadas. Por ello, y al no ser directamente motivo de estudio, los abetos es un factor secundario que recibe el nombre de bloque.

Nos interesa saber si los distintos tratamientos influyen en la producción de semillas, para ello realizamos el siguiente contraste de hipótesis:

$$H_0 : \tau_1 = \tau_2 = \tau_3$$
$$H_1 : \tau_i \neq \tau_j \text{ para algún } i \neq j$$

Es decir, contrastamos que no hay diferencia en las medias de los tres tratamientos frente a la alternativa de que al menos una media difiere de otra.

Pero, previamente hay que comprobar si la presencia del factor bloque (los abetos) está justificada. Para ello, realizamos el siguiente contraste de hipótesis:

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_10$$
$$H_1 : \beta_i \neq \beta_j \text{ para algún } i \neq j$$

En este caso lo hacemos en un archivo de texto:

Tenemos en cuenta que para que el ejercicio esté realizado de forma correcta los datos tienen que estar introducidos tal y como vienen en la imagen, es decir, las observaciones en una sola columna y a continuación especificado su tratamiento y su bloque correspondiente.

Para cargar los datos utilizamos la función `read.table` indicando el nombre del archivo (que debe de estar en el directorio de trabajo) e indicando además que tiene cabecera.

```
semillas<-read.table("supuesto3.txt", header = TRUE)
semillas
```

##	y	Tratamiento	Abeto
## 1	7	1	1
## 2	9	2	1
## 3	10	3	1
## 4	8	1	2
## 5	9	2	2
## 6	10	3	2
## 7	9	1	3
## 8	9	2	3
## 9	12	3	3
## 10	10	1	4
## 11	9	2	4
## 12	12	3	4
## 13	11	1	5
## 14	12	2	5
## 15	14	3	5
## 16	8	1	6
## 17	10	2	6
## 18	9	3	6
## 19	7	1	7
## 20	8	2	7
## 21	7	3	7
## 22	8	1	8
## 23	8	2	8
## 24	7	3	8
## 25	7	1	9
## 26	9	2	9
## 27	10	3	9
## 28	8	1	10
## 29	9	2	10
## 30	10	3	10

A continuación debemos transformar tanto la columna de los tratamientos como la de los bloques en un factor para poder realizar los cálculos posteriores

adecuadamente.

```
semillas$Tratamiento = factor(semillas$Tratamiento)
semillas$Tratamiento

## [1] 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3
## Levels: 1 2 3
```

```
semillas$Abeto = factor(semillas$Abeto)
semillas$Abeto

## [1] 1 1 1 2 2 2 3 3 3 4 4 4 5 5 5 6 6 6 7 7 7 8 8 8 9
## [26] 9 9 10 10 10
## Levels: 1 2 3 4 5 6 7 8 9 10
```

Para calcular la tabla ANOVA primero hacemos uso de la función "aov" de la siguiente forma:

```
mod = aov(y ~ Tratamiento + Abeto, data = semillas)
```

donde:

y es el nombre de la columna de las observaciones.

Tratamiento es el nombre de la columna en la que están representados los tratamientos.

Abeto es el nombre de la columna en la que están representados los bloques.

data = data.frame en el que están guardados los datos

```
mod

## Call:
## aov(formula = y ~ Tratamiento + Abeto, data = semillas)
##
## Terms:
##               Tratamiento Abeto Residuals
## Sum of Squares          16.2   54.8       15.8
## Deg. of Freedom           2     9         18
##
## Residual standard error: 0.936898
## Estimated effects may be unbalanced
```

y a continuación mostramos un resumen de los resultados con la función "summary" (verdadera tabla ANOVA):

```
summary(mod) # TABLA ANOVA
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## Tratamiento  2   16.2    8.100    9.228 0.00174 **
## Abeto        9   54.8    6.089    6.937 0.00026 ***
## Residuals   18   15.8    0.878
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Puesto que la construcción de bloques se ha diseñado para comprobar el efecto de una variable, nos preguntamos si ha sido eficaz su construcción. En caso afirmativo, la suma de cuadrados de bloques explicaría una parte sustancial de la suma total de cuadrados. También se reduce la suma de cuadrados del error dando lugar a un aumento del valor del estadístico de contraste experimental utilizado para contrastar la igualdad de medias de los tratamientos y posibilitando que se rechace la Hipótesis nula, mejorándose la potencia del contraste.

La construcción de bloques puede ayudar cuando se comprueba su eficacia pero debe evitarse su construcción indiscriminada. Ya que, la inclusión de bloques en un diseño da lugar a una disminución del número de grados de libertad para el error, aumenta el punto crítico para contrastar la Hipótesis nula y es más difícil rechazarla. La potencia del contraste es menor.

La Tabla *ANOVA*, muestra que:

- El valor del estadístico de contraste de igualdad de bloques, $F = 6.937$ deja a su derecha un p-valor menor que 0.001, menor que el nivel de significación del 5%, por lo que se rechaza la Hipótesis nula de igualdad de bloques. La eficacia de este diseño depende de los efectos de los bloques. Un valor grande de F de los bloques (6.937) implica que el factor bloque tiene un efecto grande. En este caso el diseño es más eficaz que el diseño completamente aleatorizado ya que si el cuadrado medio entre bloques es grande (6.089), el término residual será mucho menor (0.878) y el contraste principal de las medias de los tratamientos será más sensible a las diferencias entre tratamientos. Por lo tanto la inclusión del factor bloque en el modelo es acertada. Así, la producción de semillas depende del abeto.

Si los efectos de los bloques son muy pequeños, el análisis de bloque quizás no sea necesario y en caso extremo, cuando el valor de F de los bloques es próximo a 1, puede llegar a ser perjudicial, ya que el número de grados de libertad, $(I-1)(J-1)$, del denominador de la comparación de tratamientos es menor que el número de grados de libertad correspondiente, $IJ-I$, en el diseño completamente aleatorizado. Pero, ¿Cómo saber cuándo se puede prescindir de los bloques? La respuesta la tenemos en el valor de la F experimental de los bloques, se ha comprobado que si dicho valor es mayor que 3, no conviene prescindir de los

bloques para efectuar los contrastes.

- El valor del estadístico de contraste de igualdad de tratamiento, $F = 9.228$ deja a su derecha un p-valor de 0.002, menor que el nivel de significación del 5%, por lo que se rechaza la Hipótesis nula de igualdad de tratamientos. Así, los tratamientos influyen en el número de semillas. Es decir, existen diferencias significativas en el número de semillas entre los tres tratamientos.

El modelo que hemos propuesto hay que validarlo, para ello hay que comprobar si se verifican los cuatro supuestos expresados anteriormente.

1.6.1 Estudio de la Idoneidad del modelo

Como hemos dicho anteriormente, validar el modelo propuesto consiste en estudiar si las hipótesis básicas del modelo están o no en contradicción con los datos observados. Es decir si se satisfacen los supuestos del modelo: Normalidad, Independencia, Homocedasticidad. Para ello utilizamos procedimientos gráficos y analíticos.

En este modelo se ha supuesto otra hipótesis adicional: Aditividad de los efectos de tratamiento y bloque (no existe interacción entre tratamiento y bloque). Por lo que hay que contrastar la hipótesis de aditividad de los efectos de tratamiento y bloque.

1.6.2 Hipótesis de aditividad entre los bloques y tratamientos

La interacción entre el factor bloque y los tratamientos vamos a estudiarla analíticamente mediante el Test de Interacción de un grado de Tukey

Para realizar este test en R tenemos que utilizar la library "daewr" y dentro de ella la función "Tukey1df". De la siguiente forma:

- Primero hay que instalar el paquete **daewr**

Para ello, seleccionar **Paquetes/Instalar paquetes** y de la lista escoger **daewr**. O bien utilizar la siguiente orden

```
utils:::menuInstallPkgs()

## Error in install.packages(lib = .libPaths()[1L], dependencies =
NA, type = type): no packages were specified
```

Para realizar este contraste hay que utilizar la libray daewr, para ello realizamos la siguiente orden


```
library(daewr)

## Registered S3 method overwritten by 'DoE.base':
##   method      from
## factorize.factor conf.design

Tukey1df(semillas)

## Source      df      SS      MS      F      Pr>F
## A           2    16.2     8.1
## B           9    54.8    6.0889
## Error       18    15.8    71.1
## NonAdditivity 1    3.5573    3.5573    4.94    0.0401
## Residual    17   12.2427    0.7202
```

Puesto que el p-valor ($Pr > F$) es 1 no rechazamos la hipótesis nula de no interacción, es decir, no hay interacción entre los tratamientos aplicados y los abetos.

1.6.3 Hipótesis de Normalidad

La normalidad las vamos a comprobar analíticamente y gráficamente.

Analíticamente mediante el contraste de Shapiro-Wilk que es adecuado cuando las muestras son pequeñas ($n \leq 50$)

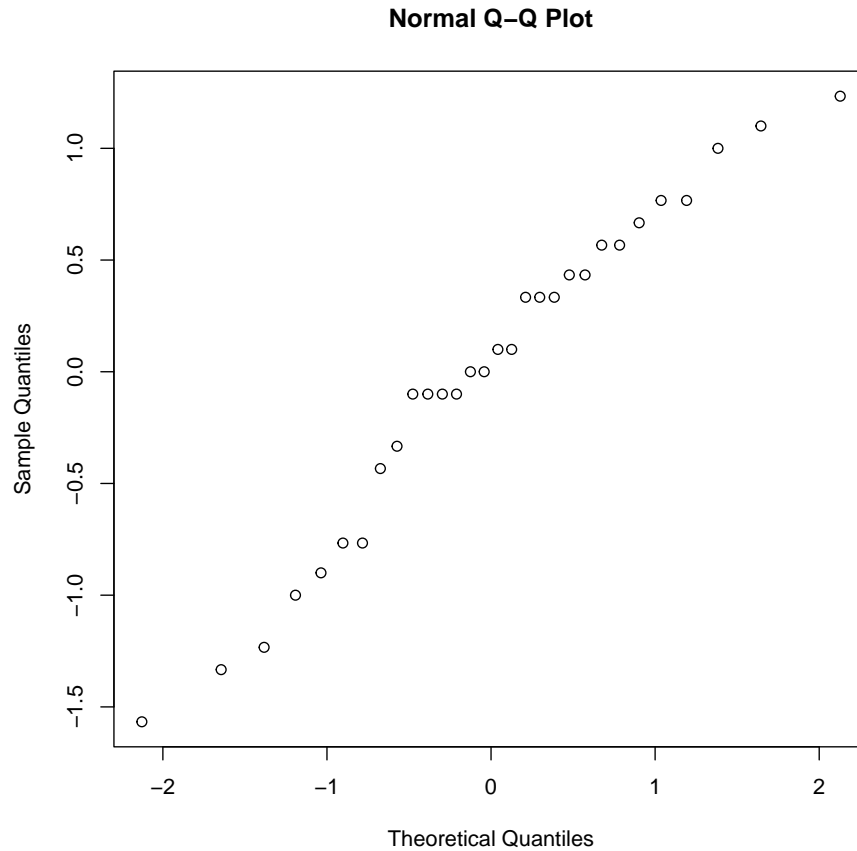
```
shapiro.test(mod$residuals)

##
## Shapiro-Wilk normality test
##
## data:  mod$residuals
## W = 0.96415, p-value = 0.3935
```

Como podemos observar tenemos un p-valor de 0.3935 que aceptaría la hipótesis de normalidad por ser mayor al 5% (nivel de significación usual).

Gráficamente mediante el gráfico probabilístico normal. Para ello utilizamos la orden "qqnorm"

```
qqnorm(mod$residuals)
```



En esta gráfica observamos que prácticamente todos los puntos se encuentran sobre la diagonal por lo tanto podemos decir que no muestra una desviación marcada de la normalidad.

Hipótesis de Homogeneidad de Varianzas Para comprobar la hipótesis de homocedasticidad utilizamos el Test de Barlett distinguiendo entre la igualdad entre varianzas del factor principal y la igualdad de varianzas del factor bloque.

En nuestro ejemplo, el test para igualdad de varianzas del factor principal sería:

```
bartlett.test(semillas$y, semillas$Tratamiento)

##
## Bartlett test of homogeneity of variances
##
## data:  semillas$y and semillas$Tratamiento
```

```
## Bartlett's K-squared = 4.1729, df = 2, p-value = 0.1241
```

El p-valor es del 0.1241 que al ser mayor del nivel significación usual del 5% no podemos rechazar la hipótesis de igualdad de varianzas en el factor principal.

De la misma manera procedemos para el factor bloque:

```
bartlett.test(semillas$y, semillas$Abeto)

##
## Bartlett test of homogeneity of variances
##
## data:  semillas$y and semillas$Abeto
## Bartlett's K-squared = 4.0723, df = 9, p-value = 0.9066
```

El p-valor es mayor que 0.05 por lo que no podemos rechazar la hipótesis de igualdad de varianzas en el factor bloque.

1.7 Hipótesis de Independencia

Comprobaremos si se satisface el supuesto de independencia entre los residuos. Para ello tenemos que representar un gráfico de los residuos tipificados frente a los pronosticados. En R obtenemos varios gráficos a la vez que están incluidos en la estimación del modelo.

Para verlos de forma correcta hacemos uso de las siguientes órdenes:

```
#layout(matrix(c(1,2,3,4),2,2))
#plot(mod)
```

Nos fijamos en el primer gráfico que representa los residuos frente a los valores ajustados y observamos que no hay ninguna tendencia sistemática. Concluimos que no hay sospechas para que se incumpla la hipótesis de independencia.

1.8 Comparaciones múltiples

Hemos probado anteriormente que se rechaza la Hipótesis nula de igualdad de tratamientos. Así, los tratamientos influyen en el número de semillas. Es decir, existen diferencias significativas en el número de semillas entre los tres tratamientos. Para saber entre que parejas de días estas diferencias son significativas aplicamos una prueba **Post-hoc**.

El contraste de Comparaciones múltiples que vamos a utilizar es el Test de

Duncan. Para poder hacer uso de él en R tenemos que instalar en primer lugar el paquete "agricolae" y dentro de él la función "duncan.test".

Destacar que este test hace las comparaciones especificándole si es para el factor principal o el factor bloque.

Comenzamos con el factor principal:

```
(duncan=duncan.test(mod, "Tratamiento" , group = T))  
  
## Error in duncan.test(mod, "Tratamiento", group = T): no se pudo encontrar la función "duncan.test"
```

En el apartado "\$groups" concluimos que los tres tratamientos difieren significativamente entre sí.

Se observa que la concentración media del número de semillas es mayor con el Tratamiento3 (10.1) y menor con el Tratamiento1 (8.3).

Para el factor bloque:

```
(duncan=duncan.test(mod, "Abeto" , group = T))  
  
## Error in duncan.test(mod, "Abeto", group = T): no se pudo encontrar la función "duncan.test"
```

Se observa que la prueba de Duncan ha agrupado los abetos 7, 8, 1, 9, 2, 6 y 10 en un mismo grupo, 1, 9, 2, 6, 10, 3 y 4, en otro grupo y un tercer está formada únicamente por el Abeto5. Inmediatamente se ve que por ejemplo el Abeto5 difiere de todos los demás, siendo en este abeto donde se produce el mayor número de semillas (12.333) y el menor en el Abeto (7.333).

2 Diseño en bloques Incompletos Aleatorizados

En los diseños en bloques Aleatorizados, puede suceder que no sea posible realizar todos los tratamientos en cada bloque. En estos casos es posible usar diseños en bloques Aleatorizados en los que cada tratamiento no está presente en cada bloque. Estos diseños reciben el nombre de diseño en bloque incompleto aleatorizado siendo uno de los más utilizados el diseño en bloque incompleto balanceado (BIB)

El diseño de bloques incompletos balanceado (BIB) compara todos los tratamientos con igual precisión.

Este diseño experimental debe verificar:

- Cada tratamiento ocurre el mismo número de veces en el diseño.
- Cada par de tratamientos ocurren juntos el mismo número de veces que cualquier otro par.

Supongamos que se tienen I tratamientos de los cuales sólo pueden experimentar K tratamientos en cada bloque ($K \leq I$). Los parámetros que caracterizan este modelo son:

- I , J y K son el número de tratamientos, el número de bloques y el número de tratamientos por bloque, respectivamente.
- R , número de veces que cada tratamiento se presenta en el diseño, es decir el número de réplicas de un tratamiento dado.
- λ , número de bloques en los que un par de tratamientos ocurren juntos.
- N , número de observaciones.

Estos parámetros deben verificar las siguientes relaciones:

$$\lambda = R \frac{K-1}{I-1}$$

donde $J \geq I$ y $N = IR = JK$

Si $J = I$ el diseño recibe el nombre de simétrico. Al igual que en el diseño en bloques completo, la asignación de los tratamientos a las unidades experimentales en cada bloque se debe realizar en forma aleatoria.

Este diseño lo estudiaremos a continuación mediante el supuesto práctico 4

Supuesto práctico 4 Se realiza un estudio para comprobar la efectividad

en el retraso del crecimiento de bacterias utilizando cuatro soluciones diferentes para lavar los envases de la leche. El análisis se realiza en el laboratorio y sólo se pueden realizar seis pruebas en un mismo día. Como los días son una fuente de variabilidad potencial, el investigador decide utilizar un diseño aleatorizado por bloques, pero al recopilar las observaciones durante seis días no ha sido posible aplicar todos los tratamientos en cada día, sino que sólo se han podido aplicar dos de las cuatro soluciones cada día. Se decide utilizar un diseño en bloques incompletos balanceado, donde $I = 4$ y $K = 2$.

Un posible diseño para estos parámetros lo proporciona la tabla correspondiente al Diseño 5 del Fichero-Adjunto, con $R = 3$, $J = 6$ y $\lambda = 1$. La disposición del diseño y las observaciones obtenidas se muestran en la siguiente tabla.

En el ejemplo:

- $N = IR = JK$. En efecto, ya que $N = 12$; $I = 4$, $J = 6$; $R = 3$ y $K = 2$.

$$\lambda = 31/3 = 1$$

El objetivo principal es estudiar la efectividad en el retraso del crecimiento de bacterias utilizando cuatro soluciones, por lo que se trata de un factor con cuatro niveles. Sin embargo, como los días son una fuente de variabilidad potencial, consideramos un factor bloque con seis niveles.

- **Variable respuesta: Número de bacterias**
- **Factor:** Soluciones que tiene cuatro niveles. Es un factor de efectos fijos ya que viene decidido qué niveles concretos se van a utilizar.
- **Bloque:** Días que tiene seis niveles. Es un factor de **efectos fijos** ya que viene decidido qué niveles concretos se van a utilizar.
- **Modelo incompleto:** Todos los tratamientos no se prueban en cada bloque.
- **Tamaño del experimento:** Número total de observaciones (12).

Podemos introducir los datos directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R.

```
bacterias = read.table("supuesto4.txt", header = TRUE)
bacterias
```

##	y	Soluciones	Dias
## 1	12	1	1
## 2	24	1	2
## 3	31	1	3
## 4	21	2	1
## 5	20	2	5
## 6	21	2	6
## 7	19	3	3
## 8	18	3	4
## 9	19	3	6
## 10	15	4	2
## 11	19	4	4
## 12	47	4	5

A continuación debemos transformar tanto la columna de los tratamientos como la de los bloques en un factor para poder realizar los cálculos posteriores adecuadamente.

```
bacterias$Soluciones = factor(bacterias$Soluciones)
bacterias$Dias = factor(bacterias$Dias)
```

Para poder analizar los datos mediante un diseño BIB debemos instalar y cargar dos paquetes de R especializados en este tipo de diseños:

```
library(daewr)
library(AlgDesign)
```

La función "BIBsize(t, k)" de la librería daewr nos permite saber si el diseño puede realizarse. Calcula los parámetros del diseño donde

- t = número de niveles del factor tratamiento.
- k = número de tratamientos por bloque.

Ejecutamos:

```
BIBsize(t = 4 , k = 2)

## Possible BIB design with b= 6 and r= 3 lambda= 1
```

El análisis de este modelo lo podemos realizar en R de dos formas:

1. Realizaremos el análisis evaluando primero el efecto de los tratamientos y después el de los bloques utilizando dos funciones
 - Para evaluar el efecto de los tratamientos, la suma de cuadrados de tratamientos debe ajustarse por bloques, por lo tanto primero se introducen los bloques y después los tratamientos.
 - Para calcular la tabla ANOVA hacemos uso de la función "aov" ($aov(y \sim A + B, data = mydataframe)$ asume suma de cuadrados tipo I) de la siguiente forma:

```
mod1 <- aov(y ~ Dias + Soluciones, data = bacterias)
```

donde:

- y = nombre de la columna de las observaciones
- Soluciones = nombre de la columna en la que están representados los tratamientos
- Dias = nombre de la columna en la que están representados los bloques
- data = data.frame en el que están guardados los datos

```

mod1

## Call:
##   aov(formula = y ~ Dias + Soluciones, data = bacterias)
##
## Terms:
##
##           Dias Soluciones Residuals
## Sum of Squares 387.6667   123.2500 396.7500
## Deg. of Freedom    5         3      3
##
## Residual standard error: 11.5
## Estimated effects may be unbalanced

```

y posteriormente mostramos un resumen de los resultados con la función "summary" (verdadera tabla ANOVA)

```

summary(mod1)

##           Df Sum Sq Mean Sq F value Pr(>F)
## Dias         5  387.7    77.53   0.586  0.720
## Soluciones   3  123.3    41.08   0.311  0.819
## Residuals    3  396.7   132.25

```

El valor del estadístico de contraste de igualdad de Soluciones, $F = 0.311$, deja a su derecha un p-valor 0.819, mayor que el nivel de significación del 5%, por lo que no se rechaza la Hipótesis Nula de igualdad de tratamientos. Por lo tanto el tipo de solución para lavar los envases de la leche no influye en el retraso del crecimiento de bacterias.

- Para evaluar el efecto de los bloques, la suma de cuadrados de bloques debe ajustarse por los tratamientos, por lo tanto primero se introducen los tratamientos y después los bloques:

```

mod2 <- aov(y ~ Soluciones + Dias, data = bacterias)
mod2

## Call:
##   aov(formula = y ~ Soluciones + Dias, data = bacterias)
##
## Terms:
##
##           Soluciones      Dias Residuals
## Sum of Squares    113.6667 397.2500 396.7500
## Deg. of Freedom         3        5      3

```



```
##  
## Residual standard error: 11.5  
## Estimated effects may be unbalanced
```

```
summary(mod2)
```

```
##           Df Sum Sq Mean Sq F value Pr(>F)  
## Soluciones   3  113.7   37.89    0.286  0.834  
## Dias         5  397.2   79.45    0.601  0.712  
## Residuals    3  396.7  132.25
```

El valor del estadístico de contraste de igualdad de Días, $F = 0.601$, deja a su derecha un p-valor 0.712, mayor que el nivel de significación del 5%, por lo que no se rechaza la Hipótesis nula de igualdad de bloques. Por lo tanto los días en los que se realiza la prueba para lavar los envases de la leche no influyen en el retraso del crecimiento de bacterias.

Con este ejemplo se ilustra el hecho de decidir si se prescinde o no de los bloques. Hay situaciones en las que, aunque los bloques no resulten significativamente diferentes no es conveniente prescindir de ellos. Pero ¿cómo saber cuándo se puede prescindir de los bloques? La respuesta la tenemos en el valor de la F de los bloques, experimentalmente se ha comprobado que si dicho valor es mayor que 3, no conviene prescindir de los bloques para efectuar los contrastes.

En esta situación si se puede prescindir del efecto de los bloques y estudiar el modelo unifactorial correspondiente, cuyo único factor es: Soluciones.

2. Realizaremos el análisis evaluando tanto para los tratamientos como para los bloques ejecutando solo una función.

Para ello necesitamos instalar y cargar el paquete “car”:

IMPORTANTE: Hemos comprobado que utilizando la versión del paquete “car” 3.0-0, encontramos un error y no permite su utilización, por lo que descargamos una versión anterior, concretamente `car_2.1-6.tar.gz`. Pinchamos en este enlace y guardamos el archivo en el escritorio

Recordad que para la utilización de un paquete es necesario instalarlo y cargarlo. Para ello:

- (a) Accedemos a Paquetes/Install package(s) from local file y elegimos el paquete descargado.
- (b) Posteriormente lo cargamos con Paquetes/Cargar paquetes

Nota: Para instalar un paquete directamente de R, procederemos de la forma siguiente

- Accedemos a la página <https://cran.r-project.org/index.html>.
- Seleccionamos **Packages**
- Seleccionamos Table of available packages, sorted by name
- Seleccionamos paquete car
- Seleccionamos Old sources car archive
- Seleccionamos el paquete car.2.1-6.tar.gz y lo guardamos en el escritorio.

Una vez instalado cargado el paquete realizamos el ANOVA

```
mod3 <- lm(y ~ Soluciones + Dias, data = bacterias)
mod3

##
## Call:
## lm(formula = y ~ Soluciones + Dias, data = bacterias)
##
## Coefficients:
## (Intercept)  Soluciones2  Soluciones3  Soluciones4      Dias2      Dias3
##      20.000      -7.000      -6.750       1.750     -1.375      8.375
##      Dias4      Dias5      Dias6
##       1.000      16.125       6.875

car::Anova(mod3, type="III")

## Anova Table (Type III tests)
##
## Response: y
##          Sum Sq Df F value Pr(>F)
## (Intercept) 533.33  1  4.0328 0.1382
## Soluciones  123.25  3  0.3106 0.8187
## Dias        397.25  5  0.6008 0.7118
## Residuals   396.75  3
```

Los resultados obtenidos coinciden con los realizados primero a los tratamientos y después a los bloques

3 Diseño de cuadrado latino

Hemos estudiado en el apartado anterior que los diseños en bloques completos aleatorizados utilizan un factor de control o variable de bloque con objeto de eliminar su influencia en la variable respuesta y así reducir el error experimental. Los diseños en cuadrados latinos utilizan dos variables de bloque para reducir el error experimental.

Un inconveniente que presentan a veces los diseños es el de requerir excesivas unidades experimentales para su realización. Un diseño en bloques completos con un factor principal y dos factores de bloque, con K_1, K_2 y K_3 niveles en cada uno de los factores, requiere $K_1 K_2 K_3$ unidades experimentales. En un experimento puede haber diferentes causas, por ejemplo de índole económico, que no permitan emplear demasiadas unidades experimentales, ante esta situación se puede recurrir a un tipo especial de diseños en bloques incompletos aleatorizados. La idea básica de estos diseños es la de fracción es decir, seleccionar una parte del diseño completo de forma que, bajo ciertas hipótesis generales, permita estimar los efectos que interesan.

Uno de los diseños en bloques incompletos aleatorizados más importante con dos factores de control es el modelo en cuadrado latino, dicho modelo requiere el mismo número de niveles para los tres factores.

En general, para K niveles en cada uno de los factores, el diseño completo en bloques aleatorizados utiliza K^2 bloques, aplicándose en cada bloque los K niveles del factor principal, resultando un total de K^3 unidades experimentales.

Los diseños en cuadrado latino reducen el número de unidades experimentales a K^2 utilizando los K^2 bloques del experimento, pero aplicando sólo un tratamiento en cada bloque con una disposición especial. De esta forma, si K fuese 4, el diseño en bloques completos necesitaría $4^3 = 64$ observaciones, mientras que el diseño en cuadrado latino sólo necesitaría $4^2 = 16$ observaciones.

Los diseños en cuadrados latinos son apropiados cuando es necesario controlar dos fuentes de variabilidad. En dichos diseños el número de niveles del factor principal tiene que coincidir con el número de niveles de las dos variables de bloque o factores secundarios y además hay que suponer que no existe interacción entre ninguna pareja de factores.

Recibe el nombre de cuadrado latino de orden K a una disposición en filas y columnas de K letras latinas, de tal forma que cada letra aparece una sola vez en cada fila y en cada columna.

En resumen, podemos decir que un diseño en cuadrado latino tiene las siguientes características:

- Se controlan tres fuentes de variabilidad, un factor principal y dos factores de bloque.
- Cada uno de los factores tiene el mismo número de niveles, K.
- Cada nivel del factor principal aparece una vez en cada fila y una vez en cada columna.
- No hay interacción entre los factores.

En el Fichero-Adjunto se muestran algunos cuadrados latinos estándares para los órdenes 3, 4, 5, 6, 7, 8 y 9.

Este diseño lo estudiaremos a continuación mediante el supuesto práctico 5

Supuesto práctico 5

Se estudia el rendimiento de un proceso químico en seis tiempos de reposo, A, B, C, D, E y F. Para ello, se consideran seis lotes de materia prima que reaccionan con seis concentraciones de ácido distintas, de manera que cada lote de materia prima en cada concentración de ácido se somete a un tiempo de reposo. Tanto la asignación de los tiempos de reposo a los lotes de materia prima, como la concentración de ácido, se hizo de forma aleatoria. Los datos del rendimiento del proceso químico se muestran en la siguiente tabla.

	Concentraciones de ácido					
Lote	1	2	3	4	5	6
Lote 1	12 A	24 B	10 C	18 D	21 E	18 F
Lote 2	21 B	26 C	24 D	16 E	20 F	21 A
Lote 3	20 C	16 D	19 E	18 F	16 A	19 B
Lote 4	22 D	15 E	14 F	19 A	27 B	17 C
Lote 5	15 E	13 F	17 A	25 B	21 C	22 D
Lote 6	17 F	11 A	12 B	22 C	14 D	20 E

El objetivo principal es estudiar la influencia de seis tiempos de reposo en el rendimiento de un proceso químico, por lo que se trata de un factor con seis niveles. Sin embargo, como los lotes de materia prima y las concentraciones son dos fuentes de variabilidad potencial, consideramos dos factores de bloque con seis niveles cada uno.

- Variable respuesta: Rendimiento.
- Factor: Tiempo de reposo que tiene seis niveles. Es un factor de efectos fijos ya que viene decidido que niveles concretos se van a utilizar.
- Bloques: Lotes y Concentraciones, ambos con seis niveles y ambos son factores de efectos fijos.
- Tamaño del experimento: Número total de observaciones (36).

Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos introducir los datos directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R.

Tenemos en cuenta que para que el ejercicio esté realizado de forma correcta los datos tienen que estar introducidos tal y como vienen en la imagen, es decir, las observaciones en una sola columna y a continuación especificado su tratamiento, su bloque y después la letra latina correspondiente.

Para cargar los datos utilizamos la función `read.table` indicando el nombre del archivo (que debe de estar en el directorio de trabajo) e indicando además que tiene cabecera.

```
latino <- read.table("supuesto5.txt", header = TRUE, dec= ",")
latino
```

##	Observaciones	Lote	Concentraciones	Tiempo_de_reposo
## 1	12	Lote1	1	A
## 2	24	Lote1	2	B
## 3	10	Lote1	3	C
## 4	18	Lote1	4	D
## 5	21	Lote1	5	E
## 6	18	Lote1	6	F
## 7	21	Lote2	1	B
## 8	26	Lote2	2	C
## 9	24	Lote2	3	D
## 10	16	Lote2	4	E
## 11	20	Lote2	5	F
## 12	21	Lote2	6	A
## 13	20	Lote3	1	C
## 14	16	Lote3	2	D
## 15	19	Lote3	3	E
## 16	18	Lote3	4	F
## 17	16	Lote3	5	A
## 18	19	Lote3	6	B
## 19	22	Lote4	1	D
## 20	15	Lote4	2	E

## 21	14 Lote4	3	F
## 22	19 Lote4	4	A
## 23	27 Lote4	5	B
## 24	17 Lote4	6	C
## 25	15 Lote5	1	E
## 26	13 Lote5	2	F
## 27	17 Lote5	3	A
## 28	25 Lote5	4	B
## 29	21 Lote5	5	C
## 30	22 Lote5	6	D
## 31	17 Lote6	1	F
## 32	11 Lote6	2	A
## 33	12 Lote6	3	B
## 34	22 Lote6	4	C
## 35	14 Lote6	5	D
## 36	20 Lote6	6	E

Para calcular la tabla ANOVA primero hacemos uso de la función "aov" de la siguiente forma:

```
mod1 <- aov(Observaciones~ Lote + Concentraciones + Tiempo_de_reposo, data = latino)
```

donde:

- Observaciones: Nombre de la columna de las observaciones
- Lote: Nombre de la columna en la que están representados los tratamientos
- Concentraciones : Nombre de la columna en la que está representado el primer factor bloque.
- Tiempo_de_reposo: Nombre de la columna en la que está representado el segundo factor bloque (letras latinas).
- data = data.frame en el que están guardados los datos

```
mod1

## Call:
##   aov(formula = Observaciones ~ Lote + Concentraciones + Tiempo_de_reposo,
##       data = latino)
##
## Terms:
```

```
##               Lote Concentraciones Tiempo_de_reposo Residuals
## Sum of Squares  99.5556          30.9429          117.8889  386.1683
## Deg. of Freedom    5              1              5          24
##
## Residual standard error: 4.011277
## Estimated effects may be unbalanced
```

y posteriormente mostramos un resumen de los resultados con la función "summary" (verdadera tabla ANOVA):

```
summary(mod1)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Lote           5   99.6   19.91   1.237  0.323
## Concentraciones 1   30.9   30.94   1.923  0.178
## Tiempo_de_reposo 5  117.9   23.58   1.465  0.238
## Residuals      24  386.2   16.09
```

Observando los valores de los p-valores, 0.281, 0.368 y 0.553; mayores respectivamente que el nivel de significación del 5%, deducimos que ningún efecto es significativo.

4 Diseño de Cuadrado Greco-Latino

El modelo en cuadrado greco-latino se puede considerar como una extensión del modelo en cuadrado latino en el que se incluye una tercera variable control o variable de bloque. En este modelo como en el diseño en cuadrado latino, todos los factores deben tener el mismo número de niveles, K , y el número de observaciones necesarias sigue siendo K^2 . Este diseño es, por tanto, una fracción del diseño completo en bloques aleatorizados con un factor principal y tres factores secundarios que requeriría K_4 observaciones.

Los cuadrados greco-latinos se obtienen por superposición de dos cuadrados latinos del mismo orden y ortogonales entre sí, uno de los cuadrados con letras latinas el otro con letras griegas. Dos cuadrados reciben el nombre de ortogonales si, al superponerlos, cada letra latina y griega aparecen juntas una sola vez en el cuadrado resultante.

En el Fichero-Adjunto se muestra una tabla de cuadrados latinos que dan lugar, por superposición de dos de ellos, a cuadrados greco-latinos. Notamos que no es posible formar cuadrados greco-latinos de orden 6.

La Tabla siguiente ilustra un cuadrado greco-latino para $K = 4$

Supuesto práctico 6

Para comprobar el rendimiento de un proceso químico en cinco tiempos de reposo, se consideran cinco lotes de materia prima que reaccionan con cinco concentraciones de ácido distintas a cinco temperaturas distintas, de manera que cada lote de materia prima con cada concentración de ácido y cada temperatura se someten a un tiempo de reposo. Tanto la asignación de los tiempos de reposo a los lotes de materia prima, como las concentraciones de ácido, y las temperaturas, se hizo de forma aleatoria. En este estudio el científico considera que tanto los lotes de materia prima, las concentraciones y las temperaturas pueden influir en el rendimiento del proceso, por lo que los considera como variables de bloque cada una con cinco niveles y decide plantear un diseño por cuadrados greco-latinos como el que muestra en la siguiente tabla.

La variable respuesta que vamos a estudiar es el rendimiento del proceso químico. El factor principal es tiempo de reposo que se presenta con cinco niveles.

- Variable respuesta: Rendimiento
- Factor: Tiempos de reposo que tiene cinco niveles. Es un factor de efectos fijos ya que viene decidido que niveles concretos se van a utilizar.
- Bloques: Lotes, Concentraciones y Temperaturas, cada uno con cinco niveles y de efectos fijos.
- Tamaño del experimento: Número total de observaciones (25).

Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos introducir los datos directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R.

En este caso lo hacemos en un archivo de texto.

Tenemos en cuenta que para que el ejercicio esté realizado de forma correcta los datos tienen que estar introducidos tal y como vienen en la imagen, es decir, las observaciones en una sola columna y a continuación especificado su tratamiento, su bloque correspondiente y después la letra latina y griega correspondiente (En este caso hemos cambiado las letras griegas como las últimas del alfabeto latino por facilidad a la hora de escribirlas).

Para cargar los datos utilizamos la función `read.table` indicando el nombre del archivo (que debe de estar en el directorio de trabajo) e indicando además que tiene cabecera.

```
greco <- read.table("supuesto6-2.txt", header = TRUE, dec= ",")
greco

##      Observaciones Lotes Concentraciones Tiempo_de_reposo Temperaturas
```


## 1	26 Lote1	1	A	Z
## 2	21 Lote1	2	B	Y
## 3	19 Lote1	3	C	X
## 4	13 Lote1	5	D	W
## 5	21 Lote1	5	E	V
## 6	22 Lote2	1	B	X
## 7	26 Lote2	2	C	W
## 8	24 Lote2	3	D	V
## 9	16 Lote2	4	E	Z
## 10	20 Lote2	5	A	Y
## 11	29 Lote3	1	C	V
## 12	26 Lote3	2	D	Z
## 13	19 Lote3	3	E	Y
## 14	18 Lote3	4	A	X
## 15	16 Lote3	5	B	W
## 16	32 Lote4	1	D	Y
## 17	15 Lote4	2	E	X
## 18	14 Lote4	3	A	W
## 19	19 Lote4	4	B	V
## 20	27 Lote4	5	C	Z
## 21	25 Lote5	1	E	W
## 22	18 Lote5	2	A	V
## 23	19 Lote5	3	B	Z
## 24	25 Lote5	4	C	Y
## 25	21 Lote5	5	D	X

A continuación debemos transformar tanto la columna de los tratamientos como la de los bloques en un factor para poder realizar los cálculos posteriores adecuadamente.

Para calcular la tabla ANOVA primero hacemos uso de la función "aov" de la siguiente forma:

```
greco <- read.table("supuesto6-2.txt", header = TRUE, dec = ",")
greco$Lote <- factor(greco$Lote)
greco$Temperaturas <- factor(greco$Temperaturas)
greco$Tiempo_de_reposo <- factor(greco$Tiempo_de_reposo)
greco$Concentraciones <- factor(greco$Concentraciones)
mod1 <- aov(Observaciones ~ Lote + Concentraciones + Tiempo_de_reposo + Temperaturas, data = greco)
mod1

## Call:
## aov(formula = Observaciones ~ Lote + Concentraciones + Tiempo_de_reposo +
## Temperaturas, data = greco)
##
## Terms:
```

```
##               Lote Concentraciones Tiempo_de_reposo Temperaturas
## Sum of Squares    9.7600      207.7607      155.0085      97.2516
## Deg. of Freedom      4          4          4          4
##               Residuals
## Sum of Squares   100.7792
## Deg. of Freedom      8
##
## Residual standard error: 3.549281
## Estimated effects may be unbalanced
```

donde:

- Observaciones: Nombre de la columna de las observaciones.
- Lote: Nombre de la columna en la que están representados los tratamientos.
- Concentraciones = Nombre de la columna en la que está representado el primer factor bloque.
- Tiempo_de_reposo = Nombre de la columna en la que está representado el segundo factor bloque (letras latinas).
- Temperaturas: Nombre de la columna en la que está representado el tercer factor bloque.
- Data: data.frame en el que están guardados los datos

y posteriormente mostramos un resumen de los resultados con la función "summary" (verdadera tabla ANOVA):

```
summary(mod1)

##               Df Sum Sq Mean Sq F value Pr(>F)
## Lote           4   9.76    2.44    0.194 0.9349
## Concentraciones 4 207.76   51.94    4.123 0.0420 *
## Tiempo_de_reposo 4 155.01   38.75    3.076 0.0825 .
## Temperaturas    4  97.25   24.31    1.930 0.1988
## Residuals      8 100.78   12.60
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Observando los valores de los p-valores, 0.150, 0.053, 0.912 y 0.020, deducimos que el único efecto significativo, al nivel de significación del 5%, es el efecto de las distintas concentraciones sobre el rendimiento del proceso químico.

5 Diseño de cuadrados de Youden

Hemos estudiado que en el diseño en cuadrado latino se tiene que verificar que los tres factores tengan el mismo número de niveles, es decir que hay el mismo número de filas, de columnas y de letras latinas. Sin embargo, puede suceder que el número de niveles disponibles de uno de los factores de control sea menor que el número de tratamientos, en este caso estaríamos ante un diseño en cuadrado latino incompleto. Estos diseños fueron estudiados por W.J. Youden y se conocen con el nombre de cuadrados de Youden.

Este diseño lo estudiaremos a continuación mediante el supuesto práctico 7.

Supuesto práctico 7

Consideremos de nuevo el experimento sobre el rendimiento de un proceso químico en el que se está interesado en estudiar seis tiempos de reposo, A, B, C, D, E y F y se desea eliminar estadísticamente el efecto de los lotes materia prima y de las concentraciones de ácido distintas. Pero supongamos que sólo se dispone de cinco tipos de concentraciones. Para analizar este experimento se decidió utilizar un cuadrado de Youden con seis filas (los lotes de materia prima), cinco columnas (las distintas concentraciones) y seis letras latinas (los tiempos de reposo). Los datos correspondientes se muestran en la siguiente tabla.

Lote	concentracion de acido				
	1	2	3	4	5
Lote 1	12	24	10	18	21
	A	B	C	D	E
Lote 2	21	26	24	16	20
	B	C	D	E	F
Lote 3	20	16	19	18	16
	C	D	E	F	A
Lote 4	22	15	14	19	27
	D	E	F	A	B
Lote 5	15	13	17	25	21
	E	F	A	B	C
Lote 6	17	11	12	22	14
	F	A	B	C	D

Tenemos en cuenta que para que el ejercicio esté realizado de forma correcta los datos tienen que estar introducidos tal y como vienen en la imagen, es decir, las observaciones en una sola columna y a continuación especificado su tratamiento, su bloque correspondiente y después la letra latina correspondiente.

Para cargar los datos utilizamos la función `read.table` indicando el nombre del archivo (que debe de estar en el directorio de trabajo) e indicando además que tiene cabecera.

```
youden <- read.table("supuesto7-1.txt", header = TRUE)
youden
```

##	Observaciones	Lote	Concentraciones	Tiempo_de_reposo
## 1	12	Lote1	1	A
## 2	24	Lote1	2	B
## 3	10	Lote1	3	C
## 4	18	Lote1	4	D
## 5	21	Lote1	5	E
## 6	21	Lote2	1	B
## 7	26	Lote2	2	C
## 8	24	Lote2	3	D
## 9	16	Lote2	4	E
## 10	20	Lote2	5	F
## 11	20	Lote3	1	C
## 12	16	Lote3	2	D
## 13	19	Lote3	3	E
## 14	18	Lote3	4	F
## 15	16	Lote3	5	A
## 16	22	Lote4	1	D
## 17	15	Lote4	2	E
## 18	14	Lote4	3	F
## 19	19	Lote4	4	A
## 20	27	Lote4	5	B
## 21	15	Lote5	1	E
## 22	13	Lote5	2	F
## 23	17	Lote5	3	A
## 24	25	Lote5	4	B
## 25	21	Lote5	5	C
## 26	17	Lote6	1	F
## 27	11	Lote6	2	A
## 28	12	Lote6	3	B
## 29	22	Lote6	4	C
## 30	14	Lote6	5	D

A continuación debemos transformar tanto la columna de los tratamientos como la de los bloques en un factor para poder realizar los cálculos posteriores

adecuadamente.

```
youden$Lote <- factor(youden$Lote)
youden$Concentraciones <- factor(youden$Concentraciones)
youden$Tiempo_de_reposo <- factor(youden$Tiempo_de_reposo)
```

Para cada factor realizamos una tabla ANOVA:

- Factor principal:

Para evaluar el efecto de los tratamientos, la suma de cuadrados de tratamientos debe ajustarse por bloques, por lo tanto primero se introducen los bloques y después los tratamientos.

Para calcular la tabla ANOVA hacemos uso de la función "aov" (asume suma de cuadrados tipo I) de la siguiente forma:

```
mod1 <- aov(Observaciones ~ Tiempo_de_reposo + Lote + Concentraciones, data = youden)
mod1

## Call:
## aov(formula = Observaciones ~ Tiempo_de_reposo + Lote + Concentraciones,
## data = youden)
##
## Terms:
##              Tiempo_de_reposo              Lote Concentraciones Residuals
## Sum of Squares              151.76667 112.73333              61.66667 282.00000
## Deg. of Freedom                  5              5                  4              15
##
## Residual standard error: 4.335897
## Estimated effects may be unbalanced
```

```
summary(mod1)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Tiempo_de_reposo  5 151.77   30.35   1.615  0.216
## Lote              5 112.73   22.55   1.199  0.356
## Concentraciones  4  61.67   15.42   0.820  0.532
## Residuals       15 282.00   18.80
```

donde:

- Observaciones: Nombre de la columna de las observaciones.
- Lote: Nombre de la columna en la que están representados los tratamientos.

- Concentraciones: Nombre de la columna en la que está representado el primer factor bloque.
- Tiempo_de_reposo: Nombre de la columna en la que está representado el segundo factor bloque (letras latinas).
- data = data.frame en el que están guardados los datos.
El p-valor, 0.532, es mayor que el nivel de significación del 5%, deducimos que el factor principal: Concentraciones no es significativo.
- Factor Bloque: Lotes.

Para evaluar el efecto del primero de los bloques, la suma de cuadrados de bloques debe ajustarse por los tratamientos, por lo tanto primero se introducen los tratamientos y después los bloques:

```
mod3 <- aov(Observaciones~ Concentraciones + Lote +Tiempo_de_reposo , data = youden)
mod3

## Call:
## aov(formula = Observaciones ~ Concentraciones + Lote + Tiempo_de_reposo,
## data = youden)
##
## Terms:
##             Concentraciones             Lote Tiempo_de_reposo Residuals
## Sum of Squares             61.66667 111.36667             153.13333 282.00000
## Deg. of Freedom              4              5              5              15
##
## Residual standard error: 4.335897
## Estimated effects may be unbalanced
```

```
summary(mod3)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Concentraciones	4	61.67	15.42	0.820	0.532
Lote	5	111.37	22.27	1.185	0.362
Tiempo_de_reposo	5	153.13	30.63	1.629	0.213
Residuals	15	282.00	18.80		

El p-valor es 0.213; mayor que el nivel de significación del 5%, deducimos que Factor Bloque:Tiempo_de_reposo no es significativo.

6 Diseños Factoriales

En muchos experimentos es frecuente considerar dos o más factores y estudiar el efecto conjunto que dichos factores producen sobre la variable respuesta. Para resolver esta situación se utiliza el Diseño Factorial.

Se entiende por diseño factorial aquel diseño en el que se investigan todas las posibles combinaciones de los niveles de los factores en cada réplica del experimento. En estos diseños, los factores que intervienen tienen la misma importancia a priori y se supone por tanto, la posible presencia de interacción. En este epígrafe vamos a considerar únicamente modelos de efectos fijos.

6.1 Diseños factoriales con dos factores

6.1.1 Modelo sin recuperacion

Supuesto práctico 8

En unos laboratorios se está investigando sobre el tiempo de supervivencia de unos animales a los que se les suministra al azar tres tipos de venenos y cuatro antídotos distintos. Se pretende estudiar si los tiempos de supervivencia de los animales varían en función de las combinaciones veneno-antídoto. Los datos que se recogen en la tabla adjunta son los tiempos de supervivencia en horas.

El objetivo principal es estudiar la influencia de tres tipos de venenos y 4 tipos de antídotos en el tiempo de supervivencia de unos determinados animales, por lo que se trata de un modelo con dos factores: el veneno (con tres niveles) y el antídoto (con cuatro niveles). La variable que va a medir las diferencias entre los tratamientos es el tiempo que sobreviven los animales. Se combinan todos los niveles de los dos factores por lo que tenemos en total doce tratamientos.

- Variable respuesta: Tiempo de supervivencia
- Factor: Tipo de veneno que tiene tres niveles. Es un factor de efectos fijos ya que viene decidido qué niveles concretos se van a utilizar.
- Factor: Tipo de antídoto que tiene cuatro niveles. Es un factor de efectos fijos ya que viene decidido qué niveles concretos se van a utilizar.
- Tamaño del experimento: Número total de observaciones (12).

Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos introducir los datos directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R.

En este caso lo hacemos en un archivo de texto:

Tenemos en cuenta que para que el ejercicio esté realizado de forma correcta los datos tienen que estar introducidos tal y como vienen en la imagen, es decir, las observaciones en una sola columna y a continuación especificado sus factores correspondientes.

Para cargar los datos utilizamos la función `read.table` indicando el nombre del archivo (que debe de estar en el directorio de trabajo) e indicando además que tiene cabecera.

```
getwd()

## [1] "C:/Users/Usuario/Desktop/respaldo/Desktop/PAQUETE R/PRACTICAS_S9"

setwd("C:/Users/Usuario/Desktop/respaldo/Desktop/PAQUETE R/PRACTICAS_S9")
factorial <- read.table("supuesto8.txt", header = TRUE)
factorial

##      Tiempo Veneno Antidoto
## 1      4.5      1      1
## 2      2.9      2      1
## 3      2.1      3      1
## 4     11.0      1      2
## 5      6.1      2      2
## 6      3.7      3      2
## 7      4.5      1      3
## 8      3.5      2      3
## 9      2.5      3      3
## 10     7.1      1      4
## 11    10.2      2      4
## 12     3.6      3      4
```

A continuación debemos transformar todas las columnas que contienen a los factores en un factor para poder realizar los cálculos posteriores adecuadamente.

```
factorial$Antidoto <- factor(factorial$Antidoto)
factorial$Veneno <- factor(factorial$Veneno)
```

Para calcular la tabla ANOVA primero hacemos uso de la función `aov` de la siguiente forma

```
mod <- aov(Tiempo ~ Veneno + Antidoto , data = factorial)
mod
```



```
## Call:
## aov(formula = Tiempo ~ Veneno + Antidoto, data = factorial)
##
## Terms:
##              Veneno Antidoto Residuals
## Sum of Squares 30.58667 39.40917 23.89333
## Deg. of Freedom      2        3        6
##
## Residual standard error: 1.995551
## Estimated effects may be unbalanced
```

y posteriormente mostramos un resumen de los resultados con la función "summary" (verdadera tabla ANOVA):

```
summary(mod)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Veneno        2  30.59   15.293    3.840 0.0844 .
## Antidoto       3  39.41   13.136    3.299 0.0995 .
## Residuals     6   23.89    3.982
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Esta Tabla ANOVA recoge la descomposición de la varianza considerando como fuente de variación los doce tratamientos o grupos que se forman al combinar los niveles de los dos factores. Mediante esta tabla se puede estudiar si varían los tiempos que sobreviven los animales en función de las combinaciones veneno-antídoto. Es decir, se pueden estudiar si existen diferencias significativas entre los tiempos medios de supervivencia con los distintos tipos de venenos y antídotos, pero no se puede estudiar si la efectividad de los antídotos es la misma para todos los venenos. Observando los p-valores, 0.084 y 0.099; mayores respectivamente que el nivel de significación del 5%, deducimos que ningún efecto es significativo. Por lo tanto, no existen diferencias en los tiempos medios de supervivencia de los animales, en función de la pareja veneno-antídoto que se les suministra.

6.2 El modelo con replicación

El modelo estadístico para este diseño es:

$$y_{ijk} = \mu + \tau_i + \beta_j + (\tau\beta)_{ij} + u_{ijk}, i = 1, 2, \dots, a, j = 1, 2, \dots, b, k = 1, \dots, r$$

donde r es el número de replicaciones y $N = abr$ es el número de observaciones.

El número de parámetros de este modelo es, como en el modelo de dos factores sin replicación, $ab + 1$ pero en este caso el número de observaciones es abr .

La descripción del diseño así como la terminología subyacente la vamos a introducir mediante el siguiente supuesto práctico.

Supuesto práctico 9

Consideremos el supuesto práctico anterior en el que realizamos dos réplicas por cada tratamiento. Los datos que se recogen en la tabla adjunta son los tiempos de supervivencia en horas de unos animales a los que se les suministra al azar tres venenos y cuatro antídotos. El objetivo es estudiar qué antídoto es el adecuado para cada veneno.

	Antídoto			
Veneno	Antídoto 1	Antídoto 2	Antídoto 3	Antídoto 4
Veneno 1	4.5	11.0	4.5	7.1
	4.3	7.2	7.6	6.2
Veneno 2	2.9	6.1	3.5	10.2
	2.3	12.4	4.0	3.8
Veneno 3	2.1	3.7	2.5	3.6
	2.3	2.9	2.2	3.3

Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos introducir los datos directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R.

En este caso lo hacemos en un archivo de texto:

tenemos en cuenta que para que el ejercicio esté realizado de forma correcta los datos tienen que estar introducidos tal y como vienen en la imagen, es decir, las observaciones en una sola columna y a continuación especificado sus factores correspondientes.

Para cargar los datos utilizamos la función `read.table` indicando el nombre del archivo (que debe de estar en el directorio de trabajo) e indicando además que tiene cabecera.

```
factorial <- read.table("supuesto8.txt", header = TRUE)
factorial
```

```
##      Tiempo Veneno Antidoto
## 1      4.5      1      1
## 2      2.9      2      1
## 3      2.1      3      1
## 4     11.0      1      2
## 5      6.1      2      2
## 6      3.7      3      2
## 7      4.5      1      3
## 8      3.5      2      3
## 9      2.5      3      3
## 10     7.1      1      4
## 11     10.2     2      4
## 12     3.6      3      4
```

A continuación debemos transformar todas las columnas que contienen a los factores en un factor para poder realizar los cálculos posteriores adecuadamente.

```
factorial$Veneno <- factor(factorial$Veneno)
factorial$Antidoto <- factor(factorial$Antidoto)
```

Para calcular la tabla ANOVA primero hacemos uso de la función "aov" de la siguiente forma:

```
mod <- aov(Tiempo ~ Veneno * Antidoto , data = factorial)
mod

## Call:
## aov(formula = Tiempo ~ Veneno * Antidoto, data = factorial)
##
## Terms:
##              Veneno Antidoto Veneno:Antidoto
## Sum of Squares  30.58667 39.40917      23.89333
## Deg. of Freedom      2      3      6
##
## Estimated effects may be unbalanced
```

```
summary(mod)

##              Df Sum Sq Mean Sq
## Veneno        2  30.59  15.293
## Antidoto       3  39.41  13.136
## Veneno:Antidoto 6  23.89   3.982
```

La Tabla ANOVA muestra las filas de Tipo_veneno, Tipo_antídoto y Tipo_veneno × Tipo_antídoto que corresponde a la variabilidad debida a los efectos de cada uno de los factores y de la interacción entre ambos.

Las preguntas que nos planteamos son: ¿Los venenos son igual de peligrosos? ¿Y los antídotos son igual de efectivos? La efectividad de los antídotos, ¿es la misma para todos los venenos? Para responder a estas preguntas, comenzamos comprobando si el efecto de los antídotos es el mismo para todos los venenos. Para ello observamos el valor del estadístico ($F_{exp} = 0.761$) que contrasta la hipótesis correspondiente a la interacción entre ambos factores ($H_0 : (\tau\beta)_{ij} = 0$). Dicho valor deja a la derecha un Sig. = 0.614, mayor que el nivel de significación 0.05. Por lo tanto la interacción entre ambos factores no es significativa y debemos eliminarla del modelo. Construimos de nuevo la Tabla ANOVA en la que sólo figurarán los efectos principales

```
mod <- aov(Tiempo ~ Veneno + Antidoto , data = factorial)
mod

## Call:
## aov(formula = Tiempo ~ Veneno + Antidoto, data = factorial)
##
## Terms:
##              Veneno Antidoto Residuals
## Sum of Squares  30.58667  39.40917   23.89333
## Deg. of Freedom      2         3         6
##
## Residual standard error: 1.995551
## Estimated effects may be unbalanced
```

```
summary(mod)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Veneno         2  30.59   15.293    3.840 0.0844 .
## Antidoto        3  39.41   13.136    3.299 0.0995 .
## Residuals       6  23.89    3.982
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Esta tabla muestra dos únicas fuentes de variación, los efectos principales de los dos factores (Tipo_veneno y Tipo_antídoto), y se ha suprimido la interacción entre ambos. Se observa que el valor de la Suma de Cuadrados del error de este modelo (73.873) se ha formado con los valores de las Sumas de cuadrados del error y de la interacción del modelo anterior ($20.363 + 53.510 = 73.873$). Observando los valores de los p-valores, 0.0046 y 0.0117 asociados a los contrastes principales, se deduce que los dos efectos son significativos a un nivel de

significación del 5%. Deducimos que ni la gravedad de los venenos es la misma, ni la efectividad de los antídotos, pero dicha efectividad no depende del tipo de veneno con el que se administre ya que la interacción no es significativa.

6.3 Diseños factoriales con tres factores

Supongamos que hay a niveles para el factor A, b niveles del factor B y c niveles para el factor C y que cada réplica del experimento contiene todas las posibles combinaciones de tratamientos, es decir contiene los abc tratamientos posibles.

6.3.1 El modelo sin replicación

Supuesto práctico 10

En una fábrica de refrescos está haciendo unos estudios en la planta embotelladora. El objetivo es obtener más uniformidad en el llenado de las botellas. La máquina de llenado teóricamente llena cada botella a la altura correcta, pero en la práctica hay variación, y la embotelladora desea entender mejor las fuentes de esta variabilidad para eventualmente reducirla. En el proceso se pueden controlar tres factores durante el proceso de llenado: El % de carbonato (factor A), la presión del llenado (factor B) y el número de botellas llenadas por minuto que llamaremos velocidad de la línea (factor C). Se consideran tres niveles para el factor A (10%, 12%, 14%), dos niveles para el factor B (25psi, 30psi) y dos niveles para el factor C (200bpm, 250bpm). Los datos recogidos de la desviación de la altura objetivo se muestran en la tabla adjunta

	Presión (B)			
	25 psi		30 psi	
	Velocidad (C)		Velocidad (C)	
	200	250	200	250
% de Carbono (A)				
10	10	3	5	-1
12	11	2	5	-3
14	2	4	-3	1

La variable respuesta de este experimento es la Desviación que se produce en la altura de llenado en las botellas de refresco, siendo dichas botellas las unidades experimentales. En estas desviaciones de la altura de llenado marcada como objetivo intervienen tres factores: Porcentaje de carbono que presenta tres niveles 10%, 12% y 14%; Presión, con dos niveles 25 psi y 30 psi y Velocidad, con dos niveles 200 y 250. Los niveles de los factores han sido fijados por el experimentador, por lo que todos los factores son de efectos fijos. Se trata de un diseño trifactorial de efectos fijos, donde el número de tratamientos es $3 \cdot 2 \cdot 2 = 12$.

Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos introducir los datos directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R.

Tenemos en cuenta que para que el ejercicio esté realizado de forma correcta los datos tienen que estar introducidos tal y como vienen en la imagen, es decir, las observaciones en una sola columna y a continuación especificado sus factores correspondientes.

Para cargar los datos utilizamos la función `read.table` indicando el nombre del archivo (que debe de estar en el directorio de trabajo) e indicando además que tiene cabecera.

```
factorial <- read.table("supuesto10.txt", header = TRUE)
factorial
```

##	Altura	Carbono	Presion	Velocidad
## 1	10	10	25	200
## 2	11	12	25	200
## 3	2	14	25	200
## 4	3	10	25	250
## 5	2	12	25	250
## 6	4	14	25	250
## 7	5	10	30	200
## 8	5	12	30	200
## 9	-3	14	30	200
## 10	-1	10	30	250
## 11	-3	12	30	250
## 12	1	14	30	250

A continuación debemos transformar la tres columnas en factores para poder realizar los cálculos posteriores adecuadamente.

```
factorial$Carbono <- factor(factorial$Carbono)
factorial$Velocidad <- factor(factorial$Velocidad)
factorial$Presion <- factor(factorial$Presion)
```

Para calcular la tabla ANOVA primero hacemos uso de la función "aov" de la siguiente forma:

```
mod=aov(Altura ~ Carbono + Presion + Velocidad + Carbono*Presion + Carbono*Velocidad + Presion*Velocidad, data=factorial)

## Error in lm.fit(x, y, offset = offset, singular.ok = singular.ok, ...): NA/NaN/Inf in 'x'
```

donde:

- Altura: Nombre de la columna de las observaciones.
- Carbono: Nombre de la columna en la que está representado el primer factor.
- Presion: Nombre de la columna en la que está representado el segundo factor.
- Velocidad: Nombre de la columna en la que está representado el tercer factor
- Carbono*Presion, Carbono*Velocidad y Presion*Velocidad hace referencia a las distintas interacciones.
- data= data.frame en el que están guardados los datos

```
mod

## Call:
##   aov(formula = Tiempo ~ Veneno + Antidoto, data = factorial)
##
## Terms:
##               Veneno Antidoto Residuals
## Sum of Squares  30.58667 39.40917  23.89333
## Deg. of Freedom      2       3       6
##
## Residual standard error: 1.995551
## Estimated effects may be unbalanced
```

```
summary(mod)

##           Df Sum Sq Mean Sq F value Pr(>F)
## Veneno      2  30.59  15.293   3.840 0.0844 .
## Antidoto     3  39.41  13.136   3.299 0.0995 .
## Residuals    6   23.89   3.982
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

La Tabla ANOVA muestra las filas de Carbono, Presión, Velocidad, Carbono*Presión, Carbono*Velocidad y Presión*Velocidad que corresponden a la variabilidad debida a los efectos de cada uno de los factores y a las interacciones de orden dos entre ambos. En dicha Tabla se indica que para un nivel de significación del 5% los efectos que no son significativos del modelo planteado son las interacciones entre los factores Carbono*Presión y Presión*Velocidad ya que los p-valores correspondientes a estos efectos son 0.125 y 0.057 mayores que el nivel de significación.

Como consecuencia de este resultado, replanteamos el modelo suprimiendo en primer lugar el efecto Carbono*Presión, cuya significación es mayor. donde los efectos deben cumplir las condiciones expuestas anteriormente. Para resolverlo suprimimos la interacción Carbono*Presión. La tabla ANOVA que corresponde a este modelo es la siguiente:

```
mod <- aov(Altura~ Carbono + Presion + Velocidad + Carbono*Velocidad + Presion*Velocidad , data = factorial )

## Error in lm.fit(x, y, offset = offset, singular.ok = singular.ok,
...): NA/NaN/Inf in 'x'

summary(mod)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Veneno        2  30.59   15.293    3.840 0.0844 .
## Antidoto       3  39.41   13.136    3.299 0.0995 .
## Residuals     6   23.89    3.982
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El efecto Presión*Velocidad sigue siendo no significativo por lo que lo suprimimos del modelo y replanteamos el siguiente modelo matemático donde los efectos deben cumplir las condiciones expuestas anteriormente. Para resolverlo suprimimos la interacción Presión*Velocidad. La tabla ANOVA que corresponde a este modelo es la siguiente:

```
mod <- aov(Altura~ Carbono + Presion + Velocidad + Carbono*Velocidad, data = factorial )

## Error in lm.fit(x, y, offset = offset, singular.ok = singular.ok,
...): NA/NaN/Inf in 'x'

summary(mod)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Veneno        2  30.59   15.293    3.840 0.0844 .
## Antidoto       3  39.41   13.136    3.299 0.0995 .
## Residuals     6   23.89    3.982
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Todos los efectos de este último modelo planteado son significativos y por lo tanto es en este modelo donde vamos a realizar el estudio. Existen diferencias significativas entre los distintos porcentajes del Carbono, los dos tipos de presión, las dos velocidades de llenado y la interacción entre el porcentaje de Carbono y la Velocidad de llenado.

Supuesto práctico 11.

Consideremos el supuesto práctico anterior en el que realizamos dos réplicas por cada tratamiento. En la Tabla adjunta se muestran los datos recogidos de la desviación de la altura objetivo de las botellas de refresco. En el proceso de llenado, la embotelladora puede controlar tres factores durante el proceso: El porcentaje de carbonato (factor A) con tres niveles (10%, 12%, 14%), la presión del llenado (factor B) con dos niveles (25psi, 30psi) y el número de botellas llenadas por minuto que llamaremos velocidad de la línea (factor C) con dos niveles (200bpm, 250bpm).

	Presión (B)			
	25 psi		30 psi	
	Velocidad (C)		Velocidad (C)	
% de Carbono (A)	200	250	200	250
10	10	3	5	-1
	20	5	9	-3
12	11	2	5	-3
	9	5	4	2
14	2	4	-3	1
	-1	7	-2	3

Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos introducir los datos directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R.

En este caso lo hacemos en un archivo de texto. Tenemos en cuenta que para que el ejercicio esté realizado de forma correcta los datos tienen que estar introducidos tal y como vienen en la imagen, es decir, las observaciones en una sola columna y a continuación especificado sus factores correspondientes.

Para cargar los datos utilizamos la función `read.table` indicando el nombre del archivo (que debe de estar en el directorio de trabajo) e indicando además que tiene cabecera.

```
factorial <- read.table("supuesto11.txt", header = TRUE)
factorial

##      Altura Carbono Presion Velocidad
## 1       10       10      25       200
## 2       20       10      25       200
## 3       11       12      25       200
## 4        9       12      25       200
## 5        2       14      25       200
```

```
## 6      -1      14      25      200
## 7       3      10      25      250
## 8       5      10      25      250
## 9       2      12      25      250
## 10      5      12      25      250
## 11      4      14      25      250
## 12      7      14      25      250
## 13      5      10      30      200
## 14      9      10      30      200
## 15      5      12      30      200
## 16      4      12      30      200
## 17     -3      14      30      200
## 18     -2      14      30      200
## 19     -1      10      30      250
## 20     -3      10      30      250
## 21     -3      12      30      250
## 22      2      12      30      250
## 23      1      14      30      250
## 24      3      14      30      250
```

A continuación debemos transformar las tres columnas en factores para poder realizar los cálculos posteriores adecuadamente.

```
factorial$Carbono <- factor(factorial$Carbono)
factorial$Velocidad <- factor(factorial$Velocidad)
factorial$Presion <- factor(factorial$Presion)
```

Para calcular la tabla ANOVA primero hacemos uso de la función "aov" de la siguiente forma:

```
mod <- aov(Altura ~ Carbono + Presion +
           Velocidad + Carbono*Presion +
           Carbono*Velocidad + Presion*Velocidad +
           Carbono*Velocidad*Presion, data = factorial)
```

donde:

- Altura: Nombre de la columna de las observaciones.
- Carbono: Nombre de la columna en la que está representado el primer factor.
- Presion: Nombre de la columna en la que está representado el segundo factor.
- Velocidad: Nombre de la columna en la que está representado el tercer factor

- Carbono*Presion, Carbono*Velocidad, Presion*Velocidad y Carbono*Velocidad*Presion hace referencia a las distintas interacciones.
- data= data.frame en el que están guardados los datos

y posteriormente mostramos un resumen de los resultados con la función "summary" (verdadera tabla ANOVA):

```
summary(mod)

##              Df Sum Sq Mean Sq F value    Pr(>F)
## Carbono      2  88.08   44.04    5.683 0.018350 *
## Presion      1 150.00  150.00   19.355 0.000866 ***
## Velocidad     1  80.67   80.67   10.409 0.007270 **
## Carbono:Presion  2  14.25    7.12    0.919 0.425122
## Carbono:Velocidad 2 230.58  115.29   14.876 0.000564 ***
## Presion:Velocidad 1   1.50    1.50    0.194 0.667799
## Carbono:Presion:Velocidad 2   1.75    0.88    0.113 0.894175
## Residuals    12  93.00    7.75
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

La Tabla ANOVA muestra las filas de Carbono, Presión, Velocidad, Carbono*Presión, Carbono*Velocidad, Presión*Velocidad y Carbono*Presión*Velocidad que corresponden a la variabilidad debida a los efectos de cada uno de los factores, a las interacciones de orden dos y orden tres entre los factores. En dicha Tabla se indica que para un nivel de significación del 5% los efectos que no son significativos del modelo planteado son las interacciones entre los factores, Carbono*Presión y Presión*Velocidad y Carbono*Presión*Velocidad ya que los p-valores correspondientes a estos efectos son 0.425, 0.668 y 0.894 mayores que el nivel de significación.

Para resolverlo suprimimos la interacción Carbono*Presión*Velocidad. La tabla ANOVA que corresponde a este modelo es la siguiente:

```
mod <- aov(Altura~ Carbono + Presion +
           Velocidad + Carbono*Presion +
           Carbono*Velocidad + Presion*Velocidad,
           data = factorial)
summary(mod)

##              Df Sum Sq Mean Sq F value    Pr(>F)
## Carbono      2  88.08   44.04    6.507 0.010038 *
## Presion      1 150.00  150.00   22.164 0.000336 ***
## Velocidad     1  80.67   80.67   11.919 0.003886 **
## Carbono:Presion  2  14.25    7.12    1.053 0.375033
## Carbono:Velocidad 2 230.58  115.29   17.035 0.000178 ***
```

```
## Presion:Velocidad 1 1.50 1.50 0.222 0.645047
## Residuals 14 94.75 6.77
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Los efectos Carbono*Presión y Presión*Velocidad siguen siendo no significativos.

Para resolverlo suprimimos la interacción Presión*Velocidad. La tabla ANOVA que corresponde a este modelo es la siguiente:

```
mod <- aov(Altura~ Carbono + Presion +
           Velocidad + Carbono*Presion +
           Carbono*Velocidad, data = factorial)
summary(mod)

##              Df Sum Sq Mean Sq F value    Pr(>F)
## Carbono        2  88.08   44.04    6.864 0.007647 **
## Presion         1 150.00  150.00   23.377 0.000218 ***
## Velocidad       1  80.67   80.67   12.571 0.002935 **
## Carbono:Presion  2  14.25    7.12    1.110 0.355049
## Carbono:Velocidad 2 230.58  115.29   17.968 0.000104 ***
## Residuals      15  96.25    6.42
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El efecto Carbono*Presión sigue siendo no significativo por lo tanto lo suprimimos y replanteamos

```
mod <- aov(Altura~ Carbono + Presion +
           Velocidad + Carbono*Velocidad,
           data = factorial)
summary(mod)

##              Df Sum Sq Mean Sq F value    Pr(>F)
## Carbono        2  88.08   44.04    6.776 0.006856 **
## Presion         1 150.00  150.00   23.077 0.000166 ***
## Velocidad       1  80.67   80.67   12.410 0.002612 **
## Carbono:Velocidad 2 230.58  115.29   17.737 6.91e-05 ***
## Residuals      17 110.50    6.50
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Todos los efectos de este último modelo planteado son significativos y por lo tanto es en este modelo donde vamos a realizar el estudio. Existen diferencias significativas entre los distintos porcentajes del Carbono, los dos tipos de

presión, las dos velocidades de llenado y la interacción entre el porcentaje de Carbono y la Velocidad de llenado.

7 Ejercicios Guiados.

Ejercicio Guiado 1

Se realiza un estudio del contenido de azufre en cinco yacimientos de carbón. Se toman muestras aleatoriamente de cada uno de los yacimientos y se analizan. Los datos del porcentaje de azufre por muestra se indican en la tabla adjunta.

Yacimientos	Porcentaje de azufre
1	151 192 108 204 214 176 117
2	169 64 90 141 101 128 159 156
3	122 132 139 133 154 104 225 149 130
4	75 126 69 62 90 120 32 73
5	80 90 124 82 72 57 118 54 130

Para un nivel de significación del 5%.

1. ¿Se puede confirmar que el porcentaje de azufre es el mismo en los cinco yacimientos?
2. Si se rechaza la hipótesis nula que las medias de porcentaje de azufre en los cinco yacimientos es la misma, determinar que medias difieren entre sí utilizando el método de comparaciones múltiples de Tukey.
3. Estudiar las hipótesis de modelo: Homocedasticidad (Homogeneidad de las varianzas por grupo), Independencia y Normalidad.

Solucion: 1

El problema planteado se modeliza a través de un diseño unifactorial totalmente aleatorizado de efectos fijos no-equilibrado.

- Variable respuesta: Contenido de Azufre
- Factor: Tipo de yacimiento con cinco niveles. Es un factor de Efectos fijos ya que viene decidido qué niveles concretos se van a utilizar.
- Modelo no-equilibrado: Los niveles de los factores tienen distinto número de elementos.
- Tamaño del experimento: Número total de observaciones, en este caso 41 unidades experimentales.

Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos realizarlo directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R. En este caso lo hacemos en un archivo de texto:

En primer lugar describimos los cinco grupos que tenemos que comparar, los cinco yacimientos, la variable respuesta es el porcentaje de azufre en estos cinco yacimientos. Los yacimientos no tienen todos el mismo número de observaciones, en total tenemos 41 observaciones. La hipótesis nula es que el porcentaje de azufre es el mismo en los cinco yacimientos. . . Es decir, no hay diferencias en los porcentajes de azufre con respecto a los distintos yacimientos y la hipótesis alternativa es que el porcentaje de azufre es diferente al menos en dos yacimientos.

Tenemos en cuenta que para que el ejercicio esté realizado de forma correcta los datos tienen que estar introducidos tal y como vienen en Figura 27, es decir, las observaciones en una sola columna y a continuación especificado su tratamiento y su bloque correspondiente.

Para cargar los datos utilizamos la función `read.table` indicando el nombre del archivo (que debe de estar en el directorio de trabajo) e indicando además que tiene cabecera.

Nota: La ruta hasta llegar al fichero varía en función del ordenador. Utilizar la orden `setwd()` para situarse en el directorio de trabajo

```
porcentaje <- read.table("guiado1-1.txt", header = TRUE)
porcentaje
```

##	Azufre	Yacimiento
## 1	151	1
## 2	192	1
## 3	108	1
## 4	204	1
## 5	214	1
## 6	176	1
## 7	117	1
## 8	169	2
## 9	64	2
## 10	90	2
## 11	141	2
## 12	101	2
## 13	128	2
## 14	159	2
## 15	156	2
## 16	122	3
## 17	132	3

```
## 18 139 3
## 19 133 3
## 20 154 3
## 21 104 3
## 22 225 3
## 23 149 3
## 24 130 3
## 25 75 4
## 26 126 4
## 27 69 4
## 28 62 4
## 29 90 4
## 30 120 4
## 31 32 4
## 32 73 4
## 33 80 5
## 34 90 5
## 35 124 5
## 36 82 5
## 37 72 5
## 38 57 5
## 39 118 5
## 40 54 5
## 41 130 5
```

Debemos transformar la variable referente a los niveles del factor fijo como factor para poder hacer los cálculos de forma adecuada:

```
porcentaje$Yacimiento<-factor(porcentaje$Yacimiento)
porcentaje$Yacimiento

## [1] 1 1 1 1 1 1 1 2 2 2 2 2 2 2 3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 4 4 5 5 5 5 5 5
## [39] 5 5 5
## Levels: 1 2 3 4 5
```

Para calcular la tabla ANOVA primero hacemos uso de la función "aov" de la siguiente forma:

```
mod <- aov(Azufre ~ Yacimiento, data = porcentaje)
```

donde:

- Azufre: Nombre de la columna de las observaciones.
- Yacimiento: Nombre de la columna en la que están representados los tratamientos.

- data= data.frame en el que están guardados los datos.

```
mod

## Call:
## aov(formula = Azufre ~ Yacimiento, data = porcentaje)
##
## Terms:
##              Yacimiento Residuals
## Sum of Squares    40432.68  42639.76
## Deg. of Freedom         4        36
##
## Residual standard error: 34.41566
## Estimated effects may be unbalanced
```

Se puede mostrar un resumen de los resultados con la función "summary" (verdadera tabla ANOVA)

```
summary(mod)

##              Df Sum Sq Mean Sq F value    Pr(>F)
## Yacimiento     4  40433   10108    8.534 5.97e-05 ***
## Residuals    36  42640    1184
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

En la Tabla ANOVA, el valor del estadístico de contraste de igualdad de medias, $F = 8.534$ deja a su derecha un p-valor menor que 0.001, menor que el nivel de significación del 5%, por lo que se rechaza la Hipótesis nula de igualdad de medias. Es decir, existen diferencias significativas en el contenido medio de azufre entre los cinco yacimientos. La pregunta que nos planteamos es si el contenido de azufre es significativamente distinto en los cinco yacimientos o sólo en alguno de ellos. Para responder a esta pregunta utilizamos algún procedimiento de comparaciones múltiples. En el apartado siguiente responderemos a esta cuestión.

2. Si se rechaza la hipótesis nula que las medias de porcentaje de azufre en los cinco yacimientos es la misma, determinar que medias difieren entre sí utilizando el método de comparaciones múltiples de Tukey.

```
mod.tukey <- TukeyHSD(mod, ordered = TRUE)
mod.tukey

## Tukey multiple comparisons of means
## 95% family-wise confidence level
## factor levels have been ordered
```



```
##
## Fit: aov(formula = Azufre ~ Yacimiento, data = porcentaje)
##
## $Yacimiento
##          diff          lwr          upr          p adj
## 5-4  8.791667 -39.217257  56.80059 0.9841364
## 2-4 45.125000  -4.275775  94.52577 0.0873389
## 3-4 62.236111  14.227188 110.24503 0.0057086
## 1-4 85.125000  33.990340 136.25966 0.0002709
## 2-5 36.333333 -11.675590  84.34226 0.2131394
## 3-5 53.444444   6.868947 100.01994 0.0177365
## 1-5 76.333333  26.542032 126.12463 0.0008288
## 3-2 17.111111 -30.897812  65.12003 0.8429902
## 1-2 40.000000 -11.134660  91.13466 0.1865081
## 1-3 22.888889 -26.902412  72.68019 0.6809794
```

Se comprueba que no se detectan diferencias significativas entre los yacimientos 1, 2 y 3 y entre los yacimientos 2, 4 y 5. Para ello nos fijamos en las Significaciones (mayores que 0.05) o en los límites de los intervalos. Dos medias se declaran iguales si el cero pertenece al intervalo de confianza construido para la diferencia de ellas.

3. Estudiar las hipótesis de modelo: Homocedasticidad (Homogeneidad de las varianzas por grupo), Independencia y Normalidad.

Hipótesis de Homocedasticidad: Test de Barlett

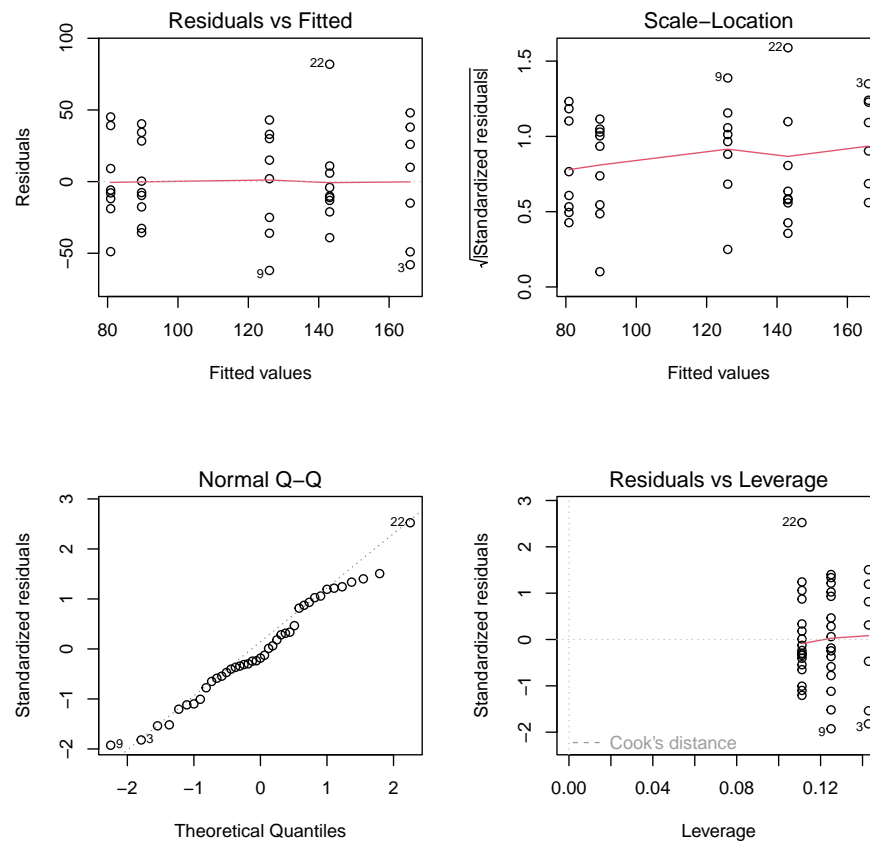
```
bartlett.test(porcentaje$Azufre, porcentaje$Yacimiento)

##
## Bartlett test of homogeneity of variances
##
## data:  porcentaje$Azufre and porcentaje$Yacimiento
## Bartlett's K-squared = 1.2655, df = 4, p-value = 0.8672
```

La salida muestra el resultado del contraste de Barlett de igualdad de varianzas en todos los grupos. El estadístico de contraste experimental, $B = 1.2655$, deja a la derecha un p-valor = 0.8672, que nos indica que no se debe rechazar la igualdad entre las varianzas.

Hipótesis de Independencia: Esta hipótesis la comprobaremos gráficamente mediante la representación de los residuos frente a los valores pronosticados por el modelo.

```
layout(matrix(c(1,2,3,4),2,2))
plot(mod)
```



En esta salida interpretamos el gráfico que se muestra en la Fila 1, Columna 1. Es decir, el gráfico el que se representan los residuos en el eje de ordenadas y los valores ajustados por el modelo en el eje de abscisas. Este gráfico no muestra ningún aspecto que haga sospechar de la hipótesis de independencia de los residuos.

La hipótesis de Normalidad la comprobaremos gráficamente y analíticamente

Gráficamente comprobaremos la normalidad mediante un histograma y el gráfico Q-Q plot

En primer lugar realizaremos el histograma, para ello
En primer lugar calculemos los residuos del modelo

```
g = mod$residuals
g
```

##	1	2	3	4	5	6
##	-15.0000000	26.0000000	-58.0000000	38.0000000	48.0000000	10.0000000
##	7	8	9	10	11	12
##	-49.0000000	43.0000000	-62.0000000	-36.0000000	15.0000000	-25.0000000
##	13	14	15	16	17	18
##	2.0000000	33.0000000	30.0000000	-21.1111111	-11.1111111	-4.1111111
##	19	20	21	22	23	24
##	-10.1111111	10.8888889	-39.1111111	81.8888889	5.8888889	-13.1111111
##	25	26	27	28	29	30
##	-5.8750000	45.1250000	-11.8750000	-18.8750000	9.1250000	39.1250000
##	31	32	33	34	35	36
##	-48.8750000	-7.8750000	-9.6666667	0.3333333	34.3333333	-7.6666667
##	37	38	39	40	41	
##	-17.6666667	-32.6666667	28.3333333	-35.6666667	40.3333333	

Calculamos la media de los residuos

```
m <- mean(g)
m
```

```
## [1] -3.462677e-16
```

Calculamos la desviación típica

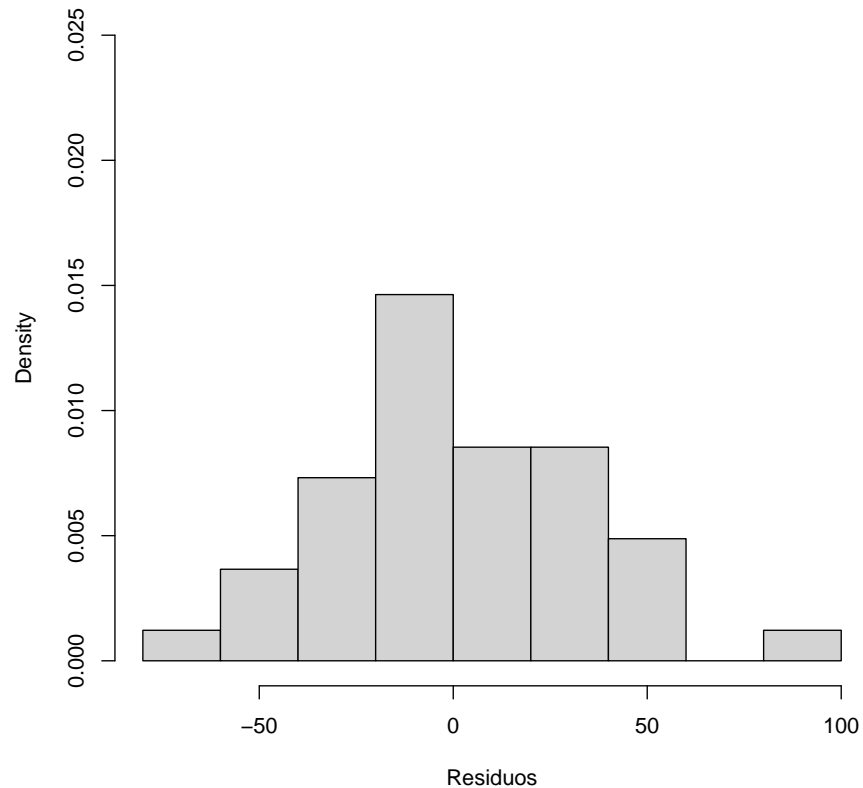
```
std <- sqrt(var(g))
std
```

```
## [1] 32.64957
```

Representamos el histograma

```
hist(g, prob = TRUE, xlab = "Residuos", ylim = c(0, 0.025),
     main = "F. dens. Normal hist. de residuos")
```

F. dens. Normal hist. de residuos



Y la curva Normal sobre el Histograma

```
curve(dnorm(x, mean = m, sd = std), col = "red",  
      lwd = 2, add = TRUE, yaxt = "n")  
  
## Error in plot.xy(xy.coords(x, y), type = type, ...): plot.new has  
## not been called yet
```

Anteriormente hemos realizado el gráfico Q-Q. Ambos gráficos no muestran desviación importante de la normalidad. Analíticamente lo vamos a comprobar mediante el contraste de Shapiro-Wilk

```
shapiro.test(mod$residuals)  
  
##  
## Shapiro-Wilk normality test
```

```
##
## data: mod$residuals
## W = 0.97902, p-value = 0.6384
```

El valor del p-valor (Sig. asintót. (bilateral)) es de 0.6384, por lo tanto no podemos rechazar la hipótesis de normalidad.

Ejercicio Guiado 2

Se realiza un estudio sobre el efecto del fotoperiodo y del genotipo en el periodo latente de infección del moho de cebada aislado AB3. Se obtienen cincuenta hojas de cuatro genotipos distintos. Cada grupo es infectado y posteriormente expuesto a diferente fotoperiodo. Los distintos fotoperiodos se trataron como bloques y se obtuvieron los siguientes datos de los totales para los bloques y tratamientos. La respuesta anotada es el número de días hasta la aparición de síntomas visibles.

	Fotoperiodo (horas de oscuridad por ciclo de 24 horas)				
Genotipo	0	2	4	8	16
Armelle	630	610	560	570	590
Golden	640	630	600	620	620
Promise	640	630	650	620	580
Emir	660	660	620	610	630

1. ¿Se puede afirmar que los diferentes genotipos no influyen en el número de días hasta la aparición de la infección? ¿Se puede concluir que los distintos fotoperiodos no afectan al tiempo de aparición de los síntomas de infección del moho?
2. En caso de que influyan significativamente alguno de los dos factores, extraer conclusiones utilizando el método de Duncan.
3. Estudiar las hipótesis de modelo: Homocedasticidad, Independencia y Normalidad.

Solución:

1. ¿Se puede afirmar que los diferentes genotipos no influyen en el número de días hasta la aparición de la infección? ¿Se puede concluir que los distintos fotoperiodos no afectan al tiempo de aparición de los síntomas de infección del moho?

En este caso se trata de un diseño en bloques completos aleatorizados. El objetivo del estudio es comparar los cuatro tipos de genotipos, por lo que se trata de un factor con cuatro niveles. Sin embargo, al realizar la medición con

los distintos fotoperiodos a los que son expuestos el moho de cebada, es posible que estos influyan sobre el periodo latente de infección del moho de cebada aislado AB3. Por ello, y al no ser directamente motivo de estudio, los fotoperiodos es un factor secundario que recibe el nombre de bloque.

Este modelo tiene que verificar los siguientes supuestos:

1. Las 20 observaciones constituyen muestras aleatorias independientes, cada una de tamaño 3, de 20 poblaciones con medias μ_{ij} , $i = 1, 2, 3, 4$ y $j = 1, 2, 3, 4, 5$.
2. Cada una de las 20 poblaciones es normal.
3. Cada una de las 20 poblaciones tiene la misma varianza.
4. Los efectos de los bloques y tratamientos son aditivos; es decir, no existe interacción entre los bloques y tratamientos. Esto significa que si hay diferencias entre dos tratamientos cualesquiera, estas se mantienen en todos los bloques (abetos).

Los tres primeros supuestos coinciden con los supuestos del modelo unifactorial, con la diferencia de que en el modelo unifactorial se examinaban I poblaciones y en este modelo se examinan IJ. El cuarto supuesto es característico del diseño en bloques. La no interacción entre los bloques y los tratamientos significa que los tratamientos tienen un comportamiento consistente a través de los bloques y que los bloques tienen un comportamiento consistente a través de los tratamientos. Expresado matemáticamente significa que la diferencia de los valores medios para dos tratamientos cualesquiera es la misma en todo un bloque y que la diferencia de los valores medios para dos bloques cualesquiera es la misma para cada tratamiento.

- Variable respuesta: Número de días hasta la aparición de síntomas visibles.
- Factor: Genotipo que tiene cuatro niveles. Es un factor de efectos fijos ya que viene decidido qué niveles concretos se van a utilizar.
- Bloque: Fotoperiodo que tiene cinco niveles. Es un factor de efectos fijos ya que viene decidido qué niveles concretos se van a utilizar.
- Modelo completo: Los cuatro tratamientos se prueban en cada bloque exactamente una vez.

Tamaño del experimento: Número total de observaciones (20). Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos introducir los datos directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R.

En este caso lo hacemos en un archivo de texto:

Para cargar los datos utilizamos la función `read.table` indicando el nombre del archivo (que debe de estar en el directorio de trabajo) e indicando además que tiene cabecera.

```
dias<-read.table("guiado2txt.txt", header = TRUE)
dias

##      Dias Fotoperiodo Genotipo
## 1    630           1         1
## 2    640           1         2
## 3    640           1         3
## 4    660           1         4
## 5    610           2         1
## 6    630           2         2
## 7    630           2         3
## 8    660           2         4
## 9    560           3         1
## 10   600           3         2
## 11   650           3         3
## 12   620           3         4
## 13   570           4         1
## 14   620           4         2
## 15   620           4         3
## 16   610           4         4
## 17   590           5         1
## 18   620           5         2
## 19   580           5         3
## 20   630           5         4
```

A continuación debemos transformar tanto la columna de los tratamientos como la de los bloques en un factor para realizar los cálculos posteriores adecuadamente.

```
dias$Fotoperiodo = factor(dias$Fotoperiodo)
dias$Fotoperiodo

## [1] 1 1 1 1 2 2 2 2 3 3 3 3 4 4 4 4 5 5 5 5
## Levels: 1 2 3 4 5
```

```
dias$Genotipo = factor(dias$Genotipo)
dias$Genotipo

## [1] 1 2 3 4 1 2 3 4 1 2 3 4 1 2 3 4 1 2 3 4
## Levels: 1 2 3 4
```

Para calcular la tabla ANOVA primero hacemos uso de la función "aov" de la siguiente forma

```
mod = aov(Dias ~ Fotoperiodo + Genotipo, data = dias)
```

donde:

- Días: Nombre de la columna de las observaciones.
- Fotoperiodo: Nombre de la columna en la que están representados los tratamientos.
- Genotipo: Nombre de la columna en la que están representados los bloques.
- data = data.frame en el que están guardados los datos

Ejecutamos

```
mod

## Call:
##   aov(formula = Dias ~ Fotoperiodo + Genotipo, data = dias)
##
## Terms:
##               Fotoperiodo Genotipo Residuals
## Sum of Squares          5030      5255      4170
## Deg. of Freedom           4         3        12
##
## Residual standard error: 18.64135
## Estimated effects may be unbalanced
```

y a continuación mostramos un resumen de los resultados con la función "summary" (verdadera tabla ANOVA):

```
summary(mod)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Fotoperiodo  4   5030  1257.5    3.619 0.0371 *
## Genotipo     3   5255  1751.7    5.041 0.0173 *
## Residuals   12   4170   347.5
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

En la Tabla ANOVA, el valor del estadístico de contraste de igualdad de medias de tratamientos, $F = 5.041$ deja a su derecha un p-valor igual a 0.017, menor que el nivel de significación del 5%, por lo que se rechaza la Hipótesis nula de igualdad de medias de tratamientos. Es decir, existen diferencias significativas

en el número de días hasta la aparición de la infección entre los cuatro genotipos.

En esta Tabla ANOVA, también se observa que el valor del estadístico de contraste de igualdad de medias de bloques, $F = 3.619$ deja a su derecha un p-valor igual a 0.037, menor que el nivel de significación del 5%, por lo que se rechaza la Hipótesis nula de igualdad de medias de bloques. Es decir, existen diferencias significativas en el número de días hasta la aparición de la infección entre los cinco tipos de fotoperiodos. Por lo tanto, se concluye que los niveles de ambos factores influyen de forma significativa en el número de días hasta la aparición de los síntomas de infección del moho.

2. En caso de que influyan significativamente alguno de los dos factores, extraer conclusiones utilizando el método de Duncan.

Primero vamos a hacer el contraste de Duncan para los tratamientos: Genotipos

```
duncan <- duncan.test(mod, "Genotipo" ,
                      main= "Número de días con diferentes genotipos")

## Error in duncan.test(mod, "Genotipo", main = "Número de días con
diferentes genotipos"): no se pudo encontrar la función "duncan.test"

duncan

## Error in eval(expr, envir, enclos): objeto 'duncan' no encontrado
```

En la tabla del factor Tipo de genotipo hay dos subconjuntos que se diferencian entre sí; el subconjunto 1 está formado por las medias del genotipo Armelle y el subconjunto 2 por las medias de los genotipos Golden, Promise y Emir. Y dentro de cada subconjunto no se aprecian diferencias significativas entre las medias. También se observa que en el genotipo Emir se produce el mayor número medio de días hasta la aparición de la infección (636) y en el genotipo Armelle se produce el menor (592).

Segundo vamos a hacer el contraste de Duncan para los bloques: Fotoperiodos

```
duncan1 <-duncan.test(mod, "Fotoperiodo",
                     main= "Número de días con diferentes fotoperiodos")

## Error in duncan.test(mod, "Fotoperiodo", main = "Número de días
con diferentes fotoperiodos"): no se pudo encontrar la función "duncan.test"

duncan1

## Error in eval(expr, envir, enclos): objeto 'duncan1' no encontrado
```

En la tabla del factor Tipo de fotoperiodo hay dos subconjuntos que se diferencian entre sí; el subconjunto 1 está formado por las medias de los Fotoperiodos 0 y 2 y el subconjunto 2 por las medias de los fotoperiodos 2, 4, 8 y 16. Y dentro de cada subconjunto no se aprecian diferencias significativas entre las medias. También se observa que en el fotoperiodo 0 se produce el mayor número medio de días hasta la aparición de la infección (642.5) y en los fotoperiodos 8 y 16 se produce el menor

3. Estudiar las hipótesis de modelo: Homocedasticidad, Independencia y Normalidad.

Estudiamos la homocedasticidad mediante el test de Barlett

```
bartlett.test(dias$Dias, dias$Fotoperiodo)

##
## Bartlett test of homogeneity of variances
##
## data: dias$Dias and dias$Fotoperiodo
## Bartlett's K-squared = 3.0629, df = 4, p-value = 0.5474
```

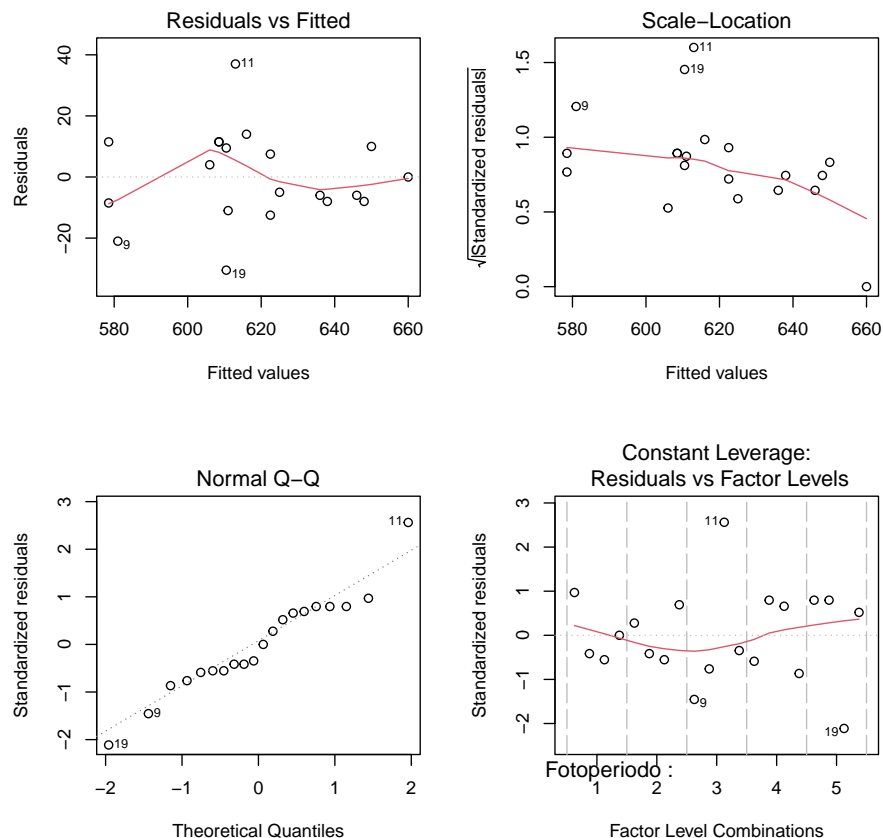
```
bartlett.test(dias$Dias, dias$Genotipo)

##
## Bartlett test of homogeneity of variances
##
## data: dias$Dias and dias$Genotipo
## Bartlett's K-squared = 1.6252, df = 3, p-value = 0.6537
```

Las Tablas muestran los resultados del contraste de Barlett de igualdad de varianzas en todos los grupos del factor genotipo y en todos los grupos del factor Fotoperiodo. Los P-valores, 0.5474 y 0.6537 indican que indican que no se debe rechazar la igualdad entre las varianzas ni el factor genotipo ni el factor fotoperiodo.

Estudiamos la independencia gráficamente

```
layout(matrix(c(1,2,3,4),2,2))
plot(mod)
```



En esta salida, como en la figura 28, interpretamos el gráfico que se muestra en la Fila 1, Columna 1. Es decir, el gráfico el que se representan los residuos en el eje de ordenadas y los valores predichos por el modelo en el eje de abscisas. Este gráfico no muestra ningún aspecto que haga sospechar de la hipótesis de independencia de los residuos.

Estudiamos la Normalidad gráficamente mediante el gráfico probabilístico normal y analíticamente mediante el contraste de Shapiro-Wilk

Observamos el gráfico que se muestra en la Fila 2, Columna 1. Es decir, el gráfico el que se representan los residuos estandarizados en el eje de ordenadas y cuantiles teóricos en el eje de abscisas. En dicho gráfico se aprecian desviaciones a la normalidad, pero el contraste ANOVA es robusto frente a desviaciones pequeñas de la normalidad. Realizaremos a continuación el contraste de Shapiro-Wilk para comprobar analíticamente la normalidad de los residuos.

```
shapiro.test(mod$residuals)

##
##  Shapiro-Wilk normality test
##
## data:  mod$residuals
## W = 0.95316, p-value = 0.4176
```

El valor del p-valor es de 0.4176, no pudiéndose rechazar la hipótesis de normalidad.

Ejercicio Guiado 3

Se realiza un estudio para determinar el efecto del nivel del agua y del tipo de planta sobre la longitud global del tallo de las plantas de guisantes. Para ello, se utilizan tres niveles de agua (bajo, medio y alto) y dos tipos de plantas (sin hojas y convencional). Se dispone para el estudio de dieciocho plantas sin hojas y dieciocho plantas convencionales. Se dividen aleatoriamente los dos tipos de plantas en tres subgrupos y después se asignan los niveles de agua aleatoriamente a los dos grupos de plantas. Los datos sobre la longitud del tallo de los guisantes (en centímetros) se muestran en la siguiente tabla:

Tipo de planta	Nivel del agua		
	Bajo	Medio	Alto
Sin hojas	69.50	96.10	121.00
	69.00	102.30	122.90
	75.00	107.50	123.10
	70.00	103.60	125.70
	74.40	100.70	125.20
	75.00	101.80	120.10
convencional	71.10	81.00	101.10
	69.20	85.80	103.20
	70.40	86.00	106.10
	73.20	87.50	109.70
	71.20	88.10	110.00
	70.90	87.60	99.00

Para un nivel de significación del 5%.

1. ¿Se puede afirmar que los distintos niveles de agua influyen en la longitud del tallo de los guisantes? ¿Y el tipo de planta?
2. ¿La efectividad del nivel del agua es la misma para los dos tipos de plantas?
3. Estudia, utilizando el método de Newman-Keuls, qué nivel de agua es más efectivo.

Solución:

1. ¿Se puede afirmar que los distintos niveles de agua influyen en la longitud del tallo de los guisantes? ¿Y el tipo de planta?
El problema planteado se modeliza a través de un diseño de dos factores con replicación. El modelo matemático es:

$$Y_{ij} = \mu + \tau + \beta_j + (\tau\beta)_{ij} + u_{ij}, i = 1, 2, 3; j = 1, 2$$

- Variable respuesta: Longitud del tallo.
- Factor A: Nivel del agua con tres niveles. Es un factor de Efectos fijos.
- Factor B: Tipo de planta con dos niveles. Es un factor de Efectos fijos.
- Tamaño del experimento: Número total de observaciones (36).

Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos introducir los datos directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R.

En este caso lo hacemos en un archivo de texto:

```
factorial= read.table("Guiado3.txt", header = TRUE)
factorial
```

##	Longitud_tallo	Nivel_agua	Tipo_planta
## 1	69.5	1	1
## 2	69.0	1	1
## 3	75.0	1	1
## 4	70.0	1	1
## 5	74.4	1	1
## 6	75.0	1	1
## 7	96.1	2	1
## 8	102.3	2	1
## 9	107.5	2	1
## 10	103.6	2	1
## 11	100.7	2	1
## 12	101.8	2	1
## 13	121.0	3	1
## 14	122.9	3	1
## 15	123.1	3	1
## 16	125.7	3	1
## 17	125.2	3	1
## 18	120.1	3	1
## 19	71.1	1	2
## 20	69.2	1	2
## 21	70.4	1	2
## 22	73.2	1	2
## 23	71.2	1	2
## 24	70.9	1	2
## 25	81.0	2	2
## 26	85.8	2	2
## 27	86.0	2	2
## 28	87.5	2	2
## 29	88.1	2	2

## 30	87.6	2	2
## 31	101.1	3	2
## 32	103.2	3	2
## 33	106.1	3	2
## 34	109.7	3	2
## 35	110.0	3	2
## 36	99.0	3	2

A continuación debemos transformar todas las columnas que contienen a los factores en un factor para poder realizar los cálculos posteriores adecuadamente.

```
factorial$agua <- factor(factorial$Nivel_agua)
factorial$agua

## [1] 1 1 1 1 1 1 2 2 2 2 2 2 3 3 3 3 3 3 1 1 1 1 1 1 2 2 2 2 2 2 3 3 3 3 3 3
## Levels: 1 2 3
```

```
factorial$planta <- factor(factorial$Tipo_planta)
factorial$planta

## [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
## Levels: 1 2
```

Para responder a este apartado debe resolverse el contraste de igualdad de medias para el factor A:

$$H_0 = \mu_1 = \mu_2 = \mu_3 = \mu \text{ vs } H_1 = \mu_i \neq \mu_j \text{ para algún } i \neq j$$

y para el factor B:

$$H_0 = \mu_1 = \mu_2 = \mu \text{ vs } H_1 = \mu_i \neq \mu_j \text{ para algún } i \neq j$$

Para calcular la tabla ANOVA en R primero hacemos uso de la función "aov" y a continuación "summary" de la siguiente forma:

```
mod = aov(Longitud_tallo ~ agua* planta ,
           data = factorial)
mod

## Call:
## aov(formula = Longitud_tallo ~ agua * planta, data = factorial)
##
## Terms:
##               agua      planta agua:planta Residuals
## Sum of Squares 10773.635 1246.090      514.145    282.250
```

```
## Deg. of Freedom      2      1      2      30
##
## Residual standard error: 3.067301
## Estimated effects may be unbalanced
```

```
summary(mod)

##              Df Sum Sq Mean Sq F value    Pr(>F)
## agua          2  10774     5387   572.56 < 2e-16 ***
## planta        1   1246     1246   132.44 1.58e-12 ***
## agua:planta    2    514      257   27.32 1.75e-07 ***
## Residuals     30    282         9
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

- El valor del estadístico de contraste de igualdad de medias del factor Nivel_agua, $F = 572.56$ deja a su derecha un p-valor menor que 0.001, menor que el nivel de significación del 5%, por lo que se rechaza la Hipótesis nula de igualdad de medias de los niveles del factor Nivel_agua. Es decir, existen diferencias significativas en la longitud del tallo de guisantes dependiendo del nivel del agua.
- El valor del estadístico de contraste de igualdad de medias del factor Tipo_planta, $F = 132.445$ deja a su derecha un p-valor menor que 0.001, menor que el nivel de significación del 5%, por lo que se rechaza la Hipótesis nula de igualdad de medias del factor Tipo_planta. Es decir, el tipo de planta afecta significativamente a la longitud del tallo de guisantes.

2.¿La efectividad del nivel del agua es la misma para los dos tipos de plantas?

Para responder a esta pregunta, realizamos el contraste de hipótesis sobre la interacción de los dos factores.

$H_0 = (\tau\beta)_{ij} = 0$ (no existe interacción) vs $H_1 = \tau\beta_{ij} \neq 0$ (existe interacción).

En la Tabla ANOVA mostrada anteriormente, el valor del estadístico de contraste de la interacción de los dos factores, $F = 27.32$ deja a su derecha un p-valor menor que 0.001, menor que el nivel de significación del 5%, por lo que se rechaza la Hipótesis nula de no interacción entre los factores. Por lo tanto la efectividad del nivel de agua no es la misma para los dos tipos de plantas. Es decir, puede ocurrir que un nivel de agua influya en el crecimiento de la longitud del tallo con un tipo de planta pero no con el otro o influya de distinta forma.

3. Estudia, utilizando el método de Newman-Keuls, qué nivel de agua es más efectivo.

```
library(agricolae)
```

```
contraste <- SNK.test(mod,"agua", console=TRUE,
                      main=" Contraste de Newman-Keuls para el factor nivel del agua")

##
## Study:  Contraste de Newman-Keuls para el factor nivel del agua
##
## Student Newman Keuls Test
## for Longitud_tallo
##
## Mean Square Error:  9.408333
##
## agua,  means
##
##   Longitud_tallo      std  r Min   Max
## 1          71.575  2.243019 12  69  75.0
## 2          94.000  8.901992 12  81 107.5
## 3         113.925 10.069949 12  99 125.7
##
## Alpha: 0.05 ; DF Error: 30
##
## Critical Range
##           2           3
## 2.557375 3.087063
##
## Means with the same letter are not significantly different.
##
##   Longitud_tallo groups
## 3          113.925    a
## 2           94.000    b
## 1           71.575    c
```

En la tabla se muestran los subgrupos formados de medias iguales al utilizar el método de Newman-Keuls. Hay tres subconjuntos que se diferencian entre sí y cada subconjunto está formado por un solo nivel de agua. También se observa que con el nivel de agua alto se produce la mayor longitud del tallo de guisantes, 113.925 cm, y con el nivel Bajo se produce el menor 71.575 cm.

8 Ejercicios propuestos

Ejercicio Propuesto 1

La convección es una forma de transferencia de calor por los fluidos debido a sus variaciones de densidad por la temperatura; las partes calientes ascienden y las frías descienden formando las corrientes de convección que hacen uniforme la temperatura del fluido. Se ha realizado un experimento para determinar las modificaciones de la densidad de fluido al elevar la temperatura en una determinada zona. Los resultados obtenidos han sido los siguientes:

Temperatura	Densidad				
100	21.8	21.9	21.7	21.6	21.7
125	21.7	21.4	21.5	21.4	
150	21.9	21.8	21.8	21.6	21.5
175	21.9	22.1	21.85	21.9	

Responder a las siguientes cuestiones:

¿Afecta la temperatura a la densidad del fluido? Determinar qué temperaturas producen modificaciones significativas en la densidad media del fluido. Estudiar las hipótesis del modelo: Homocedasticidad, independencia y normalidad. ¿Se puede afirmar que las temperaturas de 100 y 125 producen menos densidades de fluido en promedio que las temperaturas de 150 y 175?

Solucion

```
practico1<- read.table("practico1.txt", header = TRUE)
practico1

##      Densidad Temperatura
## 1      21.80          100
## 2      21.90          100
## 3      21.70          100
## 4      21.60          100
## 5      21.70          100
## 6      21.70          125
## 7      21.40          125
## 8      21.50          125
## 9      21.40          125
## 10     21.90          150
## 11     21.80          150
## 12     21.80          150
## 13     21.60          150
```

```
## 14    21.50      150
## 15    21.90      175
## 16    22.10      175
## 17    21.85      175
## 18    21.90      175
```

El problema planteado se modeliza a través de un diseño unifactorial totalmente aleatorizado de efectos fijos no-equilibrado.

- Variable respuesta: Densidad del fluido.
- Factor: Temperatura: Es un factor de Efectos fijos.
- Modelo no-equilibrado: Los niveles de los factores tienen distinto número de elementos.

¿Afecta la temperatura a la densidad del fluido?
convirtiendo la variable temperatura a factor

```
practico1$Temperatura <-factor(practico1$Temperatura)
practico1$Temperatura

## [1] 100 100 100 100 100 125 125 125 125 150 150 150 150 150 175 175 175 175
## Levels: 100 125 150 175
```

Calculando el modelo y la tabla ANOVA.

```
modp1 <- aov(Densidad ~ Temperatura, data = practico1)
modp1

## Call:
## aov(formula = Densidad ~ Temperatura, data = practico1)
##
## Terms:
##              Temperatura Residuals
## Sum of Squares    0.384375  0.256875
## Deg. of Freedom           3        14
##
## Residual standard error: 0.1354556
## Estimated effects may be unbalanced

summary(modp1)

##              Df Sum Sq Mean Sq F value  Pr(>F)
## Temperatura  3  0.3844  0.12813    6.983 0.00419 **
## Residuals    14  0.2569  0.01835
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El valor de F experimental 6.983, deja a su derecha un p-valor = 0.00419 inferior a 0.05, por lo que se rechaza la hipótesis nula de igualdad de medias. Concluyendo que existen diferencias significativas en la densidad del fluido en función de la modificación de la temperatura

2. Determinar qué temperaturas producen modificaciones significativas en la densidad media del fluido

```
modp1.tukey <- TukeyHSD(modp1, ordered = TRUE)
modp1.tukey

##    Tukey multiple comparisons of means
##      95% family-wise confidence level
##      factor levels have been ordered
##
## Fit: aov(formula = Densidad ~ Temperatura, data = practico1)
##
## $Temperatura
##           diff           lwr           upr           p adj
## 150-125 0.2200 -0.04410917 0.4841092 0.1184044
## 100-125 0.2400 -0.02410917 0.5041092 0.0807121
## 175-125 0.4375 0.15910449 0.7158955 0.0021907
## 100-150 0.0200 -0.22900452 0.2690045 0.9953076
## 175-150 0.2175 -0.04660917 0.4816092 0.1240769
## 175-100 0.1975 -0.06660917 0.4616092 0.1785203
```

Se producen diferencias significativas entre las temperaturas 175 y 125

3. Estudiar las hipótesis del modelo: Homocedasticidad, independencia y normalidad.

HOMOCEDASTICIDAD

```
bartlett.test(practico1$Densidad, practico1$Temperatura)

##
## Bartlett test of homogeneity of variances
##
## data:  practico1$Densidad and practico1$Temperatura
## Bartlett's K-squared = 0.69212, df = 3, p-value = 0.8751
```

La salida muestra el resultado del contraste de Barlett de igualdad de varianzas en todos los grupos. El estadístico de contraste experimental, B= 0.69212, deja a la derecha un p-valor = 0.8751, que nos indica que no se debe rechazar la igualdad entre las varianzas.

NORMALIDAD

```
shapiro.test(modp1$residuals)

##
##  Shapiro-Wilk normality test
##
## data:  modp1$residuals
## W = 0.94639, p-value = 0.3712
```

El valor p para la normalidad de los residuos es mayor que 0.05 por lo tanto esto implica normalidad en los los datos.

Ejercicio propuesto 2

Un laboratorio de reciclaje controla la calidad de los plásticos utilizados en bolsas. Se desea contrastar si existe variabilidad en la calidad de los plásticos que hay en el mercado. Para ello, se eligen al azar cuatro plásticos y se les somete a una prueba para medir el grado de resistencia a la degradación ambiental. De cada plástico elegido se han seleccionado ocho muestras y los resultados de la variable que mide la resistencia son los de la tabla adjunta.

Calidad plásticos	Resistencia						
Plástico A	135	175	97	169	213	171	115
Plástico B	275	170	154	133	219	187	220
Plástico C	169	239	184	222	253	179	280
Plástico D	115	105	93	85	120	74	87

¿Qué conclusiones se deducen de este experimento?

Solución:

Los cuatro tipos de plásticos analizados corresponden a una selección aleatoria de 4 conjuntos de observaciones extraídos aleatoriamente del total de diferentes tipos de plásticos que hay en el mercado, entre los cuales debemos observar si existen o no diferencias significativas. Nos encontramos por tanto ante un diseño unifactorial completamente aleatorio con efectos aleatorios .

En este modelo, se supone que las variables τ_i son variables aleatorias normales independientes con media 0 y varianza común σ_τ^2 .

```
Propuesto2 <- read.table("propuesto2.txt", header = TRUE)
Propuesto2

##      Resistencia Plastico
```

```
## 1      135      P1
## 2      175      P1
## 3       97      P1
## 4      169      P1
## 5      213      P1
## 6      171      P1
## 7      115      P1
## 8      143      P1
## 9      275      P2
## 10     170      P2
## 11     154      P2
## 12     133      P2
## 13     219      P2
## 14     187      P2
## 15     220      P2
## 16     185      P2
## 17     169      P3
## 18     239      P3
## 19     184      P3
## 20     222      P3
## 21     253      P3
## 22     179      P3
## 23     280      P3
## 24     193      P3
## 25     115      P4
## 26     120      P4
## 27     105      P4
## 28       74      P4
## 29       93      P4
## 30       87      P4
## 31       85      P4
## 32       63      P4
```

```
Propuesto2$Plastico <- factor(Propuesto2$Plastico)
Propuesto2$Plastico
```

```
## [1] P1 P1 P1 P1 P1 P1 P1 P1 P1 P2 P2 P2 P2 P2 P2 P2 P2 P3 P3 P3 P3 P3 P3 P3 P3 P4
## [26] P4 P4 P4 P4 P4 P4 P4
## Levels: P1 P2 P3 P4
```

```
modp2 <- aov(Resistencia ~ Plastico, data = Propuesto2)
modp2

## Call:
```

```
##      aov(formula = Resistencia ~ Plastico, data = Propuesto2)
##
## Terms:
##              Plastico Residuals
## Sum of Squares  69072.12  37410.75
## Deg. of Freedom      3      28
##
## Residual standard error: 36.55268
## Estimated effects may be unbalanced

summary(modp2)

##              Df Sum Sq Mean Sq F value    Pr(>F)
## Plastico      3  69072    23024   17.23 1.55e-06 ***
## Residuals    28  37411     1336
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El valor del estadístico de contraste es igual a 17.23 que deja a la derecha un p-valor de 1.55e-06, rechazando la Hipótesis nula tanto a un nivel de significación del 5% como del 1%. Podemos concluir que los datos muestran evidencias de variabilidad en la resistencia para la degradación ambiental según el tipo de plástico empleado en la fabricación de la bolsa.

Ejercicio Propuesto 3

Debido a la proliferación de los campos de golf y a la gran cantidad de agua que necesitan, un grupo de científicos estudia la calidad de varios tipos de césped para implantarlo en invierno en los campos de golf. Para ello, miden la distancia recorrida por una pelota de golf, en el campo, después de bajar por una rampa (para proporcionar a la pelota una velocidad inicial constante). El terreno del que disponen tiene mayor pendiente en la dirección norte-sur, por lo que se aconseja dividir el terreno en cinco bloques de manera que las pendientes de las parcelas individuales dentro de cada bloque sean las mismas. Se utilizó el mismo método para la siembra y las mismas cantidades de semilla. Las mediciones son las distancias desde la base de la rampa al punto donde se pararon las pelotas. En el estudio se incluyeron las variedades: *Agrostis Tenuis* (Césped muy fino y denso, de hojas cortas y larga duración), *Agrostis Canina* (Hoja muy fina, estolonífera. Forma una cubierta muy tupida), *Paspalum Notatum* (Hojas gruesas, bastas y con rizomas. Forma una cubierta poco densa), *Paspalum Vaginatum* (Césped fino, perenne, con rizomas y estolones).

	Bloque 1	Bloque 2	Bloque 3	Bloque 4	Bloque 5
Agrosty Tennis	1.30	1.60	0.50	1.20	1.10
Agrosty Canina	2.20	2.40	0.40	2.00	1.80
Paspalum Notatum	1.80	1.70	0.60	1.50	1.30
Paspalum Vaginatum	3.90	4.40	2.00	4.10	3.40

Se pide:

1. Identificar los elementos del estudio (factores, unidades experimentales, variable respuesta, etc.) y plantear detalladamente el modelo matemático utilizado en el experimento.
2. ¿Son los bloques fuente de variación?
3. Existen diferencias reales entre las distancias medias recorridas por una pelota de golf en los distintos tipos de césped?
4. Estudiar las interacciones de los factores.
5. Comprobar que se cumplen las hipótesis del modelo.
6. Utilizando el método de Duncan y Newman-Keuls, ¿qué tipo de césped ofrece menor resistencia al recorrido de las pelotas?

solucion

1. Identificar los elementos del estudio (factores, unidades experimentales, variable respuesta, etc.) y plantear detalladamente el modelo matemático utilizado en el experimento.

- Variable respuesta: Distancia.
- Factor: Tipo_Césped que tiene cuatro niveles. Es un factor de efectos fijos ya que viene decidido qué niveles concretos se van a utilizar.
- Bloque: Bloques que tiene cinco niveles. Es un factor de efectos fijos ya que viene decidido qué niveles concretos se van a utilizar.
- Modelo completo: Los cuatro tratamientos se prueban en cada bloque exactamente una vez.
- Tamaño del experimento: Número total de observaciones (20).

Este experimento se modeliza mediante un diseño en Bloques completos al azar. El modelo matemático es:

$$y_{ij} = \mu + \tau_i + \beta_j + u_{ij}, i = 1, 2, \dots, j = 1, \dots, 5$$

2. ¿Son los bloques fuente de variación?

Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos introducir los datos directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R.

En este caso lo hacemos en un archivo de texto:

Para cargar los datos utilizamos la función `read.table` indicando el nombre del archivo (que debe de estar en el directorio de trabajo) e indicando además que tiene cabecera.

Nota: La ruta hasta llegar al fichero varía en función del ordenador. Utilizar la orden `setwd()` para situarse en el directorio de trabajo

```
propuesto3<-read.table("propuesto3.txt", header = TRUE)
propuesto3$cesped=factor(propuesto3$cesped)
propuesto3$cesped

## [1] C1 C1 C1 C1 C1 C2 C2 C2 C2 C2 C3 C3 C3 C3 C3 C4 C4 C4 C4 C4
## Levels: C1 C2 C3 C4

propuesto3$bloque=factor(propuesto3$bloque)
propuesto3$bloque

## [1] B1 B2 B3 B4 B5 B1 B2 B3 B4 B5 B1 B2 B3 B4 B5 B1 B2 B3 B4 B5
## Levels: B1 B2 B3 B4 B5
```

Para calcular la tabla ANOVA primero hacemos uso de la función "aov" y a continuación mostramos un resumen de los resultados con la función "summary"

```
modp3<- aov(distancia ~ cesped+bloque, data = propuesto3)
modp3

## Call:
## aov(formula = distancia ~ cesped + bloque, data = propuesto3)
##
## Terms:
##              cesped bloque Residuals
## Sum of Squares 18.044  6.693      0.951
## Deg. of Freedom    3      4         12
##
## Residual standard error: 0.2815138
## Estimated effects may be unbalanced

summary(modp3)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## cespced      3 18.044    6.015    75.89 4.52e-08 ***
## bloque      4  6.693    1.673    21.11 2.32e-05 ***
## Residuals   12  0.951    0.079
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

En la Tabla ANOVA, el valor del estadístico de contraste de igualdad de medias de tratamientos, $F = 75.89$ deja a su derecha un p-valor igual a $4.52e-08$, menor que el nivel de significación del 5%, por lo que se rechaza la Hipótesis nula de igualdad de tratamientos. Así, los tipos de césped influyen en las distancias recorridas por las pelotas. Es decir, existen diferencias significativas en las distancias recorridas por las pelotas entre los cuatro tipos de césped.

En esta Tabla ANOVA, también se observa que el valor del estadístico de contraste de igualdad de medias de bloques, $F = 21.11$ deja a su derecha un p-valor igual a $2.32e-05$, menor que el nivel de significación del 5%, por lo que se rechaza la Hipótesis nula de igualdad de bloques. La eficacia de este diseño depende de los efectos de los bloques. En este caso este diseño es más eficaz que el diseño completamente aleatorizado y el contraste principal de las medias de los tratamientos será más sensible a las diferencias entre tratamientos. Por lo tanto la inclusión del factor bloque en el modelo es acertada. Así, las distancias recorridas por las pelotas dependen del tipo de terreno.

3. Existen diferencias reales entre las distancias medias recorridas por una pelota de golf en los distintos tipos de césped?

Esta cuestión está contestada afirmativamente en el apartado anterior, en el que la tabla ANOVA nos muestra un valor de $F = 75.89$ y un Sig. menor que 0.001

4. Estudiar las interacciones de los factores.

La interacción entre el factor bloque y los tratamientos vamos a estudiarla analíticamente mediante el Test de Interacción de un grado de Tukey

Para realizar este test en R tenemos que utilizar la library "daewr" y dentro de ella la función "Tukey1df". De la siguiente forma:

```
library(daewr)
Tukey1df(propuesto3)
```

## Source	df	SS	MS	F	Pr>F
## A	3	18.044	6.0147		
## B	4	6.693	1.6732		
## Error	12	0.951	1.268		

## NonAdditivity	1	0.6155	0.6155	20.18	9e-04
## Residual	11	0.3355	0.0305		

Puesto que el p-valor ($\Pr \hat{\chi}^2 F$) es 9e-04 rechazamos la hipótesis nula de no interacción.

5. Comprobar que se cumplen las hipótesis del modelo.

Hipótesis de Normalidad

La normalidad las vamos a comprobar analíticamente y gráficamente. Analíticamente mediante el contraste de Shapiro-Wilk que es adecuado cuando las muestras son pequeñas ($n \leq 50$)

```
shapiro.test(mod$residuals)

##
##  Shapiro-Wilk normality test
##
## data:  mod$residuals
## W = 0.97552, p-value = 0.5939
```

Como podemos observar tenemos un p-valor de 0.9437 que aceptaría la hipótesis de normalidad por ser mayor al 5% (nivel de significación usual).

Hipótesis de Homogeneidad de Varianzas

Para comprobar la hipótesis de homocedasticidad utilizamos el Test de Bartlett distinguiendo entre la igualdad entre varianzas del factor principal y la igualdad de varianzas del factor bloque.

En nuestro ejemplo, el test para igualdad de varianzas del factor principal sería:

```
bartlett.test(propuesto3$distancia, propuesto3$cesped)

##
##  Bartlett test of homogeneity of variances
##
## data:  propuesto3$distancia and propuesto3$cesped
## Bartlett's K-squared = 3.3347, df = 3, p-value = 0.3428
```

El p-valor es del 0.3428 que al ser mayor del nivel significación usual del 5% no podemos rechazar la hipótesis de igualdad de varianzas en el factor principal.

De la misma manera procedemos para el factor bloque:

```
bartlett.test(propuesto3$distancia, propuesto3$bloque)

##
## Bartlett test of homogeneity of variances
##
## data: propuesto3$distancia and propuesto3$bloque
## Bartlett's K-squared = 0.94442, df = 4, p-value = 0.9181
```

El p-valor es 0.9181, mayor que 0.05 por lo que no podemos rechazar la hipótesis de igualdad de varianzas en el factor bloque.

6. Utilizando el método de Duncan y Newman-Keuls, ¿qué tipo de césped ofrece menor resistencia al recorrido de las pelotas?

Para poder hacer uso de ambos contrastes en R tenemos que instalar en primer lugar el paquete "agricolae"

```
library(agricolae)
(duncan=duncan.test(modp3, "cesped" , group = T))

## $statistics
##   MSerror Df Mean      CV
##   0.07925 12 1.96 14.36295
##
## $parameters
##   test name.t ntr alpha
##   Duncan cesp  4  0.05
##
## $duncan
##      Table CriticalRange
## 2 3.081307      0.3879266
## 3 3.225244      0.4060478
## 4 3.312453      0.4170272
##
## $means
##   distancia      std r Min Max Q25 Q50 Q75
## C1      1.14 0.4037326 5 0.5 1.6 1.1 1.2 1.3
## C2      1.76 0.7924645 5 0.4 2.4 1.8 2.0 2.2
## C3      1.38 0.4764452 5 0.6 1.8 1.3 1.5 1.7
## C4      3.56 0.9449868 5 2.0 4.4 3.4 3.9 4.1
##
## $comparison
## NULL
##
## $groups
##   distancia groups
## C4      3.56      a
```

```
## C2      1.76      b
## C3      1.38     bc
## C1      1.14      c
##
## attr(,"class")
## [1] "group"

SNK.test(mod, "cesped", console =TRUE)

## Name:   cespced
## agua planta
```

Ejercicio Propuesto 4

Consideremos de nuevo el ejercicio propuesto 3 sobre un grupo de científicos que estudia la calidad de varios tipos de césped para implantarlo en invierno en los campos de golf. Para ello, miden la distancia recorrida por una pelota de golf, en el campo, después de bajar por una rampa (para proporcionar a la pelota una velocidad inicial constante). El terreno del que disponen tiene mayor pendiente en la dirección norte-sur, por lo que se aconseja dividir el terreno en cinco bloques de manera que las pendientes de las parcelas individuales dentro de cada bloque sean las mismas. Se utilizó el mismo método para la siembra y las mismas cantidades de semilla. Las mediciones son las distancias desde la base de la rampa al punto donde se pararon las pelotas, y al realizar dichas mediciones no se han podido obtener una para cada combinación de tipo de césped y tipo de terreno, sino que sólo se han podido realizar con tres de las variedades del césped en cada uno de los bloques de terreno. Para controlar el efecto del tipo de terreno deciden utilizar un diseño en bloques incompletos. En el estudio se incluyeron las variedades: *Agrostis Tenuis* (Césped muy fino y denso, de hojas cortas y larga duración), *Agrostis Canina* (Hoja muy fina, estolonífera. Forma una cubierta muy tupida), *Paspalum Notatum* (Hojas gruesas, bastas y con rizomas. Forma una cubierta poco densa), *Paspalum Vaginatatum* (Césped fino, perenne, con rizomas y estolones).

Variedad de césped	Bloque1	Bloque2	Bloque3	Bloque4	Bloque5	Bloque6
<i>Agrosty Tennis</i>	1.30		0.50		1.80	
<i>Agrosty Canina</i>	2.20			1.50		1.80
<i>Paspalum Notatum</i>		2.40	2.00			1.60
<i>Paspalum Vaginatatum</i>		4.40		4.10	3.40	

Se pide:

1. Identificar los elementos del estudio (factores, unidades experimentales, variable respuesta, etc.) y plantear detalladamente el modelo matemático utilizado en el experimento.
2. ¿Son los bloques fuente de variación?

3. Existen diferencias reales entre las distancias medias recorridas por una pelota de golf en los distintos tipos de césped?
4. Comprobar que se cumplen las hipótesis del modelo.
5. Utilizando el método de Newman-Keuls, ¿qué tipo de césped ofrece menor resistencia al recorrido de las pelotas?

Solución

1 Identificar los elementos del estudio (factores, unidades experimentales, variable respuesta, etc.) y plantear detalladamente el modelo matemático utilizado en el experimento.

- Variable respuesta: Distancia.
- Factor: Tipo.Césped que tiene cuatro niveles. Es un factor de efectos fijos ya que viene decidido qué niveles concretos se van a utilizar.
- Bloque: Bloques que tiene seis niveles. Es un factor de efectos fijos ya que viene decidido qué niveles concretos se van a utilizar.
- Modelo incompleto: Todos los tratamientos no se prueban en cada bloque.
- Tamaño del experimento: Número total de observaciones (12).

```
propuesto4<-read.table("propuesto4.txt", header = TRUE)
propuesto4
```

##	distancia	cesped	bloque
## 1	1.3	C1	B1
## 2	2.2	C2	B1
## 3	2.4	C3	B2
## 4	4.4	C4	B2
## 5	0.5	C1	B3
## 6	2.0	C3	B3
## 7	1.5	C2	B4
## 8	4.1	C4	B4
## 9	1.8	C1	B5
## 10	3.4	C4	B5
## 11	1.8	C2	B6
## 12	1.6	C3	B6

A continuación debemos transformar tanto la columna de los tratamientos como la de los bloques en un factor para poder realizar los cálculos posteriores adecuadamente.

```
propuesto4$Tratamiento = factor(propuesto4$cesped)
propuesto4$Tratamiento

## [1] C1 C2 C3 C4 C1 C3 C2 C4 C1 C4 C2 C3
## Levels: C1 C2 C3 C4
```

```
propuesto4$Bloque = factor(propuesto4$bloque)
propuesto4$Bloque

## [1] B1 B1 B2 B2 B3 B3 B4 B4 B5 B5 B6 B6
## Levels: B1 B2 B3 B4 B5 B6
```

2. ¿Son los bloques fuente de variación?

Para poder analizar los datos mediante un diseño BIB debemos instalar y cargar dos paquetes de R especializados en este tipo de diseños:

```
library(daewr)
```

Por lo que tenemos también que cargar e instalar el paquete colorspace

```
library(colorspace)
```

Cargamos e instalamos el paquete zoo

```
library(zoo)

##
## Attaching package: 'zoo'
## The following objects are masked from 'package:base':
##
##      as.Date, as.Date.numeric

library(daewr)
library(AlgDesign)
```

La función "BIBsize(t, k)" de la librería daewr nos permite saber si el diseño puede realizarse. Calcula los parámetros del diseño donde

- t = número de niveles del factor tratamiento.
- k = número de tratamientos por bloque. Ejecutamos:

```
BIBsize(t = 4 , k = 2)

## Posible BIB design with b= 6 and r= 3 lambda= 1
```

El análisis de este modelo lo podemos realizar en R de dos formas:

A) Realizaremos el análisis: evaluando primero el efecto de los bloques y después el de los tratamientos utilizando dos funciones

- Para evaluar el efecto de los bloques, la suma de cuadrados de bloques debe ajustarse por los tratamientos, por lo tanto primero se introducen los tratamientos y después los bloques:
- Para calcular la tabla ANOVA hacemos uso de la función "aov" (aov(y ~ A + B, data=mydataframe) asume suma de cuadrados tipo I) de la siguiente forma:

```
modp4 <- aov(distancia ~ cespel + bloque,
              data = propuesto4 )
modp4

## Call:
## aov(formula = distancia ~ cespel + bloque, data = propuesto4)
##
## Terms:
##              cespel      bloque Residuals
## Sum of Squares 12.856667  0.890833  1.062500
## Deg. of Freedom      3          5          3
##
## Residual standard error: 0.595119
## Estimated effects may be unbalanced

summary(modp4)

##              Df Sum Sq Mean Sq F value Pr(>F)
## cespel        3 12.857   4.286   12.100  0.035 *
## bloque        5  0.891   0.178    0.503  0.766
## Residuals     3  1.063   0.354
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El p-valor 0.035 menor que el nivel de significación 0.05, por lo que se rechaza la hipótesis nula de igualdad de bloques y concluimos que los bloques son una fuente de variación.

3. Existen diferencias reales entre las distancias medias recorridas por una pelota de golf en los distintos tipos de césped?

Para evaluar el efecto de los tratamientos, la suma de cuadrados de tratamientos debe ajustarse por bloques, por lo tanto primero se introducen los bloques y después los tratamientos.

```
modp42 <- aov(distancia ~ bloque + cespced, data = propuesto4 )
```

donde:

- distancia = nombre de la columna de las observaciones
- cespced = nombre de la columna en la que están representados los tratamientos
- bloque = nombre de la columna en la que están representados los bloques
- data = data.frame en el que están guardados los datos

```
modp42

## Call:
## aov(formula = distancia ~ bloque + cespced, data = propuesto4)
##
## Terms:
##               bloque cespced Residuals
## Sum of Squares  6.6000  7.1475    1.0625
## Deg. of Freedom      5      3        3
##
## Residual standard error: 0.595119
## Estimated effects may be unbalanced

summary(modp42)

##              Df Sum Sq Mean Sq F value Pr(>F)
## bloque         5  6.600   1.3200   3.727 0.1540
## cespced         3  7.148   2.3825   6.727 0.0759 .
## Residuals      3  1.063   0.3542
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El valor del estadístico de contraste de igualdad de tipo de césped, $F = 3.727$, deja a su derecha un p-valor 0.1540, mayor que el nivel de significación del 5%, por lo que no se rechaza la Hipótesis Nula de igualdad de tratamientos. Por lo tanto no hay diferencias reales entre las distancias medias recorridas por una pelota de golf en los distintos tipos de césped ya que el p-valor es mayor que 0.05.

B) Realizaremos el análisis evaluando tanto los tratamientos como los bloques ejecutando solo una función.

Para ello necesitamos instalar y cargar el paquete "car": Una vez instalado cargado el paquete realizamos el ANOVA

```
modp43 <- lm(distancia ~ cespced + bloque,
              data = propuesto4 )
modp43

##
## Call:
## lm(formula = distancia ~ cespced + bloque, data = propuesto4)
##
## Coefficients:
## (Intercept)      cespcedC2      cespcedC3      cespcedC4      bloqueB2      bloqueB3
##          1.4375          0.6250          0.8250          2.5500          0.2750         -0.6000
##      bloqueB4      bloqueB5      bloqueB6
##         -0.2250         -0.1125         -0.4625

car::Anova(modp43, type="III")

## Anova Table (Type III tests)
##
## Response: distancia
##           Sum Sq Df F value  Pr(>F)
## (Intercept)  2.7552  1  7.7794 0.06847 .
## cespced      7.1475  3  6.7271 0.07589 .
## bloque      0.8908  5  0.5031 0.76563
## Residuals    1.0625  3
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Los resultados obtenidos coinciden con los realizados primero a los bloque y después a los tratamientos.

4. Comprobar que se cumplen las hipótesis del modelo.

Los supuestos que han de verificarse son Normalidad, Homocedasticidad e Independencia además del supuesto de aditividad entre los tratamientos y los bloques.

Hipótesis de normalidad

Comprobamos la hipótesis de normalidad mediante el análisis de la normalidad de los residuos. Para ello, hacemos uso del test de Shapiro-Wilks:

```
shapiro.test(modp43$residuals)

##
##  Shapiro-Wilk normality test
##
## data:  modp43$residuals
## W = 0.94151, p-value = 0.5178
```

El p-valor, 0.5178, es mayor que el nivel de significación del 5%, aceptándose la hipótesis de normalidad.

Homogeneidad de varianzas

En este caso hacemos uso del Test de Barlett para contrastar la igualdad entre varianzas del factor.

```
bartlett.test(propuesto4$distancia, propuesto4$cesped)

##
##  Bartlett test of homogeneity of variances
##
## data:  propuesto4$distancia and propuesto4$cesped
## Bartlett's K-squared = 0.76822, df = 3, p-value = 0.8571
```

El p-valor es 0.8571 por lo tanto no se puede rechazar la hipótesis de homogeneidad de las varianzas y se concluye que la cuatro variedades de césped tienen varianzas homogéneas.

```
bartlett.test(propuesto4$distancia, propuesto4$bloque)

##
##  Bartlett test of homogeneity of variances
##
## data:  propuesto4$distancia and propuesto4$bloque
## Bartlett's K-squared = 3.1837, df = 5, p-value = 0.6717
```

El p-valor es mayor que 0.05, por lo tanto no se puede rechazar la hipótesis de homogeneidad de las varianzas los bloques.

5. Utilizando el método de Newman-Keuls, ¿qué tipo de césped ofrece menor resistencia al recorrido de las pelotas?

```
library(agricolae)
Newman_Keuls<- SNK.test(modp43,"cesped", console=TRUE,
                        main="")
```

```
##
## Study:
##
## Student Newman Keuls Test
## for distancia
##
## Mean Square Error: 0.3541667
##
## cespced, means
##
##      distancia      std r Min Max
## C1  1.200000 0.6557439 3 0.5 1.8
## C2  1.833333 0.3511885 3 1.5 2.2
## C3  2.000000 0.4000000 3 1.6 2.4
## C4  3.966667 0.5131601 3 3.4 4.4
##
## Alpha: 0.05 ; DF Error: 3
##
## Critical Range
##      2      3      4
## 1.546391 2.030514 2.344854
##
## Means with the same letter are not significantly different.
##
##      distancia groups
## C4  3.966667      a
## C3  2.000000      b
## C2  1.833333      b
## C1  1.200000      b
```

En la tabla se muestran los subgrupos formados de medias iguales al utilizar el método de Newman-Keuls. Hay dos subconjuntos que se diferencian entre sí. Por una parte el formado por el tipo de césped C4 y por otra parte el subgrupo formado por los tipos de césped: C3, C2 y C1. También se observa que el tipo de césped donde la distancia recorrida es más grande es el C4 con una distancia de 3.96667 y la distancia más corta es de 1.2 en el césped C1.

Ejercicio Propuesto 5

Un investigador quiere evaluar la productividad de cuatro variedades de aguacates, A, B, C y D. Para ello decide realizar el ensayo en un terreno que posee un gradiente de pendiente de oriente a occidente y además, diferencias en la disponibilidad de Nitrógeno de norte a sur, para controlar los efectos de la pendiente y la disponibilidad de Nitrógeno, utilizó un diseño de cuadrado latino, los datos corresponden a la producción en kg/parcela.

Responder a las siguientes cuestiones:

1. ¿Se puede afirmar que la productividad media de las cuatro variedades de aguacate es la misma?
2. ¿Qué supuestos han de verificarse?
3. ¿Se obtiene la misma producción con las cuatro variedades de aguacate?
En caso negativo, analizar mediante el procedimiento de Tukey y Newman-Keuls, con qué variedad de aguacate hay mayor producción.

Solución: 1. ¿Se puede afirmar que la productividad media de las cuatro variedades de aguacate es la misma?

El análisis de la productividad de las variedades de aguacate corresponde al análisis de un factor con 4 niveles. Dado que en el estudio intervienen dos fuentes de variación: la Disponibilidad de Nitrógeno y la Pendiente, se consideran dos factores de bloque, cada uno de ellos con 4 niveles.

Se pretende, entonces dar respuesta al contraste:

$$H_0 : \mu_A = \mu_B = \mu_C = \mu_D$$

$$H_1 : \mu_i \neq \mu_j \text{ para alguna } i \neq j$$

- Variable respuesta: Productividad
- Factor: Variedad de aguacate. Es un factor de efectos fijos ya que desde el principio se establecen los niveles concretos que se van a analizar.
- Bloques: Disponibilidad de Nitrógeno y Pendiente, ambos con 4 niveles y ambos de efectos fijos.
- Tamaño del experimento: Número total de observaciones (42) .

```
propuesto5 <- read.table("propuesto5.txt", header = TRUE)
propuesto5
```

##	productividad	aguacate	nitrogeno	pendiente
## 1	785	D	N1	P1
## 2	730	A	N1	P2
## 3	700	C	N1	P3
## 4	595	B	N1	P4
## 5	855	A	N2	P1
## 6	775	B	N2	P2
## 7	760	D	N2	P3
## 8	710	C	N2	P4
## 9	950	C	N3	P1
## 10	885	D	N3	P2
## 11	795	B	N3	P3
## 12	780	A	N3	P4
## 13	945	B	N4	P1

## 14	950	C	N4	P2
## 15	880	A	N4	P3
## 16	835	D	N4	P4

A continuación debemos transformar tanto la columna de los tratamiento como la de los bloques en un factor para poder realizar los cálculos posteriores adecuadamente.

```
propuesto5$latina = factor(propuesto5$aguacate)
propuesto5$latina

## [1] D A C B A B D C C D B A B C A D
## Levels: A B C D
```

```
propuesto5$Bloque1 = factor(propuesto5$nitrogeno)
propuesto5$Bloque1

## [1] N1 N1 N1 N1 N2 N2 N2 N2 N3 N3 N3 N3 N4 N4 N4 N4
## Levels: N1 N2 N3 N4
```

```
propuesto5$Bloque2 = factor(propuesto5$pendiente)
propuesto5$Bloque2

## [1] P1 P2 P3 P4 P1 P2 P3 P4 P1 P2 P3 P4 P1 P2 P3 P4
## Levels: P1 P2 P3 P4
```

Para calcular la tabla ANOVA primero hacemos uso de la función "aov" de la siguiente forma:

```
modp5 <- aov(productividad ~ aguacate + nitrogeno +
              pendiente, data = propuesto5 )
```

donde:

- productividad: Nombre de la columna de las observaciones
- aguacate : Nombre de la columna en la que están representados los tratamientos
- nitrogeno : Nombre de la columna en la que está representado el primer factor bloque
- pendiente: Nombre de la columna en la que está representado el segundo factor bloque (letras latinas)
- data = data.frame en el que están guardados los datos

```

modp5 <- aov(productividad ~ aguacate + nitrogeno +
             pendiente, data = propuesto5 )
modp5

## Call:
## aov(formula = productividad ~ aguacate + nitrogeno + pendiente,
## data = propuesto5)
##
## Terms:
##               aguacate nitrogeno pendiente Residuals
## Sum of Squares   5556.25  92518.75  52556.25    112.50
## Deg. of Freedom      3         3         3         6
##
## Residual standard error: 4.330127
## Estimated effects may be unbalanced

```

y posteriormente mostramos un resumen de los resultados con la función "summary" (verdadera tabla ANOVA):

```

summary(modp5)

##              Df Sum Sq Mean Sq F value    Pr(>F)
## aguacate      3   5556    1852    98.78 1.70e-05 ***
## nitrogeno     3  92519   30840  1644.78 3.92e-09 ***
## pendiente     3  52556   17519   934.33 2.13e-08 ***
## Residuals     6    113      19
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Observando los valores de los p-valores, 1.70e-05, 3.92e-09 y 2.13e-08; menores respectivamente que el nivel de significación del 5%, deducimos que los tres efectos son significativos. Tanto las variedades de aguacates utilizadas, como la pendiente del terreno y la disponibilidad de nitrógeno influyen en la productividad de los aguacates

2. ¿Qué supuestos han de verificarse?

Los supuestos que han de verificarse en un diseño de cuadrados latinos son Normalidad, Homocedasticidad e Independencia además del supuesto de aditividad entre filas, columnas y tratamientos (es decir, que no haya interacciones entre los mismos).

Hipótesis de normalidad

Comprobamos la hipótesis de normalidad mediante el análisis de la normalidad de los residuos. Para ello, hacemos uso del test de Shapiro-Wilks:

```
shapiro.test(mod$residuals)

##
##  Shapiro-Wilk normality test
##
## data:  mod$residuals
## W = 0.97552, p-value = 0.5939
```

Observamos el contraste de Shapiro-Wilk que es adecuado cuando las muestras son pequeñas ($n \leq 50$) y es una alternativa más potente que el test de Kolmogorov-Smirnov. El p-valor (0.07616) es mayor que el nivel de significación del 5%, aceptándose la hipótesis de normalidad.

Homogeneidad de varianzas

En este caso hacemos uso del Test de Barlett para contrastar la igualdad entre varianzas del factor.

```
bartlett.test(propuesto5$productividad, propuesto5$aguacate)

##
##  Bartlett test of homogeneity of variances
##
## data:  propuesto5$productividad and propuesto5$aguacate
## Bartlett's K-squared = 3.355, df = 3, p-value = 0.3401
```

El p-valor es 0.3401 por lo tanto no se puede rechazar la hipótesis de homogeneidad de las varianzas y se concluye que la cuatro variedades tienen varianzas homogéneas.

```
bartlett.test(propuesto5$productividad, propuesto5$nitrogeno)

##
##  Bartlett test of homogeneity of variances
##
## data:  propuesto5$productividad and propuesto5$nitrogeno
## Bartlett's K-squared = 0.56609, df = 3, p-value = 0.9041
```

```
bartlett.test(propuesto5$productividad, propuesto5$pendiente)

##
##  Bartlett test of homogeneity of variances
```



```
##
## data:  propuesto5$productividad and propuesto5$pendiente
## Bartlett's K-squared = 0.4168, df = 3, p-value = 0.9368
```

Los p-valores son mayores que 0.05, por lo tanto no se puede rechazar la hipótesis de homogeneidad de las varianzas en ninguno de los factores bloques

3. ¿Se obtiene la misma producción con las cuatro variedades de aguacate? En caso negativo, analizar mediante el procedimiento de Tukey y Newman-Keuls, con qué variedad de aguacate hay mayor producción.

```
TukeyHSD(mod, "aguacate", ordered = TRUE)

## Error in TukeyHSD.aov(mod, "aguacate", ordered = TRUE): 'which'
## specified no factors
```

Como se puede observar, todos los intervalos de confianza construidos para las diferencias entre las producciones medias de las variedades no contienen al 0, excepto el correspondiente a la pareja de variedades de aguacates A y D. Lo que significa que todas las producciones medias pueden considerarse distintas estadísticamente excepto las producciones medias correspondientes a las variedades A y D. Se deduce que únicamente no se observan diferencias significativas entre las producciones de las variedades de aguacates A y D (P-valor = 0.4289199).

La tabla de comparaciones múltiples muestra los intervalos simultáneos construidos por el método de Tukey para cada posible combinación de variedades de aguacates. En la tabla se muestra un resumen de las comparaciones de cada tratamiento con los restantes. Es decir, aparecen comparadas dos a dos las cuatro medias de los tratamientos.

```
plot(TukeyHSD(mod, "aguacate"))

## Error in h(simpleError(msg, call)): error in evaluating the argument
## 'x' in selecting a method for function 'plot': 'which' specified no
## factors
```

```
Newman_Keuls <- SNK.test(mod,"aguacate", console=TRUE)

## Name:  aguacate
##  agua planta
```

El contraste de Newman-keuls muestra una tabla de Subconjuntos homogéneos. En nuestro estudio sobre las producciones de aguacates se observan que hay tres subgrupos homogéneos, al primer subgrupo pertenece la Variedad C, con una producción de 827.50 media kg/parcela, el segundo las variedades D y A, con

producciones medias de 816.25 y 811.25 Kg(parcela, respectivamente y el tercero la Variedad B, con una producción media de 777.50 kg/parcela. Por lo tanto, se observa que la producción media mayor se obtiene con la Variedad C (827.5 Kg/ parcela) y la menor con la Variedad B (777.50 Kg/parcela).

Ejercicio Propuesto 6

Consideremos de nuevo el ejercicio propuesto 5 del investigador que quiere evaluar la productividad de cuatro variedades de aguacate, A, B, C y D. Para ello, decide realizar el ensayo en un terreno que posee un gradiente de pendiente de oriente a occidente y además, diferencias en la disponibilidad de Nitrógeno de norte a sur. Se seleccionan cuatro disponibilidades de nitrógeno, pero sólo dispone de tres gradientes de pendiente. Para controlar estas posibles fuentes de variabilidad, el investigador decide utilizar un diseño en cuadrado de Youden con cuatro filas, las cuatro disponibilidades de Nitrógeno (N1, N2, N3, N4), tres columnas, los tres gradientes de pendientes (P1, P2, P3) y cuatro letras latinas, las variedades de aguacates (A, B, C, D). Los datos corresponden a la producción en kg/parcela.

Nitrógeno	Pendiente		
	P1	P2	P3
N1	D	A	C
	956	820	689
N2	A	B	D
	867	975	680
N3	C	D	B
	850	775	699
N4	B	C	A
	950	870	980

Responder a las siguientes cuestiones:

1. Estudiar cuál es el tipo de diseño adecuado a este experimento y escribir el modelo matemático asociado.
2. ¿Se puede afirmar que la productividad media de las cuatro variedades de aguacate es la misma?
3. Comprobar la hipótesis de homocedasticidad
4. ¿Se obtiene la misma producción con las cuatro variedades de aguacate? En caso negativo, analizar mediante el procedimiento de Duncan, con qué variedad de aguacate hay mayor producción.

Solución:

1. Estudiar cuál es el tipo de diseño adecuado a este experimento y escribir el modelo matemático asociado.

El análisis de la productividad de las variedades de aguacate corresponde al análisis de un factor con 4 niveles. Dado que en el estudio intervienen dos fuentes de variación: la Disponibilidad de Nitrógeno y la Pendiente, se consideran dos factores de bloque, el primero con 4 niveles y el segundo con tres niveles.

Se pretende, entonces dar respuesta al contraste:

$$H_0 : \mu_A = \mu_B = \mu_C = \mu_D$$

$$H_1 : \mu_i \neq \mu_j$$

- Variable respuesta: Productividad.
 - Factor: Variedad de aguacate. Es un factor de efectos fijos ya que desde el principio se establecen los niveles concretos que se van a analizar.
 - Bloques: Disponibilidad de Nitrógeno y Pendiente, con 4 y 3 niveles, respectivamente y ambos de efectos fijos.
 - Tamaño del experimento: Número total de observaciones: 12.
2. ¿Se puede afirmar que la productividad media de las cuatro variedades de aguacate es la misma?

```
propuesto6<-read.table("propuesto6.txt", header = TRUE)
propuesto6
```

##	Productividad	Nitrogeno	Pendiente	Variedad
## 1	756	N1	P1	D
## 2	720	N1	P2	A
## 3	689	N1	P3	C
## 4	596	N2	P1	A
## 5	855	N2	P2	B
## 6	780	N2	P3	D
## 7	750	N3	P1	C
## 8	975	N3	P2	D
## 9	899	N3	P3	B
## 10	950	N4	P1	B
## 11	870	N4	P2	C
## 12	880	N4	P3	A

A continuación debemos transformar tanto la columna de los tratamientos como la de los bloques en un factor para poder realizar los cálculos posteriores adecuadamente.

```
propuesto6$Nitrogeno <- factor(propuesto6$Nitrogeno)
propuesto6$Nitrogeno

## [1] N1 N1 N1 N2 N2 N2 N3 N3 N3 N4 N4 N4
## Levels: N1 N2 N3 N4
```

```
propuesto6$Pendiente <- factor(propuesto6$Pendiente)
propuesto6$Pendiente

## [1] P1 P2 P3 P1 P2 P3 P1 P2 P3 P1 P2 P3
## Levels: P1 P2 P3
```

```
propuesto6$Variedad <- factor(propuesto6$Variedad)
propuesto6$Variedad

## [1] D A C A B D C D B B C A
## Levels: A B C D
```

Para poder analizar los datos mediante un diseño BIB debemos instalar y cargar los paquetes de R especializados en este tipo de diseños:

```
library(daewr)
library(AlgDesign)
```

La función "BIBsize(t , k)" de la librería daewr nos permite saber si el diseño puede realizarse. Calcula los parámetros del diseño donde

- t = número de niveles del factor tratamiento.
- k = número de tratamientos por bloque.

```
BIBsize(t = 4 , k = 3)

## Posible BIB design with b= 4 and r= 3 lambda= 2
```

El análisis de este modelo lo podemos realizar en R de dos formas:

1. Realizaremos el análisis evaluando primero el efecto del tratamiento y después el de los bloques utilizando tres funciones

Para cada factor realizamos una tabla ANOVA:

- Factor principal: Variedad

Para evaluar el efecto de los tratamientos, la suma de cuadrados de tratamientos debe ajustarse por bloques, por lo tanto primero se introducen los bloques y después los tratamientos.

Para calcular la tabla ANOVA hacemos uso de la función "aov" (asume suma de cuadrados tipo I) de la siguiente forma:

```
modp6 <- aov(Productividad ~ Pendiente + Nitrogeno +
              Variedad, data = propuesto6)
```

donde:

- Productividad: Nombre de la columna de las observaciones.
- Variedad: Nombre de la columna en la que están representados los tratamientos.
- Pendiente: Nombre de la columna en la que está representado el primer factor bloque.
- Nitrogeno: Nombre de la columna en la que está representado el segundo factor bloque (letras latinas).
- data = data.frame en el que están guardados los datos.

```
modp6

## Call:
##   aov(formula = Productividad ~ Pendiente + Nitrogeno + Variedad,
##       data = propuesto6)
##
## Terms:
##               Pendiente Nitrogeno Variedad Residuals
## Sum of Squares   16952.00   73454.00  47805.25    3012.75
## Deg. of Freedom         2         3         3         3
##
## Residual standard error: 31.6899
## Estimated effects may be unbalanced
```

```
summary(modp6)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Pendiente     2  16952    8476   8.44 0.0586 .
## Nitrogeno     3  73454   24485  24.38 0.0131 *
## Variedad      3  47805   15935  15.87 0.0241 *
## Residuals     3   3013    1004
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El p-valor, 0.0241, es menor que el nivel de significación del 5%, deducimos que el factor principal: Variedades del aguacate es significativo.

- Factor Bloque: Pendiente

Para evaluar el efecto del primero de los bloques, la suma de cuadrados de bloques debe ajustarse por los tratamientos, por lo tanto primero se introducen los tratamientos y después los bloques:

```
modp62 <- aov(Productividad~ Variedad + Nitrogeno +
               Pendiente, data = propuesto6)
modp62

## Call:
## aov(formula = Productividad ~ Variedad + Nitrogeno + Pendiente,
## data = propuesto6)
##
## Terms:
##              Variedad Nitrogeno Pendiente Residuals
## Sum of Squares  50344.67  70914.58  16952.00   3012.75
## Deg. of Freedom      3         3         2         3
##
## Residual standard error: 31.6899
## Estimated effects may be unbalanced

summary(modp62)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Variedad      3  50345   16782    16.71 0.0224 *
## Nitrogeno     3  70915   23638    23.54 0.0138 *
## Pendiente     2  16952    8476     8.44 0.0586 .
## Residuals     3   3013    1004
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El p-valor, 0.0586, es mayor que el nivel de significación del 5%, deducimos que el Factor Bloque: Pendiente no es significativo.

- Factor Bloque: Nitrogeno

Para evaluar el efecto del segundo bloque, la suma de cuadrados de bloques debe ajustarse también por los tratamientos, por lo tanto primero se introducen los tratamientos y después los bloques:

```
modp63 <- aov(Productividad ~ Variedad + Pendiente +
               Nitrogeno, data = propuesto6)
modp63
```

```
## Call:
## aov(formula = Productividad ~ Variedad + Pendiente + Nitrogeno,
## data = propuesto6)
##
## Terms:
##             Variedad Pendiente Nitrogeno Residuals
## Sum of Squares 50344.67 16952.00 70914.58 3012.75
## Deg. of Freedom      3      2      3      3
##
## Residual standard error: 31.6899
## Estimated effects may be unbalanced

summary(modp63)

##           Df Sum Sq Mean Sq F value Pr(>F)
## Variedad   3 50345    16782   16.71 0.0224 *
## Pendiente  2 16952     8476    8.44 0.0586 .
## Nitrogeno  3 70915    23638   23.54 0.0138 *
## Residuals  3  3013     1004
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El p-valor es 0.0138; menor que el nivel de significación del 5%, deducimos que Factor Bloque: Nitrógeno es significativo.

2. Realizaremos el análisis evaluando tanto los tratamientos como los bloques ejecutando solo una función. Para ello necesitamos instalar y cargar el paquete "car"

```
modp64 <- lm(Productividad ~ Variedad + Nitrogeno +
              Pendiente, data = propuesto6)
modp64

##
## Call:
## lm(formula = Productividad ~ Variedad + Nitrogeno + Pendiente,
## data = propuesto6)
##
## Coefficients:
## (Intercept)  VariedadB  VariedadC  VariedadD  NitrogenoN2  NitrogenoN3
##      634.38      133.12       -6.75      127.62      -24.63      108.62
## NitrogenoN4  PendienteP2  PendienteP3
##      176.50       92.00       49.00
```

```
car::Anova(mod4, type="III")

## Error in car::Anova(mod4, type = "III"): objeto 'mod4' no encontrado
```

Los resultados obtenidos coinciden con los realizados primero a los bloques y después al tratamiento.

3. Comprobar la hipótesis de homocedasticidad

```
bartlett.test(propuesto6$Productividad, propuesto6$Variedad)

##
## Bartlett test of homogeneity of variances
##
## data: propuesto6$Productividad and propuesto6$Variedad
## Bartlett's K-squared = 1.8021, df = 3, p-value = 0.6145
```

```
bartlett.test(propuesto6$Productividad, propuesto6$Pendiente)

##
## Bartlett test of homogeneity of variances
##
## data: propuesto6$Productividad and propuesto6$Pendiente
## Bartlett's K-squared = 0.49855, df = 2, p-value = 0.7794
```

```
bartlett.test(propuesto6$Productividad, propuesto6$Nitrogeno)

##
## Bartlett test of homogeneity of variances
##
## data: propuesto6$Productividad and propuesto6$Nitrogeno
## Bartlett's K-squared = 3.8699, df = 3, p-value = 0.2759
```

Los p-valores del factor tratamiento, Variedad de aguacate (0.6145), del factor bloque Pendiente (0.779) y del factor bloque Nitrógeno (0.2759) son mayores que 0.05, por lo tanto no se puede rechazar la hipótesis de homogeneidad de la varianza del tratamiento y de los bloques.

4. ¿Se obtiene la misma producción con las cuatro variedades de aguacate? En caso negativo, analizar mediante el procedimiento de Duncan, con qué variedad de aguacate hay mayor producción.

Para realizar el contraste de comparaciones múltiples hay que cargar e instalar el paquete agricolae


```

library(agricolae)
(duncan=duncan.test(modp64, "Variedad" , group = T))

## $statistics
##      MSerror Df Mean      CV
##    1004.25   3   810 3.912334
##
## $parameters
##      test  name.t ntr alpha
##    Duncan Variedad   4  0.05
##
## $duncan
##      Table CriticalRange
##    2 4.500659      82.34484
##    3 4.515652      82.61915
##    4 4.472854      81.83611
##
## $means
##      Productividad      std r Min Max   Q25 Q50   Q75
##    A      732.0000 142.37977 3 596 880 658.0 720 800.0
##    B      901.3333  47.54296 3 855 950 877.0 899 924.5
##    C      769.6667  92.08873 3 689 870 719.5 750 810.0
##    D      837.0000 120.11245 3 756 975 768.0 780 877.5
##
## $comparison
## NULL
##
## $groups
##      Productividad groups
##    B      901.3333      a
##    D      837.0000     ab
##    C      769.6667     bc
##    A      732.0000      c
##
## attr(,"class")
## [1] "group"

```

En la tabla se muestran los subgrupos formados de medias iguales al utilizar el método de Duncan. Hay tres subconjuntos que se diferencian entre sí. Por una parte el formado por la variedad de aguacate B y D, el subgrupo formado por D y C y el formado por A y C. También se observa que la mayor productividad de aguacate es la del tipo B, con una producción de 901.3333 Kg por parcela y la menor el tipo A, 732.0000 kg por parcela.

Ejercicio Propuesto 7

En un invernadero se está estudiando el crecimiento de determinadas plantas, para ello se quiere controlar los efectos del terreno, abono, insecticida y semilla. El estudio se realiza con cuatro tipos de semillas diferentes que se plantan en cuatro tipos de terreno, se les aplican cuatro tipos de abonos y cuatro tipos de insecticidas. La asignación de los tratamientos a las plantas se realiza de forma aleatoria. Para controlar estas posibles fuentes de variabilidad se decide plantear un diseño por cuadrados greco-latinos como el que se muestra en la siguiente tabla, donde las letras griegas corresponden a los cuatro tipos de semilla y las latinas a los abonos.

Responder a las siguientes cuestiones:

1. Estudiar cuál es el tipo de diseño adecuado a este experimento y escribir el modelo matemático asociado.
2. ¿Se puede afirmar que el crecimiento de las plantas es el mismo para los cuatro tipos de abonos? ¿Y con los distintos insecticidas?
3. ¿Existen diferencias significativas en el crecimiento de las plantas con las distintas semillas? ¿Y el tipo de tierra influye en dicho crecimiento?
4. ¿Con qué tipo de semilla se produce el mayor crecimiento de las plantas?
5. ¿El crecimiento de las plantas es el mismo utilizando al mismo tiempo los abonos A y B que utilizando los abonos C y D?

Solución:

1. Estudiar cuál es el tipo de diseño adecuado a este experimento y escribir el modelo matemático asociado.

Es un diseño en cuadrado greco-latino, la variable respuesta $y_{(ij(hp))}$ viene descrita por la siguiente ecuación

$$y_{(ij(hp))} = \mu + \tau_i + \beta_j + \gamma_k + \delta_p + \varepsilon_{(ij(hp))}; i, j, k, h, p = 1, 2, 3, 4$$

La variable respuesta que vamos a estudiar es el crecimiento de determinadas plantas. El factor principal es tipo de abono que se presenta con cuatro niveles.

- Variable respuesta: Crecimiento
- Factor: Tipo.abonos que tiene cuatro niveles. Es un factor de efectos fijos ya que viene decidido que niveles concretos se van a utilizar.
- Bloques: Insecticidas, Terrenos y Semillas, cada uno con cuatro niveles y de efectos fijos.
- Tamaño del experimento: Número total de observaciones (16).

Para realizar este supuesto en R debemos introducir primero los datos de forma correcta. Podemos introducir los datos directamente en R de forma manual o introducirlos previamente en un archivo de texto o Excel y leerlos en R. En este caso lo hacemos en un archivo de texto.

```
propuesto7 <- read.table("propuesto7.txt", header = TRUE)
propuesto7
```

	Crecimiento	Tipo_abono	Tipo_semilla	Tipo_insecticida	Tipo_terreno
## 1	6	C	beta	Insecticida1	Terreno1
## 2	12	B	alfa	Insecticida2	Terreno1
## 3	13	A	delta	Insecticida3	Terreno1
## 4	13	D	gamma	Insecticida4	Terreno1
## 5	6	B	gamma	Insecticida1	Terreno2
## 6	10	C	delta	Insecticida2	Terreno2
## 7	16	D	alfa	Insecticida3	Terreno2
## 8	11	A	beta	Insecticida4	Terreno2
## 9	7	D	delta	Insecticida1	Terreno3
## 10	5	A	gamma	Insecticida2	Terreno3
## 11	5	B	beta	Insecticida3	Terreno3
## 12	7	C	alfa	Insecticida4	Terreno3
## 13	11	A	alfa	Insecticida1	Terreno4
## 14	11	D	beta	Insecticida2	Terreno4
## 15	8	C	gamma	Insecticida3	Terreno4
## 16	9	B	delta	Insecticida4	Terreno4

A continuación debemos transformar tanto la columna de los tratamientos como la de los bloques en un factor para poder realizar los cálculos posteriores adecuadamente.

```
propuesto7$Tipo_abono <- factor(propuesto7$Tipo_abono)
propuesto7$Tipo_abono
```

```
## [1] C B A D B C D A D A B C A D C B
## Levels: A B C D
```

```
propuesto7$Tipo_semilla <- factor(propuesto7$Tipo_semilla)
propuesto7$Tipo_semilla
```

```
## [1] beta alfa delta gamma gamma delta alfa beta delta gamma beta alfa
## [13] alfa beta gamma delta
## Levels: alfa beta delta gamma
```

```
propuesto7$Tipo_insecticida <- factor(propuesto7$Tipo_insecticida)
propuesto7$Tipo_insecticida

## [1] Insecticida1 Insecticida2 Insecticida3 Insecticida4 Insecticida1
## [6] Insecticida2 Insecticida3 Insecticida4 Insecticida1 Insecticida2
## [11] Insecticida3 Insecticida4 Insecticida1 Insecticida2 Insecticida3
## [16] Insecticida4
## Levels: Insecticida1 Insecticida2 Insecticida3 Insecticida4
```

```
propuesto7$Tipo_terreno <- factor(propuesto7$Tipo_terreno)
propuesto7$Tipo_terreno

## [1] Terreno1 Terreno1 Terreno1 Terreno1 Terreno2 Terreno2 Terreno2 Terreno2
## [9] Terreno3 Terreno3 Terreno3 Terreno3 Terreno4 Terreno4 Terreno4 Terreno4
## Levels: Terreno1 Terreno2 Terreno3 Terreno4
```

2. ¿Se puede afirmar que el crecimiento de las plantas es el mismo para los cuatro tipos de abonos? ¿Y con los distintos insecticidas?

Primero vamos a obtener la tabla ANOVA

```
modp7 <- aov(Crecimiento~ Tipo_abono + Tipo_semilla +
              Tipo_insecticida + Tipo_terreno,
              data = propuesto7)
modp7

## Call:
## aov(formula = Crecimiento ~ Tipo_abono + Tipo_semilla + Tipo_insecticida +
## Tipo_terreno, data = propuesto7)
##
## Terms:
##              Tipo_abono Tipo_semilla Tipo_insecticida Tipo_terreno Residuals
## Sum of Squares      42.25      31.25      20.75      64.25      1.25
## Deg. of Freedom        3         3         3         3         3
##
## Residual standard error: 0.6454972
## Estimated effects may be unbalanced

summary(modp7)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Tipo_abono    3  42.25  14.083   33.8 0.00820 **
## Tipo_semilla  3  31.25  10.417   25.0 0.01266 *
## Tipo_insecticida 3  20.75   6.917   16.6 0.02260 *
## Tipo_terreno  3  64.25  21.417   51.4 0.00445 **
## Residuals    3   1.25   0.417
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Son significativos los efectos de todos los factores Para determinar si el crecimiento de las plantas es el mismo para los cuatro tipos de abonos y con los distintos insecticidas, vamos a realizar contrastes de comparaciones múltiples, por ejemplo de Duncan para los tipos de abonos y Newman Keuls para los insecticidas. Recordar que hay que cargar e instalar el paquete agricolae

```
library(agricolae)
(duncan=duncan.test(modp7, "Tipo_abono" , group = T))

## $statistics
##      MSerror Df  Mean      CV
## 0.4166667  3 9.375 6.885304
##
## $parameters
##      test      name.t ntr alpha
##  Duncan Tipo_abono   4  0.05
##
## $duncan
##      Table CriticalRange
## 2 4.500659      1.452581
## 3 4.515652      1.457420
## 4 4.472854      1.443607
##
## $means
##      Crecimiento      std r Min Max   Q25  Q50  Q75
## A      10.00 3.464102 4  5  13  9.50 11.0 11.50
## B       8.00 3.162278 4  5  12  5.75  7.5  9.75
## C       7.75 1.707825 4  6  10  6.75  7.5  8.50
## D      11.75 3.774917 4  7  16 10.00 12.0 13.75
##
## $comparison
## NULL
##
## $groups
##      Crecimiento groups
## D      11.75      a
## A      10.00      b
## B       8.00      c
## C       7.75      c
##
## attr("class")
## [1] "group"
```

El mayor crecimiento de las plantas se produce con el Abono D siendo la altura que alcanza de 11.75 y la menor altura (7.75) la alcanza con el tipo de abono C

```
Newman_Keuls <- SNK.test(modp7,"Tipo_insecticida",
                          console=TRUE,
                          main=" Contraste de Newman-Keuls para el factor tipo de insecticida")

##
## Study:  Contraste de Newman-Keuls para el factor tipo de insecticida
##
## Student Newman Keuls Test
## for Crecimiento
##
## Mean Square Error:  0.4166667
##
## Tipo_insecticida,  means
##
##           Crecimiento      std r Min Max
## Insecticida1         7.5 2.380476 4   6  11
## Insecticida2         9.5 3.109126 4   5  12
## Insecticida3        10.5 4.932883 4   5  16
## Insecticida4        10.0 2.581989 4   7  13
##
## Alpha: 0.05 ; DF Error: 3
##
## Critical Range
##           2           3           4
## 1.452581 1.907336 2.202606
##
## Means with the same letter are not significantly different.
##
##           Crecimiento groups
## Insecticida3         10.5      a
## Insecticida4         10.0      a
## Insecticida2          9.5      a
## Insecticida1          7.5      b
```

La menor altura (7.5) la alcanza cuando se le suministra el Insecticida 1.

3. ¿Existen diferencias significativas en el crecimiento de las plantas con las distintas semillas? ¿Y el tipo de tierra influye en dicho crecimiento?

Aplicamos el contraste de comparaciones múltiple de Tukey

```
mod.tukey<- TukeyHSD(modp7, "Tipo_semilla", ordered = TRUE)
mod.tukey

##    Tukey multiple comparisons of means
##      95% family-wise confidence level
##      factor levels have been ordered
##
## Fit: aov(formula = Crecimiento ~ Tipo_abono + Tipo_semilla + Tipo_insecticida + Tipo_terr
##
## $Tipo_semilla
##           diff           lwr           upr           p adj
## beta-gamma  0.25 -1.9526064  2.452606  0.9412719
## delta-gamma 1.75 -0.4526064  3.952606  0.0901983
## alfa-gamma   3.50  1.2973936  5.702606  0.0139151
## delta-beta   1.50 -0.7026064  3.702606  0.1305850
## alfa-beta    3.25  1.0473936  5.452606  0.0171739
## alfa-delta   1.75 -0.4526064  3.952606  0.0901983
```

Comprobamos que únicamente hay diferencias significativas entre los tipos de semillas alfa-gamma (p-valor = 0.0139) y entre alfa-beta (p-valor = 0.0171).

Para determinar si el tipo de terreno influye en el crecimiento aplicamos el método de comparaciones múltiples LSD.

```
LSD.test(modp7,"Tipo_terreno", p.adj="bonferroni", console=TRUE)

##
## Study: modp7 ~ "Tipo_terreno"
##
## LSD t Test for Crecimiento
## P value adjustment method: bonferroni
##
## Mean Square Error:  0.4166667
##
## Tipo_terreno, means and individual ( 95 %) CI
##
##           Crecimiento      std r      LCL      UCL Min Max
## Terreno1      11.00 3.366502 4 9.97287 12.02713   6 13
## Terreno2      10.75 4.112988 4 9.72287 11.77713   6 16
## Terreno3       6.00 1.154701 4 4.97287  7.02713   5  7
## Terreno4       9.75 1.500000 4 8.72287 10.77713   8 11
##
## Alpha: 0.05 ; DF Error: 3
## Critical Value of t: 6.231543
##
## Minimum Significant Difference: 2.844297
##
```

```
## Treatments with the same letter are not significantly different.
##
##          Crecimiento groups
## Terreno1      11.00      a
## Terreno2      10.75      a
## Terreno4       9.75      a
## Terreno3       6.00      b
```

Hay dos grupos de terrenos que difieren significativamente, por un lado el grupo formado por los tipos de terrenos 1,2 y 4 y por otra parte el grupo formado un solo tipo de terreno, el tipo de terreno 3. El mayor crecimiento de las plantas se produce en el tipo de terreno 1 con un crecimiento de 11 u.c.

4. ¿Con qué tipo de semilla se produce el mayor crecimiento de las plantas?

```
Newman_Keuls1 <- SNK.test(modp7,"Tipo_semilla",
                          console=TRUE,
                          main=" Contraste de Newman-Keuls para el factor tipo de semilla ")

##
## Study:  Contraste de Newman-Keuls para el factor tipo de semilla
##
## Student Newman Keuls Test
## for Crecimiento
##
## Mean Square Error:  0.4166667
##
## Tipo_semilla,  means
##
##          Crecimiento      std r Min Max
## alfa          11.50 3.696846 4    7  16
## beta           8.25 3.201562 4    5  11
## delta          9.75 2.500000 4    7  13
## gamma          8.00 3.559026 4    5  13
##
## Alpha: 0.05 ; DF Error: 3
##
## Critical Range
##          2          3          4
## 1.452581 1.907336 2.202606
##
## Means with the same letter are not significantly different.
##
##          Crecimiento groups
## alfa          11.50      a
## delta          9.75      b
```



```
## beta      8.25      c
## gamma     8.00      c
```

El mayor crecimiento de la planta se produce con el tipo de semilla alfa (11.50 u.c.)

5. ¿El crecimiento de las plantas es el mismo utilizando al mismo tiempo los abonos A y B que utilizando los abonos C y D?

```
TukeyHSD(modp7, "Tipo_abono", ordered = TRUE)

##    Tukey multiple comparisons of means
##      95% family-wise confidence level
##      factor levels have been ordered
##
## Fit: aov(formula = Crecimiento ~ Tipo_abono + Tipo_semilla + Tipo_insecticida + Tipo_terro)
##
## $Tipo_abono
##      diff      lwr      upr      p adj
## B-C 0.25 -1.95260644  2.452606  0.9412719
## A-C 2.25  0.04739356  4.452606  0.0472540
## D-C 4.00  1.79739356  6.202606  0.0094879
## A-B 2.00 -0.20260644  4.202606  0.0643491
## D-B 3.75  1.54739356  5.952606  0.0114235
## D-A 1.75 -0.45260644  3.952606  0.0901983
```

Hay diferencias significativas en el crecimiento utilizando los abonos C y D (P-valor = 0.0094), pero no las hay con los abonos A y B (P-valor= 0.064).

Ejercicio Propuesto 8

Se realiza un estudio sobre el efecto que produce la descarga de aguas residuales de una planta sobre la ecología del agua natural de un río. En el estudio se utilizaron dos lugares de muestreo. Un lugar está aguas arriba del punto en el que la planta introduce aguas residuales en la corriente; el otro está aguas abajo. Se tomaron muestras durante un periodo de cuatro semanas y se obtuvieron los datos sobre el número de diatomeas halladas. Los datos se muestran en la tabla adjunta:

	Semanas			
Lugar	Semana 1	Semana 2	Semana 3	Semana 4
Aguas	78 94	620 760	204 333	890 655
Arriba	43 58	420 913	98 89	763 562
Aguas	79 87	546 652	45 69	254 86
abajo	145 522	76 94	59 62	789 267

Responder a las siguientes cuestiones:

1. Identificar el diseño adecuado a este experimento, escribir el modelo matemático y explicar los distintos elementos que intervienen.
2. Estudiar si la semana y el lugar son factores determinantes en el número de diatomeas halladas en el agua del río. ¿Hay posibilidad que una semana sea más recomendable en un lugar del río en concreto y no lo sea en el otro lugar?
3. Estudiar en qué semana se producen menos contaminación en el río, utilizando el método de Duncan.
4. Estudiar en qué lugar del río se producen menos diatomeas.

Solución: Responder a las siguientes cuestiones: 1. Identificar el diseño adecuado a este experimento, escribir el modelo matemático y explicar los distintos elementos que intervienen.

Es un diseño factorial de dos factores con replicación, la variable respuesta y_{ijk} viene descrita por la siguiente ecuación

Modelo estadístico del diseño factorial de dos factores con replicación

$$y_{(ijk)} = \mu + \tau_i + \beta_j + (\tau\beta)_{ij} + u_{ijk}, i = 1, 2; j = 1, 2, 3, 4; k = 1, 2, 3, 4$$

En este caso lo hacemos en un archivo de texto:

- Variable respuesta: n_diatomeas
- Factor: semana que tiene cuatro niveles. Es un factor de efectos fijos ya que viene decidido que niveles concretos se van a utilizar.

- Factor: lugar con dos niveles y de efectos fijos.
- Tamaño del experimento: Número total de observaciones (32).

```
propuesto8<-read.table("propuesto8.txt", header = TRUE)
propuesto8
```

##	n_diatomeas	semana	lugar
## 1	78	Semana1	Aguas_arriba
## 2	94	Semana1	Aguas_arriba
## 3	620	Semana2	Aguas_arriba
## 4	760	Semana2	Aguas_arriba
## 5	204	Semana3	Aguas_arriba
## 6	333	Semana3	Aguas_arriba
## 7	890	Semana4	Aguas_arriba
## 8	655	Semana4	Aguas_arriba
## 9	43	Semana1	Aguas_arriba
## 10	58	Semana1	Aguas_arriba
## 11	420	Semana2	Aguas_arriba
## 12	913	Semana2	Aguas_arriba
## 13	98	Semana3	Aguas_arriba
## 14	89	Semana3	Aguas_arriba
## 15	763	Semana4	Aguas_arriba
## 16	562	Semana4	Aguas_arriba
## 17	79	Semana1	Aguas_abajo
## 18	87	Semana1	Aguas_abajo
## 19	546	Semana2	Aguas_abajo
## 20	652	Semana2	Aguas_abajo
## 21	45	Semana3	Aguas_abajo
## 22	69	Semana3	Aguas_abajo
## 23	254	Semana4	Aguas_abajo
## 24	86	Semana4	Aguas_abajo
## 25	145	Semana1	Aguas_abajo
## 26	522	Semana1	Aguas_abajo
## 27	76	Semana2	Aguas_abajo
## 28	94	Semana2	Aguas_abajo
## 29	59	Semana3	Aguas_abajo
## 30	62	Semana3	Aguas_abajo
## 31	789	Semana4	Aguas_abajo
## 32	267	Semana4	Aguas_abajo

A continuación debemos transformar tanto la columna de los tratamiento como la de los bloques en un factor para poder realizar los cálculos posteriores adecuadamente.

```
propuesto8$semana <- factor(propuesto8$semana)
propuesto8$semana

## [1] Semana1 Semana1 Semana2 Semana2 Semana3 Semana3 Semana4 Semana4 Semana1
## [10] Semana1 Semana2 Semana2 Semana3 Semana3 Semana4 Semana4 Semana1 Semana1
## [19] Semana2 Semana2 Semana3 Semana3 Semana4 Semana4 Semana1 Semana1 Semana2
## [28] Semana2 Semana3 Semana3 Semana4 Semana4
## Levels: Semana1 Semana2 Semana3 Semana4
```

```
propuesto8$lugar <- factor(propuesto8$lugar)
propuesto8$lugar

## [1] Aguas_arriba Aguas_arriba Aguas_arriba Aguas_arriba Aguas_arriba
## [6] Aguas_arriba Aguas_arriba Aguas_arriba Aguas_arriba Aguas_arriba
## [11] Aguas_arriba Aguas_arriba Aguas_arriba Aguas_arriba Aguas_arriba
## [16] Aguas_arriba Aguas_abajo Aguas_abajo Aguas_abajo Aguas_abajo
## [21] Aguas_abajo Aguas_abajo Aguas_abajo Aguas_abajo Aguas_abajo
## [26] Aguas_abajo Aguas_abajo Aguas_abajo Aguas_abajo Aguas_abajo
## [31] Aguas_abajo Aguas_abajo
## Levels: Aguas_abajo Aguas_arriba
```

2. Estudiar si la semana y el lugar son factores determinantes en el número de diatomeas halladas en el agua del río. ¿Hay posibilidad que una semana sea más recomendable en un lugar del río en concreto y no lo sea en el otro lugar?

Primero vamos a obtener la tabla ANOVA

```
modp8 <- aov(n_diatomeas ~ semana + lugar + semana * lugar,
             data = propuesto8)
modp8

## Call:
## aov(formula = n_diatomeas ~ semana + lugar + semana * lugar,
## data = propuesto8)
##
## Terms:
##          semana          lugar semana:lugar Residuals
## Sum of Squares 1236724.7 235984.5    330818.3 914922.0
## Deg. of Freedom      3          1          3        24
##
## Residual standard error: 195.2479
## Estimated effects may be unbalanced

summary(modp8)
```

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## semana      3 1236725  412242  10.814 0.00011 ***
## lugar       1  235985  235985    6.190 0.02018 *
## semana:lugar  3  330818  110273    2.893 0.05616 .
## Residuals   24  914922   38122
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Son significativos los efectos de los dos factores, pero no es significativo la interacción. Ambos son factores determinantes en el número de diatomeas hal-ladas en el agua del río.

Para determinar si hay algún día de la semana más recomendable que otro, vamos a aplicar un contraste de comparaciones múltiples

Recordar que hay que cargar e instalar el paquete agricolae

```
LSD.test(modp8,"semana", p.adj="bonferroni", console=TRUE)

##
## Study: modp8 ~ "semana"
##
## LSD t Test for n_diatomeas
## P value adjustment method: bonferroni
##
## Mean Square Error:  38121.75
##
## semana, means and individual ( 95 %) CI
##
##           n_diatomeas      std r      LCL      UCL Min Max
## Semana1      138.250 157.90933 8  -4.222082 280.7221  43 522
## Semana2      510.125 299.51887 8 367.652918 652.5971  76 913
## Semana3      119.875  99.41319 8 -22.597082 262.3471  45 333
## Semana4      533.250 295.20344 8 390.777918 675.7221  86 890
##
## Alpha: 0.05 ; DF Error: 24
## Critical Value of t: 2.875094
##
## Minimum Significant Difference: 280.6781
##
## Treatments with the same letter are not significantly different.
##
##           n_diatomeas groups
## Semana4      533.250      a
## Semana2      510.125      a
## Semana1      138.250      b
```

```
## Semana3      119.875      b
```

El número de diatomeas (533.250) es mayor en la semana 4 y menor en la semana 3 (119.875). Hay dos subgrupos diferenciados. Un subgrupo está formado por las semanas 2 y 4 y el otro subgrupo por las semanas 1 y 3.

```
SNK.test(modp8,"lugar", console=TRUE)

##
## Study: modp8 ~ "lugar"
##
## Student Newman Keuls Test
## for n_diatomeas
##
## Mean Square Error:  38121.75
##
## lugar,  means
##
##           n_diatomeas      std  r Min Max
## Aguas_abajo      239.50 245.9295 16  45 789
## Aguas_arriba     411.25 324.0622 16  43 913
##
## Alpha: 0.05 ; DF Error: 24
##
## Critical Range
##           2
## 142.4721
##
## Means with the same letter are not significantly different.
##
##           n_diatomeas groups
## Aguas_arriba     411.25      a
## Aguas_abajo      239.50      b
```

Se produce menos contaminación en aguas abajo, con un número de diatomeas de 239.50. Y como hemos dicho anteriormente hay diferencias significativas entre el número de diatomeas que se producen en aguas arribas del río y las que se producen en aguas abajo del río

3. Estudiar en qué semana se producen menos contaminación en el río, utilizando el método de Duncan.

```
(duncan=duncan.test(modp8, "semana" , group = T))

## $statistics
##      MSError Df      Mean      CV
```

```
## 38121.75 24 325.375 60.00704
##
## $parameters
## test name.t ntr alpha
## Duncan semana 4 0.05
##
## $duncan
## Table CriticalRange
## 2 2.918793 201.4859
## 3 3.065610 211.6208
## 4 3.159874 218.1279
##
## $means
## n_diatomeas std r Min Max Q25 Q50 Q75
## Semana1 138.250 157.90933 8 43 522 73.00 83.0 106.75
## Semana2 510.125 299.51887 8 76 913 338.50 583.0 679.00
## Semana3 119.875 99.41319 8 45 333 61.25 79.0 124.50
## Semana4 533.250 295.20344 8 86 890 263.75 608.5 769.50
##
## $comparison
## NULL
##
## $groups
## n_diatomeas groups
## Semana4 533.250 a
## Semana2 510.125 a
## Semana1 138.250 b
## Semana3 119.875 b
##
## attr(,"class")
## [1] "group"
```

El número de diatomeas (533.250) es mayor en la semana 4 y menor en la semana 3 (119.875)

Ejercicio Propuesto 9

La cotinina es uno de los principales metabolitos de la nicotina. Actualmente se le considera el mejor indicador de la exposición al humo de tabaco. Se ha realizado un estudio con distintas marcas de tabaco distinguiendo principalmente entre negro y rubio para detectar las posibles diferencias en el nivel de nicotina de personas expuestas al humo de tabaco. Para ello, se han analizado personas de distintas edades (niños, jóvenes y adultos) y se ha distinguido entre mujeres y hombres. Se han obtenido los datos de la siguiente tabla sobre el nivel de nicotina en miligramos por mililitro.

	Sexo			
	Hombres		Mujeres	
	Tabaco		Tabaco	
Edades	Rubio	Negro	Rubio	Negro
Niños	110	360	230	141
	240	125	219	123
Jóvenes	112	252	655	873
	239	455	432	256
Adultos	652	354	653	198
	451	701	259	343

Responder a las siguientes cuestiones:

- Identificar el diseño adecuado a este experimento, escribir el modelo matemático y explicar los distintos elementos que intervienen.
- Contrastar la hipótesis nula de no interacción entre los factores. Adecuar el modelo al resultado de las interacciones y contrastar los efectos principales.
- ¿Hay diferencias significativas en el nivel de nicotina en las distintas edades? ¿En qué edad el nivel de nicotina es mayor?
- ¿El tipo de tabaco es un factor determinante en el nivel de nicotina?
- Comparar el nivel medio de nicotina entre las mujeres y los hombres. ¿Se detectan diferencias significativas?

Solución:

1. Identificar el diseño adecuado a este experimento, escribir el modelo matemático y explicar los distintos elementos que intervienen

Es un diseño factorial de tres factores con replicación, la variable respuesta y_{ijkl} viene descrita por la siguiente ecuación

Modelo estadístico del diseño factorial de cuatro factores con replicación

$$y_{ijkl} = \mu + \tau_i + \beta_j + \gamma_k + (\tau\beta)_{ij} + (\tau\gamma)_{ik} + (\beta\gamma)_{jk} + (\tau\beta\gamma)_{ijk} + u_{ijkl}, i = 1, 2; j = 1, 2; k = 1, 2, 3; l = 1, 2$$

- Variable respuesta: Nivel_nicotina
- Factor: sexo que tiene dos niveles. Es un factor de efectos fijos ya que viene decidido que niveles concretos se van a utilizar.

- Factor: Tabaco con dos niveles y de efectos fijos.
- Factor: Edades con tres niveles y de efectos fijos.
- Tamaño del experimento: Número total de observaciones (24).

```
propuesto9<-read.table("propuesto9.txt", header = TRUE)
propuesto9
```

```
##      Nivel_nicotina      Sexo Tabaco  Edades
## 1             110 Hombres  Rubio   Niños
## 2             360 Hombres  Negro   Niños
## 3             230 Mujeres  Rubio   Niños
## 4             141 Mujeres  Negro   Niños
## 5             240 Hombres  Rubio   Niños
## 6             125 Hombres  Negro   Niños
## 7             219 Mujeres  Rubio   Niños
## 8             123 Mujeres  Negro   Niños
## 9             112 Hombres  Rubio  Jóvenes
## 10            252 Hombres  Negro  Jóvenes
## 11            655 Mujeres  Rubio  Jóvenes
## 12            873 Mujeres  Negro  Jóvenes
## 13            239 Hombres  Rubio  Jóvenes
## 14            455 Hombres  Negro  Jóvenes
## 15            432 Mujeres  Rubio  Jóvenes
## 16            256 Mujeres  Negro  Jóvenes
## 17            652 Hombres  Rubio  Adultos
## 18            354 Hombres  Negro  Adultos
## 19            653 Mujeres  Rubio  Adultos
## 20            198 Mujeres  Negro  Adultos
## 21            451 Hombres  Rubio  Adultos
## 22            701 Hombres  Negro  Adultos
## 23            259 Mujeres  Rubio  Adultos
## 24            343 Mujeres  Negro  Adultos
```

A continuación debemos transformar tanto la columna de los tratamientos como la de los bloques en un factor para poder realizar los cálculos posteriores adecuadamente.

```
propuesto9$Sexo <- factor(propuesto9$Sexo)
propuesto9$Sexo
```

```
## [1] Hombres Hombres Mujeres Mujeres Hombres Hombres Mujeres Mujeres Hombres
## [10] Hombres Mujeres Mujeres Hombres Hombres Mujeres Mujeres Hombres Hombres
## [19] Mujeres Mujeres Hombres Hombres Mujeres Mujeres
## Levels: Hombres Mujeres
```

```
propuesto9$Tabaco <- factor(propuesto9$Tabaco)
propuesto9$Tabaco

## [1] Rubio Negro Rubio Negro Rubio Negro Rubio Negro Rubio Negro Rubio Negro Rubio Negro Rubio Negro
## [13] Rubio Negro Rubio Negro Rubio Negro Rubio Negro Rubio Negro Rubio Negro Rubio Negro Rubio Negro
## Levels: Negro Rubio
```

```
propuesto9$Edades <- factor(propuesto9$Edades)
propuesto9$Edades

## [1] Niños Niños Niños Niños Niños Niños Niños Niños Jóvenes
## [10] Jóvenes Jóvenes Jóvenes Jóvenes Jóvenes Jóvenes Jóvenes Adultos Adultos
## [19] Adultos Adultos Adultos Adultos Adultos Adultos
## Levels: Adultos Jóvenes Niños
```

2. Contrastar la hipótesis nula de no interacción entre los factores. Adecuar el modelo al resultado de las interacciones y contrastar los efectos principales.

Primero vamos a obtener la tabla ANOVA

```
modp9 <- aov(Nivel_nicotina ~ Sexo + Tabaco + Edades +
             Sexo * Tabaco + Sexo * Edades +
             Tabaco * Edades + Sexo * Tabaco * Edades,
             data = propuesto9)
modp9

## Call:
## aov(formula = Nivel_nicotina ~ Sexo + Tabaco + Edades + Sexo *
## Tabaco + Sexo * Edades + Tabaco * Edades + Sexo * Tabaco *
## Edades, data = propuesto9)
##
## Terms:
##              Sexo      Tabaco      Edades  Sexo:Tabaco  Sexo:Edades
## Sum of Squares   4565.0     210.0 306192.3     38160.4    227044.1
## Deg. of Freedom      1        1        2           1          2
##              Tabaco:Edades  Sexo:Tabaco:Edades  Residuals
## Sum of Squares      41848.1              5.2  448698.5
## Deg. of Freedom        2              2        12
##
## Residual standard error: 193.3689
## Estimated effects may be unbalanced

summary(modp9)

##              Df Sum Sq Mean Sq F value Pr(>F)
```

```
## Sexo          1    4565    4565    0.122 0.7328
## Tabaco        1     210     210    0.006 0.9415
## Edades        2  306192  153096    4.094 0.0441 *
## Sexo:Tabaco   1   38160   38160    1.021 0.3323
## Sexo:Edades   2  227044  113522    3.036 0.0857 .
## Tabaco:Edades  2   41848   20924    0.560 0.5857
## Sexo:Tabaco:Edades 2      5      3    0.000 0.9999
## Residuals     12 448698   37392
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El único efecto significativo son las distintas edades. Al no ser significativa ninguna de las interacciones, realizamos de nuevo el ANOVA sufriendo la interacción entre los tres factores

```
modp91 <- aov(Nivel_nicotina ~ Sexo + Tabaco +
              Edades+ Sexo * Tabaco + Sexo *
              Edades + Tabaco * Edades, data = propuesto9)
modp91

## Call:
## aov(formula = Nivel_nicotina ~ Sexo + Tabaco + Edades + Sexo *
## Tabaco + Sexo * Edades + Tabaco * Edades, data = propuesto9)
##
## Terms:
##              Sexo      Tabaco      Edades  Sexo:Tabaco  Sexo:Edades
## Sum of Squares   4565.0     210.0  306192.3     38160.4    227044.1
## Deg. of Freedom      1         1        2           1          2
##              Tabaco:Edades Residuals
## Sum of Squares      41848.1  448703.7
## Deg. of Freedom        2        14
##
## Residual standard error: 179.0259
## Estimated effects may be unbalanced
```

De nuevo el único efecto significativo es la Edad. Realizamos de nuevo el ANOVA suprimiendo la interacción de orden 2 Sexo*Tabaco

```
modp92 <- aov(Nivel_nicotina ~ Sexo + Tabaco +
              Edades+ Sexo * Edades + Tabaco * Edades,
              data = propuesto9)
modp92

## Call:
## aov(formula = Nivel_nicotina ~ Sexo + Tabaco + Edades + Sexo *
## Edades + Tabaco * Edades, data = propuesto9)
```

```
##
## Terms:
##              Sexo    Tabaco    Edades  Sexo:Edades  Tabaco:Edades  Residuals
## Sum of Squares   4565.0     210.0 306192.3    227044.1      41848.1  486864.1
## Deg. of Freedom      1         1      2          2          2        15
##
## Residual standard error: 180.16
## Estimated effects may be unbalanced
```

De nuevo el mismo resultado, suprimimos la interacción Tabaco*Edades

```
modp93 <- aov(Nivel_nicotina ~ Sexo + Tabaco +
              Edades+ Sexo * Edades, data = propuesto9)
modp93

## Call:
## aov(formula = Nivel_nicotina ~ Sexo + Tabaco + Edades + Sexo *
##      Edades, data = propuesto9)
##
## Terms:
##              Sexo    Tabaco    Edades  Sexo:Edades  Residuals
## Sum of Squares   4565.0     210.0 306192.3    227044.1  528712.2
## Deg. of Freedom      1         1      2          2        17
##
## Residual standard error: 176.354
## Estimated effects may be unbalanced
```

Son significativos el efecto de la Edad y el efecto de la interacción del Sexo por la Edad

3. ¿Hay diferencias significativas en el nivel de nicotina en las distintas edades?
- ¿En qué edad el nivel de nicotina es mayor?
- ¿Hay diferencias significativas en el nivel de nicotina en las distintas edades?

```
summary(modp9)

##              Df Sum Sq Mean Sq F value Pr(>F)
## Sexo          1   4565    4565    0.122  0.7328
## Tabaco         1    210     210    0.006  0.9415
## Edades         2 306192  153096    4.094  0.0441 *
## Sexo:Tabaco    1  38160   38160    1.021  0.3323
## Sexo:Edades    2 227044  113522    3.036  0.0857 .
## Tabaco:Edades  2  41848   20924    0.560  0.5857
## Sexo:Tabaco:Edades 2      5      3    0.000  0.9999
## Residuals     12 448698   37392
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Si hay diferencias significativas del nivel de nicotina en las distintas edades (P-valor es 0.0441)

¿En qué edad el nivel de nicotina es mayor?

Realizamos un contraste de comparaciones múltiples. Por ejemplo LSD

```
library(agricolae)
LSD.test(modp93,"Edades", p.adj="bonferroni", console=TRUE)

##
## Study: modp93 ~ "Edades"
##
## LSD t Test for Nivel_nicotina
## P value adjustment method: bonferroni
##
## Mean Square Error: 31100.72
##
## Edades, means and individual ( 95 %) CI
##
##      Nivel_nicotina      std r      LCL      UCL Min Max
## Adultos      451.375 194.80755 8 319.82686 582.9231 198 701
## Jóvenes      409.250 251.74349 8 277.70186 540.7981 112 873
## Niños        193.500  85.57202 8  61.95186 325.0481 110 360
##
## Alpha: 0.05 ; DF Error: 17
## Critical Value of t: 2.654996
##
## Minimum Significant Difference: 234.1095
##
## Treatments with the same letter are not significantly different.
##
##      Nivel_nicotina groups
## Adultos      451.375      a
## Jóvenes      409.250     ab
## Niños        193.500      b
```

El nivel de nicotina es mayor en los Adultos (451.375 miligramos por mililitro).

4. ¿El tipo de tabaco es un factor determinante en el nivel de nicotina?

```
summary(modp9)
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## Sexo          1    4565     4565   0.122 0.7328
## Tabaco        1     210      210   0.006 0.9415
## Edades        2  306192  153096   4.094 0.0441 *
## Sexo:Tabaco    1   38160   38160   1.021 0.3323
## Sexo:Edades    2  227044  113522   3.036 0.0857 .
## Tabaco:Edades  2   41848   20924   0.560 0.5857
## Sexo:Tabaco:Edades 2      5      3   0.000 0.9999
## Residuals     12 448698   37392
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

El tipo de tabaco no es determinante en el nivel de nicotina (P-valor = 0.9415)

5. Comparar el nivel medio de nicotina entre las mujeres y los hombres. ¿Se detectan diferencias significativas?

No hay diferencias significativas en el nivel de nicotina entre las mujeres y los hombres (P-valor = 0.7328)

Apliquemos un contraste de comparaciones múltiples para ver donde es mayor

```
SNK.test(modp93,"Sexo", console=TRUE)
```

```
##
## Study: modp93 ~ "Sexo"
##
## Student Newman Keuls Test
## for Nivel_nicotina
##
## Mean Square Error: 31100.72
##
## Sexo, means
##
##      Nivel_nicotina      std  r Min Max
## Hombres      337.5833 198.3548 12 110 701
## Mujeres      365.1667 239.1971 12 123 873
##
## Alpha: 0.05 ; DF Error: 17
##
## Critical Range
##      2
## 151.8987
```

```
##
## Means with the same letter are not significantly different.
##
##      Nivel_nicotina groups
## Mujeres      365.1667      a
## Hombres      337.5833      a
```

Las mujeres y los hombres forman un único grupo donde no se aprecian diferencias significativas. Es mayor el nivel de nicotina entre las mujeres (365.1667 miligramos por mililitro)