

Práctica 4: Experimentos con uno y dos factores

1ª parte

Modelos Lineales

Grados en Estadística y Matemáticas

- 1 Experimentos con un factor
- 2 Diseño Completamente Aleatorizado
- 3 Diseño de Bloques al Azar
- 4 Diseño de Cuadrados Latinos

Modelos de Diseño de Experimentos

Recordemos que en los modelos de Diseño de Experimentos aparecen dos elementos comunes:

- La **variable respuesta** (Y): la característica que puede ser observada, medida y analizada. Supondremos que es cuantitativa continua y con distribución normal.
- El **factor o factores**: un factor es una variable cuyo posible efecto sobre la respuesta se quiere estudiar; los niveles de un factor son los valores que puede tomar el factor. El factor será tratado como una variable cualitativa.

El objetivo del análisis estadístico que vamos a realizar a través de estos modelos es determinar la posible influencia del factor o factores en la variable respuesta Y . Sólo consideraremos modelos de efectos fijos.

Diseño Completamente Aleatorizado

Toma de datos

La asignación de los k niveles del factor o factores, que pueden influir en la variable respuesta, a cada una de las unidades experimentales consideradas se hace totalmente al azar.

Como resultado obtendremos k muestras independientes, no necesariamente del mismo tamaño:

- Si todas las muestras tienen el mismo tamaño se dice que el diseño es **balanceado**
- Si las muestras tienen distinto tamaño se dice que el diseño es **no balanceado**

Ejercicio 1

Una operación de llenado tiene tres máquinas idénticas que se ajustan para vaciar una cantidad específica de un producto en recipientes de igual tamaño. Con el propósito de identificar diferencias entre las cantidades (en litros) vaciadas por cada máquina, se toman muestras aleatorias, en forma periódica, de cada una. Para un periodo particular se observaron los datos que aparecen en el fichero `maquina.dat`.

Se trata de determinar si existen diferencias estadísticamente significativas en las cantidades promedio vaciadas por las tres máquinas.

Modelo matemático

El modelo matemático que se utiliza para analizar estos datos es el siguiente:

$$Y_{ij} = \mu + \tau_i + \varepsilon_{ij} \quad i = 1, \dots, k \quad j = 1, \dots, n_i \quad N = \sum_{i=1}^k n_i$$

donde

- Y_{ij} : representa el valor de Y para la observación j del nivel i del factor
- μ : es la media poblacional de Y si el factor no estuviera presente.
- τ_i : efecto que sobre la media poblacional de Y tiene el nivel i del factor.
- ε_{ij} : componente aleatoria en la medición de Y_{ij} (no observable y con media 0).

Si denotamos $\mathcal{E} = (\mathcal{E}_{11}, \dots, \mathcal{E}_{kn_k})^t$, para aplicar toda la metodología vista en clase asumiremos:

$$\mathcal{E} \sim \mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{I}_N)$$

Tabla ANOVA

El contraste de hipótesis más importante para este modelo es el siguiente:

$$H_0 : \tau_1 = \tau_2 = \cdots = \tau_k$$

$$H_1 : \tau_i \neq \tau_j \text{ para algún par de niveles}(i, j)$$

Para resolverlo se utiliza la tabla Análisis de la Varianza vista en clase, donde se comparan las fuentes de variabilidad del experimento:

F.V.	g.l.	Suma de cuadrados	M.C.	F
Debida a H_0	$k - 1$	$Q_1 = \sum_{i=1}^k n_i \bar{Y}_i^2 - N \bar{Y}_{..}^2$	$\frac{Q_1}{k - 1}$	$F = \frac{(N - k)Q_1}{(k - 1)Q_0}$
Error	$N - k$	$Q_0 = \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}^2 - \sum_{i=1}^k n_i \bar{Y}_i^2$	$\frac{Q_0}{N - k}$	

Comparaciones múltiples

Si aceptamos H_1 , habremos detectado diferencias entre los tratamientos en cuanto al comportamiento de la variable respuesta Y . Para saber en qué tratamientos residen dichas diferencias hay que aplicar un procedimiento de comparaciones múltiples o comparaciones dos a dos donde se realizan, de modo simultáneo, todos los contrastes posibles.

$$H_0 : \tau_i = \tau_j$$

$$H_1 : \tau_i \neq \tau_j$$

Para resolverlos hay gran cantidad de métodos: LSD, Scheffé, **Tukey**,...

Diseño de Bloques al Azar

Toma de datos

Supongamos que estamos estudiando la posible influencia que un factor (con k niveles) puede ejercer sobre una variable cuantitativa continua Y .

- Además del factor, existe **otra fuente de variabilidad en el experimento** con influencia sobre el mismo.
- De acuerdo a esta fuente de variabilidad, se divide a las unidades experimentales en s grupos denominados **bloques**.
- Dentro de cada bloque se elige, de forma aleatoria, una unidad experimental para cada uno de los niveles del factor.

Las observaciones aparecen agrupadas por bloques, obteniendo como resultado k muestras relacionadas, una para cada nivel del factor.

Ejercicio 2

Un fabricante desea emplear un nuevo tipo de aleación en la producción de filtros y quiere determinar la velocidad de agitación adecuada ya que piensa que ésta puede influir en la resistencia que alcanza el producto final.

En el local de la fabrica se cuenta con 4 hornos, cada uno de los cuales tiene sus propias características de operación, lo que los convierte en posibles fuentes de variabilidad.

Cada horno puede operarse a las tres velocidades propuestas por un experto (10 rpm, 15 rpm y 20 rpm). Los datos de resistencia (en Kg/cm^2) se encuentran en el archivo `filtros.dat`

Ejercicio 2 (continuación)

Los datos podrían disponerse de este modo

Velocidad de Agitación	H1	H2	H3	H4
10 rpm	8.3	9.4	8.5	7.5
15 rpm	14.4	17.8	15.8	13.4
20 rpm	11.1	13.2	11.8	10.1

Las unidades experimentales (aleación para fabricar un filtro) se han agrupado en **bloques: los hornos donde se produce la aleación**.

De esta forma, para cada **tratamiento (velocidad de agitación)** tenemos 4 datos de resistencia, uno correspondiente a cada horno. Son tres muestras relacionadas.

Modelo matemático

El modelo matemático que se utiliza para analizar estos datos es el siguiente:

$$Y_{ij} = \mu + \tau_i + \delta_j + \varepsilon_{ij} \quad i = 1, \dots, r \quad j = 1, \dots, s$$

donde

- Y_{ij} : representa el valor de Y para el nivel i y el bloque j .
- μ : es la media poblacional de Y sin ninguna fuente de variabilidad presente.
- τ_i : efecto que sobre la media poblacional de Y tiene el nivel i del factor.
- δ_j : efecto sobre la media poblacional de Y por pertenecer al bloque j .
- ε_{ij} : componente aleatoria en la medición de Y_{ij} (no observable y con media 0).

Si denotamos $\mathcal{E} = (\mathcal{E}_{11}, \dots, \mathcal{E}_{rs})^t$, para aplicar toda la metodología vista en clase asumiremos:

$$\mathcal{E} \sim \mathcal{N}_{rs}(\mathbf{0}, \sigma^2 \mathbf{I}_{rs})$$

Tabla ANOVA

El contraste de hipótesis que se plantea es el siguiente:

$$H_0 : \tau_1 = \tau_2 = \dots = \tau_r$$

H_1 : los efectos de los tratamientos no son iguales

Para resolverlo se utiliza la tabla Análisis de la Varianza vista en clase, donde se comparan las fuentes de variabilidad del experimento:

F.V.	g.l.	Suma de cuadrados	M.C.	F
Tr. (τ)	$r - 1$	$Q_1 = \sum_{i=1}^r s\bar{y}_i^2 - ks\bar{y}_{..}^2$	$\frac{Q_1}{r - 1}$	$F_T = \frac{(s - 1)Q_1}{Q_0}$
Bl. (δ)	$s - 1$	$Q'_1 = \sum_{j=1}^s r\bar{y}_{.j}^2 - rs\bar{y}_{..}^2$	$\frac{Q'_1}{s - 1}$	$F_B = \frac{(r - 1)Q'_1}{Q_0}$
Error	$(r - 1)(s - 1)$	$Q_0 = \sum_{i=1}^r \sum_{j=1}^s y_{ij}^2 - \sum_{i=1}^r s\bar{y}_i^2 - \sum_{j=1}^s r\bar{y}_{.j}^2 + rs\bar{y}_{..}^2$	$\frac{Q_0}{(r - 1)(s - 1)}$	

Comparaciones múltiples

Si aceptamos H_1 , habremos detectado diferencias entre los tratamientos en cuanto al comportamiento de la variable Y . Para saber en qué tratamientos residen dichas diferencias hay que aplicar un procedimiento de comparaciones múltiples o comparaciones dos a dos, donde se realizan, de modo simultáneo, todos los contrastes posibles

$$H_0 : \tau_i = \tau_j$$

$$H_1 : \tau_i \neq \tau_j$$

Para resolverlos hay gran cantidad de métodos: LSD, Scheffé, **Tukey**,...

Diseño de Cuadrados Latinos

Toma de datos

Supongamos que estamos estudiando la posible influencia que un factor (con k niveles) puede ejercer sobre una variable cuantitativa continua Y .

- Además del factor, existen **otras dos fuentes de variabilidad en el experimento** con influencia sobre el mismo.
- Si dividimos la población en bloques de acuerdo a ambas fuentes de variabilidad, el número de observaciones necesarias para abarcar todas las combinaciones posibles podría ser demasiado grande.
- Cada fuente de variabilidad se divide en k niveles con los cuales se forma un cuadrado $k \times k$.
- Cada tratamiento del factor aparece una sola vez en cada fila y cada columna.

Ejercicio 3

Un investigador quiere evaluar la productividad de cuatro variedades de aguacate, Reed, Hass, Pinkerton y Bacon. 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, estableció 4 zonas de acuerdo a la pendiente y otras 4 en base a la disponibilidad de Nitrógeno, fijando así 16 parcelas. Posteriormente realizó un plantación de cada variedad en cada una de las pendientes y para cada una de las diferentes concentraciones de Nitrógeno.

Los datos correspondientes a la producción (en kg/parcela) se encuentran en el archivo `aguacate.txt`.

Ejercicio 3 (continuación)

Veamos como se disponen los datos en el cuadrado latino que ha sido elegido por el experimentador.

		Pendiente			
		1	2	3	4
Nitrógeno	1	Reed 785	Hass 730	Pinkerton 700	Bacon 595
	2	Hass 855	Bacon 775	Reed 760	Pinkerton 710
	3	Pinkerton 950	Reed 885	Bacon 795	Hass 780
	4	Bacon 945	Pinkerton 950	Hass 880	Reed 835

Modelo matemático

El modelo matemático que se utiliza para analizar estos datos es el siguiente:

$$Y_{ijh} = \mu + \tau_i + \delta_j + \gamma_h + \varepsilon_{ijh}, \quad i, j = 1, \dots, k$$

h varía según marca el cuadrado latino elegido. En este modelo:

- Y_{ijh} : representa el valor de Y en el nivel i , fila j , columna h .
- μ : es la media poblacional de Y sin ninguna fuente de variabilidad presente.
- τ_i : efecto que sobre la media poblacional de Y tiene el nivel i del factor.
- δ_j : efecto sobre la media poblacional de Y por estar en la fila j .
- γ_h : efecto sobre la media poblacional de Y por estar en la columna h .
- ε_{ijh} : componente aleatoria en la medición de Y_{ijh} (no observable y con media 0).

Si denotamos $\mathcal{E} = (\mathcal{E}_{11}, \dots, \mathcal{E}_{k^2})^t$, para aplicar toda la metodología vista en clase asumiremos:

$$\mathcal{E} \sim \mathcal{N}_{k^2}(\mathbf{0}, \sigma^2 \mathbf{I}_{k^2})$$

Tabla ANOVA

El contraste de hipótesis que se plantea es el siguiente:

$$H_0 : \tau_1 = \tau_2 = \dots = \tau_k$$

H_1 : los efectos de los tratamientos no son iguales

Tabla de Análisis de la Varianza

F.V.	g.l.	Suma de cuadrados	M.C.	F
Tr. (τ)	$k - 1$	$Q_1 = \sum_{i=1}^k k\bar{Y}_{i..}^2 - k^2\bar{Y}^2_{...}$	$\frac{Q_1}{k - 1}$	$F_T = \frac{(k - 2)Q_1}{Q_0}$
Fil. (δ)	$k - 1$	$Q'_1 = \sum_{j=1}^k k\bar{Y}^2_{.j} - k^2\bar{Y}^2_{...}$	$\frac{Q'_1}{k - 1}$	$F_A = \frac{(k - 2)Q'_1}{Q_0}$
Col. (γ)	$k - 1$	$Q''_1 = \sum_{j=1}^k k\bar{Y}^2_{.j} - k^2\bar{Y}^2_{...}$	$\frac{Q''_1}{k - 1}$	$F_B = \frac{(k - 2)Q''_1}{Q_0}$
Error	$(k - 2)(k - 1)$	$Q_0 = \sum_{i,j,h=1}^k Y^2_{ijh} - \sum_{i=1}^k k\bar{Y}^2_{i..} - \sum_{j=1}^s k\bar{Y}^2_{.j} - \sum_{j=1}^s k\bar{Y}^2_{.j} + 2k^2\bar{Y}^2_{...}$	$\frac{Q_0}{(k - 2)(k - 1)}$	

Hipótesis teóricas

Para utilizar este modelo hemos de asumir **normalidad**, **homocedasticidad** y **no autocorrelación**, que se comprobarán a través del análisis de los residuos.

Comparaciones múltiples

Si aceptamos H_1 , habremos detectado diferencias entre los tratamientos en cuanto al comportamiento de la variable Y . Para saber en qué tratamientos residen dichas diferencias hay que aplicar un procedimiento de comparaciones múltiples o comparaciones dos a dos, es decir realizaremos de modo simultáneo, todos los contrastes posibles

$$H_0 : \tau_i = \tau_j$$

$$H_1 : \tau_i \neq \tau_j$$

Ejercicio 4

Desde el incremento en los precios de la gasolina, se han desarrollado varios dispositivos que se colocan en los carburadores de los automóviles. Una empresa selecciona tres de estos dispositivos para someterlos a prueba. La empresa desea compararlos con los carburadores estándar, con el propósito de determinar si existe un incremento apreciable de kilómetros por litro de gasolina con el uso de estos dispositivos.

La compañía selecciona cinco automóviles para el experimento. Para controlar la variación, se planea utilizar el mismo conductor para todo el experimento. Se observan los datos del archivo carbu.dat.

Estudia las posibles diferencias de consumo entre los carburadores.

Ejercicio 5

El conjunto de datos `iris` proporciona las medidas en centímetros de las variable *longitud del sépalo*, *anchura del sépalo*, *longitud del pétalo* y *anchura del pétalo*, respectivamente, para 50 flores de tres especies of iris: setosa, versicolor y virginica.

Se trata de detectar si hay diferencias significativas entre las tres especies en cada una de las cuatro medidas. Los datos se encuentran en el archivo `iris.dat`

Ejercicio 6

Una empresa de bebidas está interesada en evaluar el impacto de diferentes métodos de carbonatación en la efervescencia de sus refrescos. Se han identificado cuatro métodos de carbonatación que se aplicarán a un lote de refrescos:

- A: Carbonatación por inyección de CO_2
- B: Carbonatación natural (fermentación)
- C: Carbonatación por absorción
- D: Carbonatación por mezcla de gases

El objetivo del estudio es determinar si hay diferencias significativas en la efervescencia de los refrescos, medida a través de la cantidad de dióxido de carbono (CO_2) disuelto, al aplicar los diferentes métodos de carbonatación.

Para controlar la variabilidad entre los diferentes lotes de refrescos y el tiempo de almacenamiento, cada método de carbonatación se aplicará una vez para cada lote de refresco y en cada tiempo de almacenamiento. Los datos se encuentran en el archivo `carbonatacion.txt`. Analiza estos datos para determinar las posibles diferencias en contenido de dióxido de carbono de refrescos carbonatados con los 4 métodos.