

Diseño y Análisis de Experimentos

Diseño Completamente Aleatorizado

Diseño Completamente Aleatorizado

Caso: Pozos Geotérmicos

- ♦ Una compañía de explotación geotérmica desea ensayar dos nuevos diseños de trépanos para perforación:
 - » Uno en base a una nueva aleación.
 - » Uno con una nueva geometría de los elementos de corte de la roca.
- ♦ La variable respuesta es la duración de las herramientas en metros perforados de pozo hasta el desgaste, determinado en base al ralentamiento del avance.
- ♦ Las unidades experimentales son los tramos de pozo elegidos para ensayar cada trepano:
 - » Estos tramos son asignados al azar a los tres tipos de herramientas ensayados: la actual y las dos nuevas.

Diseño Completamente Aleatorizado

Caso: Pozos Geotérmicos

Duración [m]	Actual	Nueva Aleación	Nueva Geometría
	104	115	180
	128	108	146
	140	128	244
	156	180	155
	145	104	197
	71	164	145
	70	147	110
Prom.	116.3	135.1	168.1
s	35.23	29.29	43.41

- ◆ Modelo:

$$\tilde{y}_{ia} = \mu_i + \varepsilon_{ia} = \mu + \alpha_i + \tilde{\varepsilon}_{ia}$$

$$i = 1 \dots p \quad a = 1 \dots n_i \quad N = \sum n_i$$

- ◆ Condición de unicidad de efectos:

$$\sum_{i=1}^p \alpha_i = 0$$

- ◆ Supuestos:

$$S1) E\tilde{\varepsilon}_{ia} = 0$$

$$S2) V\tilde{\varepsilon}_{ia} = \sigma^2 \quad (\text{homocedasticidad})$$

$$S3) \tilde{\varepsilon}_{ia} \text{ independientes}$$

$$S4) \tilde{\varepsilon}_{ia} \text{ Normales}$$

Diseño Completamente Aleatorizado

Factores Fijos vs Aleatorios

- ◆ Factores fijos:
 - » Los p niveles en el experimento son todos los que existen, o bien son elegidos por un procedimiento sistemático.
 - » Ejemplo: Cantidad de fertilizante: 10, 20, 30, 50.
 - » La inferencia se aplica a los niveles seleccionados. Si el factor es continuo puede asumirse continuidad en el efecto, pero es un supuesto.

- ◆ Factores aleatorios:
 - » Los p niveles en el experimento se eligen entre un conjunto mayor de P niveles posibles al azar.
 - » La inferencia se aplica a todos los niveles del factor.

Diseño Completamente Aleatorizado

Inferencia

- ♦ El objetivo es ensayar:

$$H_0) \mu_1 = \mu_2 = \dots = \mu_p \quad \text{o} \quad H_0) \alpha_1 = \alpha_2 = \dots = \alpha_p = 0$$

- ♦ Evaluar hipótesis múltiples es complicado, pero R.A. FISHER encontró una hipótesis equivalente:

$$H_0) \sigma_\alpha^2 = 0$$

- ♦ Utilizando la varianza de las medias, de ahí el nombre: Análisis de la Varianza.
 - » La varianza de las medias se define, según la convención de W. COCHRAN:

$$\sigma_\alpha^2 = \frac{1}{p-1} \sum \alpha_i^2$$

- » Cuando el factor es fijo no es una verdadera varianza, sino un instrumento ingenioso para convertir la hipótesis múltiple en simple.

Diseño Completamente Aleatorizado

Inferencia

- ♦ En busca de una estimación de: $\sigma_\alpha^2 = \frac{1}{p-1} \sum \alpha_i^2$
- ♦ Analizamos su análogo empírico: $s_{\bar{y}}^2 = \frac{1}{p-1} \sum (\bar{y}_i - \bar{y}_\bullet)^2$
- ♦ Sin embargo se demuestra que, tanto para factor fijo como aleatorio su esperanza es:

$$\mathcal{E}s_{\bar{y}}^2 = \sigma_\alpha^2 + \frac{1}{n} \sigma^2$$

- » La varianza de los promedios no estima solamente a la varianza de las medias, sino que se entromete el ruido experimental.
- ♦ Si conociéramos σ^2 podríamos inferir sobre σ_α^2 mediante:

$$\text{Bajo } H_0 : \frac{(p-1)s_{\bar{y}}^2}{\frac{1}{n}\sigma^2} = \chi_{p-1}^2$$

- » Pero cada experimento tiene un nivel de ruido propio desconocido.

Diseño Completamente Aleatorizado

Inferencia

- ♦ Es necesario conseguir una estimación independiente del ruido.
- ♦ Como dentro de cada grupo hay replicas se puede estimar el ruido en cada grupo y amalgamar esos estimadores:

$$s^2 = \frac{1}{\nu} \sum_{i=1}^p \nu_i s_i^2$$

$$\nu_i = n_i - 1$$

$$\nu = \sum \nu_i$$

$$\text{con } s_i^2 = \frac{1}{\nu_i} \sum_{a=1}^{n_i} (y_{ia} - \bar{y}_i)^2$$

- ♦ Evidentemente:

$$\mathbb{E}s^2 = \frac{1}{\nu} \sum_{i=1}^p \nu_i \mathbb{E}s_i^2 = \frac{1}{\nu} \sum_{i=1}^p \nu_i \sigma^2 = \sigma^2$$

Diseño Completamente Aleatorizado

Inferencia

- Entonces por comparación de ambos estimadores podemos inferir sobre σ_α^2

$$ns_{\bar{y}}^2 = \frac{n}{p-1} \sum (\bar{y}_i - \bar{y}.)^2 \qquad \mathbb{E}ns_{\bar{y}}^2 = \sigma^2 + n\sigma_\alpha^2$$

$$s^2 = \frac{1}{v} \sum v_i s_i^2 \qquad \mathbb{E}s^2 = \sigma^2$$

- » Si $ns_{\bar{y}}^2 \gg s^2$ inferimos que $\sigma_\alpha^2 > 0$ y rechazamos H_0
- » Como ambos estimadores son independientes por serlo \bar{y}_i de s_i^2 :
 - » $\frac{ns_{\bar{y}}^2}{s^2} : F_{p-1, N-p}$ que permite inferir.
 - » También se puede estimar σ_α^2 : $s_\alpha^2 = s_{\bar{y}}^2 - \frac{1}{n} s^2$
 - » Como los estimadores son independientes puede dar negativa. En este caso puede aceptarse H_0 , aunque si el cociente queda en la extremidad improbable de la distribución F , seguramente se debe a error de especificación del modelo.

Diseño Completamente Aleatorizado

Inferencia

- ♦ La inferencia se hace ordenadamente mediante la tabla ANOVA:

<i>Efecto</i>	<i>Estimador</i>	Q	ν	$CM = Q / \nu$	$\mathbb{E}CM$
α_i	$\bar{y}_i - \bar{y}_\bullet$	$C_\alpha - C$	$p - 1$		$\sigma^2 + n\sigma_\alpha^2$
ε_{ia}	$y_{ia} - \bar{y}_i$	$C_\varepsilon - C_\alpha$	$N - p$		σ^2
Total	$y_{ia} - \bar{y}_\bullet$	$C_\varepsilon - C$	$N - 1$		

$$C = N \bar{y}_\bullet^2$$

$$C_\alpha = \sum_{i=1}^p n_i \bar{y}_i^2$$

$$C_\varepsilon = \sum_{i=1}^p \sum_{a=1}^{n_i} y_{ia}^2$$

Diseño Completamente Aleatorizado

Caso: Pozos Geotérmicos

<i>Efecto</i>	Q	ν	$CM = Q / \nu$	$\mathbb{E}CM$
α_i	9645	2	4823	$\sigma^2 + n\sigma_\alpha^2$
ε_{ia}	23901	18	1328	σ^2
Total	33546	20		

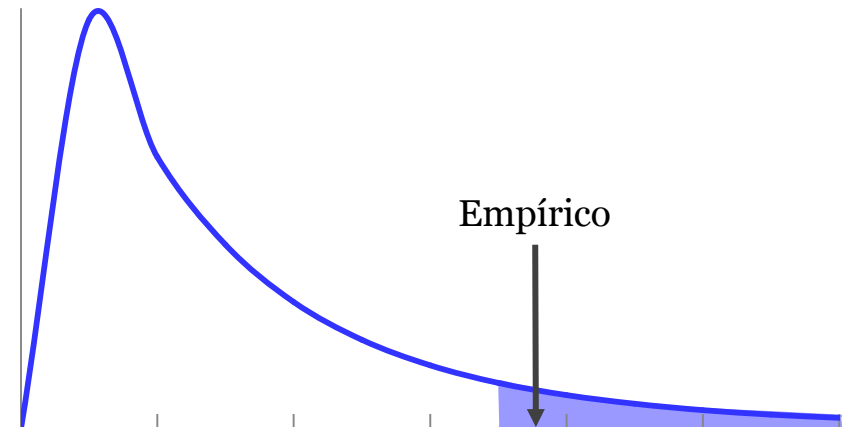
$$C = 410.760$$

$$C_\alpha = 420.406$$

$$C_\varepsilon = 444.307$$

$$F = \frac{4823}{1328} = 3,632 \quad \Rightarrow \quad \alpha^* = 4,7\%$$

$$F_{1;18;95\%} = 3,55 \quad \Rightarrow \quad \text{se rechaza } H_0$$



Comparaciones Múltiples

Comparaciones Múltiples

Caso: Pozos Geotérmicos

- ♦ El análisis de varianza da un resultado general: la hipótesis de FISHER.
- ♦ Cuando los factores son fijos interesa hacer inferencias específicas:
- ♦ Por ejemplo:
 - » Saber si los nuevos trépanos son mejores: $\mathcal{G}' = \frac{\mu_2 + \mu_3}{2} - \mu_1$
 - » Saber cual de los trépanos nuevos es mejor: $\mathcal{G}'' = \mu_2$ vs μ_3
- ♦ Estas comparaciones no pueden ensayarse por los métodos tradicionales por dos razones:
 - » El nivel de significación conjunto de múltiples ensayos es mayor que el de cada ensayo individual.
 - » Muchas veces las comparaciones son sugeridas por los resultados (a posteriori) y esto altera la distribución del estadístico.

Comparaciones Múltiples

- ♦ Las comparaciones se pueden plantear:
 - » A priori: antes de ver los datos.
 - » A posteriori: sugeridas por los datos.
- ♦ Los siguientes métodos cubren los casos mas frecuentes:

Comparaciones	a Priori	STUDENT
		BONFERRONI
		DUNNETT
	a Posteriori	TUKEY
		SCHEFFÉ

Comparaciones Múltiples

Método de Student

- ◆ Definición de comparación o contraste:

$$\mathcal{T} = \sum_{i=1}^p c_i \mu_i \quad \text{es una comparación sii} \quad \sum c_i = 0$$

- ◆ Se desea ensayar la hipótesis: $H_0) \mathcal{T} = \mathcal{T}_0$

» Donde habitualmente: $\mathcal{T}_0 = 0$

- ◆ Estimador:

$$\hat{\mathcal{T}} = \sum_{i=1}^p c_i \bar{y}_i$$

$$\mathbb{E} \hat{\mathcal{T}} = \sum_{i=1}^p c_i \mathbb{E} \bar{y}_i = \sum_{i=1}^p c_i \mu_i = \mathcal{T}$$

$$\mathcal{V} \hat{\mathcal{T}} = \sum_{i=1}^p c_i^2 \mathcal{V} \bar{y}_i = \sum_{i=1}^p \frac{c_i^2}{n_i} \sigma^2$$

Comparaciones Múltiples

Método de Student

- ♦ Estimando esta varianza mediante:

$$s_{\mathcal{C}}^2 = \sum_{i=1}^p \frac{c_i^2}{n_i} s^2$$

- ♦ Bajo el supuesto de normalidad $\mathcal{C} : \mathcal{N}$ entonces:

$$t_{\nu} = \frac{\hat{\mathcal{C}} - \mathcal{C}_0}{\sqrt{s^2 \sum \frac{c_i^2}{n_i}}}$$

- ♦ es t-Student con ν : grados de libertad de s^2

- ♦ Cuando se desea ensayar más de una comparación surge un problema: El resultado de una comparación puede influir en otros, violando el nivel de riesgo tolerado. Es decir, la segunda comparación no es verdaderamente a priori.

Comparaciones Múltiples

Método de Student

- ♦ Def: \mathcal{C}' y \mathcal{C}'' ortogonales $\Leftrightarrow \sum_{i=1}^p c'_i c''_i = 0$
- ♦ Teorema: $\hat{\mathcal{C}}'$ y $\hat{\mathcal{C}}''$ independientes $\Leftrightarrow \mathcal{C}'$ y \mathcal{C}'' ortogonales
- ♦ Entonces se pueden ensayar varias comparaciones por el método de Student siempre que sean ortogonales dos a dos.
- ♦ Existen $p-1$ comparaciones ortogonales.
- ♦ Para resolver el problema de la amplificación del riesgo se ajusta el nivel de significación:

$$1 - \alpha = (1 - \alpha_c)^{p-1}$$

Comparaciones Múltiples

Método de Student – Caso: Pozos Geotérmicos

- Las comparaciones son ortogonales:

	1	2	3	
\mathcal{C}'	-1	$\frac{1}{2}$	$\frac{1}{2}$	
\mathcal{C}''	0	-1	1	
$\mathcal{C}' \bullet \mathcal{C}''$	0	$-\frac{1}{2}$	$\frac{1}{2}$	= 0

$$\alpha_c = 1 - (1 - \alpha)^{3-1} = 1 - (1 - 5\%)^2 = 2,5\%$$

$$H_0') \frac{\mu_2 + \mu_3}{2} - \mu_1 \leq 0$$

$$\hat{\mathcal{C}}' = \frac{135,1 + 168,1}{2} - 116,3 = 35,4$$

$$\sum \frac{c_i'^2}{n_i} = \frac{1}{n} \sum c_i'^2 = \frac{1}{7} \left[(-1)^2 + \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right] = \frac{1,5}{7}$$

$$t_{18} = \frac{35,4 - 0}{\sqrt{\frac{1,5 \cdot 1328}{7}}} = 2,096$$

$$\Rightarrow \alpha^* = 2,52\% \Rightarrow \text{Rech } H_0'$$

$$H_0'') \mu_3 - \mu_2 = 0$$

$$\hat{\mathcal{C}}'' = 168,1 - 135,1 = 33,4$$

$$\sum \frac{c_i''^2}{n_i} = \frac{1}{n} \sum c_i''^2 = \frac{1}{7} \left[(-1)^2 + 1^2 \right] = \frac{2}{7}$$

$$t_{18} = \frac{33 - 0}{\sqrt{\frac{2 \cdot 1328}{7}}} = 1,649$$

$$\Rightarrow \alpha^* = 11\% \Rightarrow \text{No Rech } H_0''$$

Comparaciones Múltiples

Método de Bonferroni

- ♦ R.A. FISHER emplea la desigualdad de BONFERRONI para ensayar comparaciones no ortogonales a priori mediante una modificación de α :

$$H_0) \mu_1 = \mu_2 = \dots = \mu_p \Rightarrow \forall j = 1 \dots r : H_0^j) \mathcal{C}^j = 0$$

$$P(\text{rech } H_0 / H_0) \leq \alpha$$

$$P(\text{rech } H_0^j / H_0^j) \leq \alpha_c$$

- ♦ Se desea ensayar todas estas comparaciones con el mismo nivel de significación que la hipótesis general:

$$P(\text{rech } H_0^1 \cup H_0^2 \cup \dots \cup H_0^r) \leq \alpha$$

- ♦ Por la desigualdad de BONFERRONI:

$$P(\text{rech } H_0^1 \cup H_0^2 \cup \dots \cup H_0^r) \leq P(\text{rech } H_0^1) + P(\text{rech } H_0^2) + \dots + P(\text{rech } H_0^r)$$

$$\Rightarrow \alpha^* \leq \sum \alpha_{c_j}^* \leq r \bar{\alpha}_c^*$$

- ♦ Esto sugiere hacer el ensayo de STUDENT con:

$$\alpha_c = \alpha / r$$

Comparaciones Múltiples

Método de Bonferroni – Caso: Pozos Geotérmicos

- ♦ Interesa comparar todos los pares de medias:

$$\alpha_c = \frac{\alpha}{3} = \frac{5\%}{3} = 1,67\%$$

$$H_0) \mu_1 = \mu_2$$

$$t_{18} = \frac{135,1 - 116,3}{\sqrt{\frac{2.1328}{7}}} = \frac{18,9}{19,48} = 0,968 \Rightarrow \alpha^* = 35\% \Rightarrow \text{No rech } H_0$$

$$H_0) \mu_1 = \mu_3$$

$$t_{18} = \frac{168,1 - 116,3}{\sqrt{\frac{2.1328}{7}}} = \frac{51,9}{19,48} = 2,662 \Rightarrow \alpha^* = 1,6\% < 1,67\% = \alpha_c \Rightarrow \text{Rech } H_0$$

$$H_0) \mu_2 = \mu_3$$

$$t_{18} = \frac{168,1 - 135,1}{\sqrt{\frac{2.1328}{7}}} = \frac{33}{19,48} = 1,694 \Rightarrow \alpha^* = 11\% > 1,67\% = \alpha_c \Rightarrow \text{No rech } H_0$$

Comparaciones Múltiples

Método de Dunnett

- ♦ Fue diseñado para comparar todos los tratamientos contra un testigo o control

$$H_0) \mu_i = \mu_1 \quad 1: \text{grupo testigo}$$

- ♦ Estas comparaciones no son ortogonales. Podrían ensayarse por el método de BONFERRONI, pero C. DUNNETT (1955) inventó un estadístico más potente:

$$D_{p,v} = \frac{\bar{y}_i - \bar{y}_1}{\sqrt{2s^2/n}}$$

- ♦ La distribución de D esta tabulada a simple y doble extremidad.

Comparaciones Múltiples

Método de Dunnett – Caso: Pozos Geotérmicos

$$H_0) \mu_1 \leq \mu_1$$

$$D_{3;18} = \frac{135,1 - 116,3}{\sqrt{\frac{2.1328}{7}}} = \frac{18,9}{19,48} = 0,968 < 2,4 = D_{3;18;95\%} \Rightarrow \text{No rech } H_0$$

$$H_0) \mu_3 \leq \mu_1$$

$$D_{3;18} = \frac{168,1 - 116,3}{\sqrt{\frac{2.1328}{7}}} = \frac{51,9}{19,48} = 2,662 > 2,4 = D_{3;18;95\%} \Rightarrow \text{Rech } H_0$$

Comparaciones Múltiples

Método de Tukey

- ♦ J. TUKEY (1949) encontró un método para comparar todas las diferencias de medias, que puede usarse a posteriori (y por lo tanto también a priori):

$$H_0) \mu_i = \mu_j$$

$$q_{p,v} = \frac{\bar{y}_i - \bar{y}_j}{\sqrt{s^2/n}}$$

- ♦ Es menos potente que el de Student.
- ♦ La distribución de q está tabulada a simple y doble extremidad.

Comparaciones Múltiples

Método de Tukey – Caso: Pozos Geotérmicos

$$H_0) \mu_1 \leq \mu_2$$

$$q_{3;18} = \frac{135,1 - 116,3}{\sqrt{\frac{1328}{7}}} = \frac{18,9}{13,77} = 1,37 < 3,61 = q_{3;18;95\%} \Rightarrow \text{No rech } H_0$$

$$H_0) \mu_3 \leq \mu_1$$

$$q_{3;18} = \frac{168,1 - 116,3}{\sqrt{\frac{1328}{7}}} = \frac{51,9}{13,77} = 3,77 > 3,61 = q_{3;18;95\%} \Rightarrow \text{Rech } H_0$$

- ♦ Obs: Rechaza más holgado que el método de BONFERRONI, pero menos que el de DUNNETT.

Comparaciones Múltiples

Método de Scheffé

- ♦ H. SCHEFFÉ (1953) encontró un método para ensayar comparaciones generales a posteriori (y por lo tanto también a priori):

$$H_0) \mathcal{T} = \mathcal{T}_0$$

$$F_{p-1, \nu} = \frac{(\hat{\mathcal{T}} - \mathcal{T}_0)^2}{(p-1)s^2 \sum \frac{c_i^2}{n_i}}$$

- ♦ Se puede usar para diferencias de medias pero es conservador, en cambio el método de TUKEY es exacto para ese caso.
- ♦ El método de SCHEFFÉ es bastante más conservador que el de STUDENT pues:

$$t_\nu = \frac{\hat{\mathcal{T}} - \mathcal{T}_0}{\sqrt{s^2 \sum \frac{c_i^2}{n_i}}} \Rightarrow F_{1, \nu} = \frac{(\hat{\mathcal{T}} - \mathcal{T}_0)^2}{s^2 \sum \frac{c_i^2}{n_i}}$$

Comparaciones Múltiples

Método de Scheffé – Caso: Pozos Geotérmicos

- ♦ Analizaremos dos comparaciones no ortogonales:

$$H'_0) \frac{\mu_2 + \mu_3}{2} - \mu_1 \leq 0$$

$$\hat{G}' = \frac{135,1 + 168,1}{2} - 116,3 = 35,4$$

$$\sum \frac{c_i'^2}{n_i} = \frac{1}{n} \sum c_i'^2 = \frac{1}{7} \left[(-1)^2 + \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right] = \frac{1,5}{7}$$

$$F_{2;18} = \frac{(35,4 - 0)^2}{2 \cdot 1328 \cdot \frac{1,5}{7}} = 2,196$$

$$\Rightarrow \alpha^* = 14\% \Rightarrow \text{No rech } H'_0$$

$$H''_0) \mu_3 - \mu_1 = 0$$

$$\hat{G}'' = 168,1 - 116,3 = 51,9$$

$$\sum \frac{c_i''^2}{n_i} = \frac{1}{n} \sum c_i''^2 = \frac{1}{7} \left[(-1)^2 + 1^2 \right] = \frac{2}{7}$$

$$F_{2;18} = \frac{(51,9 - 0)^2}{2 \cdot 1328 \cdot \frac{2}{7}} = 3,544$$

$$\Rightarrow \alpha^* = 5\% \Rightarrow \text{No Rech } H''_0$$

Incumplimiento de Supuestos

Modelo y Supuestos

- ◆ Modelo:

$$\tilde{y}_{ia} = \mu_i + \varepsilon_{ia} = \mu + \alpha_i + \tilde{\varepsilon}_{ia} \quad i = 1 \dots p \quad a = 1 \dots n_i \quad N = \sum n_i$$

- ◆ Supuestos:

$$S1) E \tilde{\varepsilon}_{ia} = 0$$

$$S2) V \tilde{\varepsilon}_{ia} = \sigma^2 \quad (\text{homocedasticidad})$$

$$S3) \tilde{\varepsilon}_{ia} \text{ independientes (no autocorrelación)}$$

$$S4) \tilde{\varepsilon}_{ia} \text{ Normales}$$

- ◆ S1 es parte de la definición de los efectos. No puede fallar.
- ◆ La aleatorización asegura S3.
- ◆ En cambio S2 y S4 pueden fallar. Cuando fallan en general lo hacen juntos: distribuciones anormales y distintas, con distintas varianzas.

Heterocedasticidad

- ♦ Es poco frecuente en experimentos agropecuarios, pero se da en experimentos industriales y comerciales.

- ♦ Diagnóstico:
 - » Método de W. COCHRAN (1941).
 - » Método de H. LEVENE (1960).

- ♦ Resolución:
 - » Transformaciones homogeneizantes.
 - » Métodos no paramétricos: sólo en casos extremos porque son poco potentes.

Heterocedasticidad

Método de Cochran

$$H_0) \sigma_1^2 = \sigma_2^2 = \dots = \sigma_p^2$$

- ♦ W. COCHRAN encontró la distribución del estadístico:

$$H_p = \frac{s_{MAX}^2}{s_1^2 + s_2^2 + \dots + s_p^2}$$

Heterocedasticidad

Método de Levene

- ♦ La desigualdad de Chebishev establece una relacion entre la varianza y las desviaciones absolutas para cualquier variable aleatoria con varianza finita:

$$P(|\tilde{x} - \mu| \geq k) < \sigma^2 / k^2$$

- ♦ Esto sugiere comparar las medias entre grupos de la variable:

$$|y_{ia} - \bar{y}_i|$$

- ♦ Lo cual puede hacerse mediante análisis de la varianza.
- ♦ Si bien esto parece un argumento circular, el ensayo F es suficientemente robusto como para evaluar esta hipótesis de manera aproximada, y eventualmente corregir la heterocedasticidad para una evaluación mas exacta de la hipótesis original.

Heterocedasticidad

Resolución

- ♦ Se aplican transformaciones homogeneizantes:

VARIABLES ECONÓMICAS DE CORTE TRANSVERSAL	$\sigma = k\mu$	$y' = \ln y$
NÚMERO DE EVENTOS DISCRETOS EN UN CONTINUO	$\sigma = k\sqrt{\mu}$	$y' = \sqrt{y}$
DISTRIBUCIÓN GAMMA		$y' = \sqrt[3]{y}$
PORCIONES (BINOMIAL)	$\sigma = \sqrt{\mu(1-\mu)}$	$y' = \arcsen\sqrt{y}$

- ♦ O la transformación general de BOX & COX:

$$y' = \frac{y^\lambda - 1}{\lambda} \quad \text{o} \quad y' = \ln y \quad \text{para } \lambda = 0$$

- » Buscando λ en $[-1;1]$ que optimice la verosimilitud o s^2 .

Anormalidad

- ♦ Es poco frecuente en experimentos agropecuarios, pero se da en experimentos industriales y comerciales.
- ♦ En general viene acompañada de la heterocedasticidad, y las soluciones son comunes.
- ♦ Diagnóstico:
 - » Q-Q plot de residuos.
 - » Modelado de la perturbación aleatoria.
- ♦ Resolución:
 - » Transformaciones normalizantes.
 - » Métodos no paramétricos: sólo en casos extremos porque son poco potentes.

Anormalidad

Q-Q plot

- ♦ Bajo la hipótesis de normalidad:

- » Los residuos estandarizados son una muestra de una variable normal estándar:

$$r_i = \frac{y_{ia} - \bar{y}_i}{s}$$

- » Se ordenan en forma creciente: ${}_o r_1; {}_o r_2; {}_o r_3; \dots {}_o r_N$

- » Y se calcula la función de distribución teórica: $F_N({}_o r_i / 0;1)$

- ♦ Por otro lado se estima la función de distribución empírica:

$$\hat{F}({}_o r_i) = \frac{i}{N+1} \quad \text{o con mayor precisión segun O. Mermoz: } \hat{F}({}_o r_i) = \frac{i-0,3227}{N+0,3546}$$

- ♦ Ambas deberían coincidir.

- ♦ El Q-Q plot es un grafico de los cuantiles teóricos y empíricos:

$${}_o r_i \quad \text{vs} \quad \hat{z}_i = z_{\hat{F}_i}$$

- » Si los puntos están alineados hay normalidad.

Anormalidad

Resolución

- ♦ Las mismas transformaciones homogeneizantes son generalmente normalizantes.
 - » En el caso de la transformación de BOX & COX puede diferir el parámetro λ óptimo para homogeneizar o normalizar. Se busca una solución de compromiso.
- ♦ Los métodos no paramétricos (WILCOXON, KRUSKAL & WALLIS, etc) son el último recurso dado que tienen poca potencia.

Principios del Diseño Experimental

Principios del Diseño Experimental

- ◆ Unidad experimental:
 - » Sujeto sobre el que se miden las variables.
- ◆ Principio de Replicación:
 - » Repetir los tratamientos o factores en varias unidades experimentales para posibilitar y perfeccionar la inferencia estadística.
- ◆ Principio de Aleatorización:
 - » Asignar los tratamientos o factores a las unidades experimentales al azar.
- ◆ Principio de Control Local:
 - » Cuando el ruido experimental es excesivo se aplica esta técnica.
 - » Consiste en clasificar las unidades experimentales en bloques según un factor que influye en la variable respuesta de manera de controlar el ruido.

Diseño en Bloques

Experimentos en Bloques

Caso: Abrasivos

Granulometría	Adhesivo 1	Adhesivo 2	Adhesivo 3	Adhesivo 4	Adhesivo 5	Prom.
Fina	70 60 79	67 71 70	76 65 89	80 88 82	78 88 96	77,3
Media	54 61 62	62 51 59	62 70 72	82 72 65	77 75 74	66,5
Gruesa	50 49 52	49 44 61	46 58 59	44 57 60	67 80 73	56.6
Promedio	59,7	59,3	66,3	70,0	78,7	

Variable respuesta: Porción de carga abrasiva retenida luego de la tarea de lijado normalizada.

Experimentos en Bloques

Tabla ANOVA

<i>Efecto</i>	Q	ν	$CM = Q / \nu$
α_i	$C_\alpha - C$	$p - 1$	
β_j	$C_\beta - C$	$q - 1$	
ε_{ija}	$C_\varepsilon - C_\alpha - C_\beta + C$	$pqn - p - q + 1$	
Total	$C_\varepsilon - C$	$pqn - 1$	

- ♦ Donde las Q siguen las siguientes fórmulas de cálculo:

$$C = pqn \bar{y}_{..}^2$$

$$C_\alpha = qn \sum \bar{y}_{i.}^2$$

$$C_\beta = pn \sum \bar{y}_{.j}^2$$

$$C_\varepsilon = \sum y_{ija}^2$$

Experimentos en Bloques

Inferencia

- ♦ El objetivo es ensayar:

$$H_0) \alpha_1 = \alpha_2 = \dots = \alpha_p$$

- ♦ Y el efecto de los bloques, si se desea ver la efectividad del control local:

$$H_0) \beta_1 = \beta_2 = \dots = \beta_q$$

- ♦ Se convierten estas hipótesis múltiples en simples según la idea de FISHER:

$$H_0) \sigma_\alpha^2 = 0 \qquad H_0) \sigma_\beta^2 = 0$$

- ♦ Utilizando la varianza de las medias, que según la convención de W. Cochran son:

$$\sigma_\alpha^2 = \frac{1}{p-1} \sum \alpha_i^2 \qquad \sigma_\beta^2 = \frac{1}{q-1} \sum \beta_i^2$$

- » No son verdaderas varianzas salvo que los factores sean aleatorios.

Experimentos en Bloques

Inferencia

- Utilizando este artilugio se puede demostrar que las esperanzas de los cuadrados completan la tabla ANOVA de esta manera:

<i>Efecto</i>	Q	ν	$CM = Q / \nu$	\mathcal{ECM}
α_i	$C_\alpha - C$	$p - 1$		$\sigma^2 + qn\sigma_\alpha^2$
β_j	$C_\beta - C$	$q - 1$		$\sigma^2 + pn\sigma_\beta^2$
ε_{ija}	$C_\varepsilon - C_\alpha - C_\beta + C$	$pqn - p - q + 1$		σ^2
Total	$C_\varepsilon - C$	$pqn - 1$		

Diseño Factorial a 2 Factores Fijos Cruzados

Experimentos a 2 Factores Cruzados Fijos

Caso: Panel de Internet

Duración cuestionario	Nivel Socioeconómico					
	ABC			DE		
10 min	14,0	13,5	11,5	11,5	7,0	15,5
	11,0	14,5	13,5	10,5	9,5	13,0
20 min	14,0	9,0	14,0	7,0	8,0	8,5
	12,5	11,5	16,0	5,0	7,0	9,0
30 min	9,5	13,0	13,5	3,5	5,5	6,0
	10,5	8,0	11,0	4,5	3,5	6,5

α : duración cuest.

β : nivel socioeconómico

$p = 3$

$q = 2$

$n = 6$

Variable respuesta: Tasa de respuesta = Porción de la gente invitada a responder la encuesta que la completa.

Cada dato corresponde a 200 invitaciones.

Experimentos

Modelo Factorial a 2 Factores Cruzados

- ♦ El modelo lineal:

$$\tilde{y}_{ija} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + \tilde{\varepsilon}_{ija} \quad i=1\dots p \quad j=1\dots q \quad a=1\dots n$$

- ♦ Condiciones de unicidad de efectos:

$$\sum_{i=1}^p \alpha_i = 0 \quad \sum_{j=1}^q \beta_j = 0 \quad \forall i, j: \sum_{i=1}^p \alpha\beta_{ij} = 0 = \sum_{j=1}^q \alpha\beta_{ij}$$

- ♦ Supuestos:

$$S1) E\tilde{\varepsilon}_{ija} = 0$$

$$S2) V\tilde{\varepsilon}_{ija} = \sigma^2 \quad (\text{homocedasticidad})$$

$$S3) \tilde{\varepsilon}_{ija} \quad \text{independientes}$$

$$S4) \tilde{\varepsilon}_{ija} \quad \text{Normales}$$

Experimentos

Tabla ANOVA

<i>Efecto</i>	Q	ν	$CM = Q / \nu$
α_i	$C_\alpha - C$	$p - 1$	
β_j	$C_\beta - C$	$q - 1$	
$\alpha\beta_{ij}$	$C_{\alpha\beta} - C_\alpha - C_\beta + C$	$(p - 1)(q - 1)$	
ε_{ija}	$C_\varepsilon - C_{\alpha\beta}$	$pq(n - 1)$	
Total	$C_\varepsilon - C$	$pqn - 1$	

- ♦ Donde las Q siguen las fórmulas de cálculo:

$$C = pqn \bar{y}_{..}^2$$

$$C_\alpha = qn \sum \bar{y}_{i.}^2$$

$$C_\beta = pn \sum \bar{y}_{.j}^2$$

$$C_{\alpha\beta} = n \sum \bar{y}_{ij}^2$$

$$C_\varepsilon = \sum y_{ija}^2$$

Experimentos Inferencia

- ♦ El objetivo es ensayar:

$$H_0) \alpha_1 = \alpha_2 = \dots = \alpha_p \quad H_0) \beta_1 = \beta_2 = \dots = \beta_q$$

- ♦ Evaluar hipótesis múltiples es complicado, pero Fisher encontró la forma de convertirlas en hipótesis simples:

$$H_0) \sigma_\alpha^2 = 0 \quad H_0) \sigma_\beta^2 = 0$$

- ♦ Utilizando la varianza de las medias, que según la convención de W. Cochran son:

$$\sigma_\alpha^2 = \frac{1}{p-1} \sum \alpha_i^2 \quad \sigma_\beta^2 = \frac{1}{q-1} \sum \beta_i^2 \quad \sigma_{\alpha\beta}^2 = \frac{1}{(p-1)(q-1)} \sum \alpha\beta_{ij}^2$$

- » No son verdaderas varianzas salvo que los factores sean aleatorios.

Experimentos Inferencia

- ♦ Utilizando este artilugio se puede demostrar que las esperanzas de los cuadrados completan la tabla ANOVA de esta manera:

<i>Efecto</i>	Q	ν	$CM = Q / \nu$	\mathcal{ECM}
α_i	$C_\alpha - C$	$p - 1$		$\sigma^2 + qn\sigma_\alpha^2$
β_j	$C_\beta - C$	$q - 1$		$\sigma^2 + pn\sigma_\beta^2$
$\alpha\beta_{ij}$	$C_{\alpha\beta} - C_\alpha - C_\beta + C$	$(p - 1)(q - 1)$		$\sigma^2 + n\sigma_{\alpha\beta}^2$
ε_{ija}	$C_\varepsilon - C_{\alpha\beta}$	$pq(n - 1)$		σ^2
Total	$C_\varepsilon - C$	$pqn - 1$		

Diseño Factorial a 2 Factores Cruzados Fijos o Aleatorios

Experimentos a 2 Factores Cruzados Fijos o Aleatorios

Caso: Hilo de Poliéster

\bar{y}_{ij}	Turno 1	Turno 2	Turno 3
Maq. 1	3,786	3,810	3,702
Maq. 2	3,620	3,674	3,286
Maq. 3	3,656	3,148	3,180

α : máquina

β : turno

$p = 3$ $P = 8$

$q = 3$ $Q = 3$

$n = 5$

$$\sum y_{ija}^2 = 568,1979$$

Variable respuesta: Resistencia del hilo.

Experimentos a 2 Factores Cruzados

Tabla ANOVA

<i>Efecto</i>	Q	ν	$\mathbb{E}Q / \nu$
α_i	$C_\alpha - C$	$p - 1$	$\sigma^2 + (1 - q/Q)n\sigma_{\alpha\beta}^2 + qn\sigma_\alpha^2$
β_j	$C_\beta - C$	$q - 1$	$\sigma^2 + (1 - p/P)n\sigma_{\alpha\beta}^2 + pn\sigma_\beta^2$
$\alpha\beta_{ij}$	$C_{\alpha\beta} - C_\alpha - C_\beta + C$	$(p - 1)(q - 1)$	$\sigma^2 + n\sigma_{\alpha\beta}^2$
ε_{ija}	$C_\varepsilon - C_{\alpha\beta}$	$pq(n - 1)$	σ^2
Total	$C_\varepsilon - C$	$pqn - 1$	

- Donde las Q siguen las fórmulas de cálculo:

$$C = pqn \bar{y}_{..}^2$$

$$C_\alpha = qn \sum \bar{y}_{i.}^2$$

$$C_\beta = pn \sum \bar{y}_{.j}^2$$

$$C_{\alpha\beta} = n \sum \bar{y}_{ij}^2$$

$$C_\varepsilon = \sum y_{ija}^2$$

Diseño Factorial a 2 Factores Anidados

Experimentos a 2 Factores Anidados

Caso: Restaurants

Grandes ciudades			Ciudades chicas		
G1	G2	G3	C1	C2	C3
72	71	43	42	10	23
52	78	64	55	37	51
56	65	60	57	31	33
42	67	73	46	43	54
55,5	70,3	60,0	50,0	30,3	40,3
61,9			40,2		

α : tamaño ciudad

β : local

$p = 2$ $P = 2$

$q = 3$ $Q \rightarrow \infty$

$n = 4$

Variable respuesta: Porción de las mesas que piden vino con la comida.

Experimentos a 2 Factores Anidados

Tabla ANOVA

<i>Efecto</i>	Q	ν	$\mathbb{E}Q / \nu$
α_i	$C_\alpha - C$	$p - 1$	$\sigma^2 + (1 - q/Q)n\sigma_\beta^2 + qn\sigma_\alpha^2$
β_j	$C_{\alpha\beta} - C_\alpha$	$p(q - 1)$	$\sigma^2 + n\sigma_\beta^2$
ε_{ija}	$C_\varepsilon - C_{\alpha\beta}$	$pq(n - 1)$	σ^2
Total	$C_\varepsilon - C$	$pqn - 1$	

- ♦ Donde las Q siguen las fórmulas de cálculo:

$$C = pqn \bar{y}_{..}^2$$

$$C_\alpha = qn \sum \bar{y}_{i.}^2$$

$$C_{\alpha\beta} = n \sum \bar{y}_{ij}^2$$

$$C_\varepsilon = \sum y_{ija}^2$$

Diseño Factorial a 3 Factores

Experimentos a 3 Factores Cruzados

Caso: Film de polietileno

♦ Promedios:

		Ctrl	COC	EVOH
M A Q 1	Ctrl	501,8	297	255,4
	Enfr.	311,8	200,2	156,8
	E. ráp.	197,2	145,2	108,8
M A Q 2	Ctrl	499,2	302	250,4
	Enfr.	306	205,6	173
	E. ráp.	225,4	181,4	126,6

α : Enfriamiento $p = 3$ $P = 3$

β : Bi-laminado $q = 3$ $Q = 3$

γ : Máquina $r = 2$ $R = 5$

$n = 4$

$$\sum y_{ijka}^2 = 6.709.748$$

$$\bar{y}_{...} = 246,867$$

	Ctrl	COC	EVOH
Ctrl	500	300	253
Enfr.	309	203	165
E. ráp.	211	163	118

	Maq 1	Maq 2
Ctrl	352	351
Enfr.	223	228
E. ráp.	150	178

	Maq 1	Maq 2
Ctrl	337	344
COC	214	230
EVOH	174	183

Experimentos a 3 Factores Cruzados

Tabla ANOVA

<i>Efecto</i>	Q	ν	$\mathbb{E}Q / \nu$
α_i	$C_\alpha - C$	$p-1$	$\sigma^2 + q'r'n\sigma_{\alpha\beta\gamma}^2 + q'rn\sigma_{\alpha\beta}^2 + qr'n\sigma_{\alpha\gamma}^2 + qrn\sigma_\alpha^2$
β_j	$C_\beta - C$	$q-1$	$\sigma^2 + p'r'n\sigma_{\alpha\beta\gamma}^2 + p'rn\sigma_{\alpha\beta}^2 + pr'n\sigma_{\beta\gamma}^2 + prn\sigma_\beta^2$
γ_k	$C_\gamma - C$	$r-1$	$\sigma^2 + p'q'n\sigma_{\alpha\beta\gamma}^2 + p'qn\sigma_{\alpha\gamma}^2 + pq'n\sigma_{\beta\gamma}^2 + pqn\sigma_\gamma^2$
$\alpha\beta_{ij}$	$C_{\alpha\beta} - C_\alpha - C_\beta + C$	$(p-1)(q-1)$	$\sigma^2 + r'n\sigma_{\alpha\beta\gamma}^2 + rn\sigma_{\alpha\beta}^2$
$\alpha\gamma_{ik}$	$C_{\alpha\gamma} - C_\alpha - C_\gamma + C$	$(p-1)(r-1)$	$\sigma^2 + q'n\sigma_{\alpha\beta\gamma}^2 + qn\sigma_{\alpha\gamma}^2$
$\beta\gamma_{jk}$	$C_{\beta\gamma} - C_\beta - C_\gamma + C$	$(q-1)(r-1)$	$\sigma^2 + p'n\sigma_{\alpha\beta\gamma}^2 + pn\sigma_{\beta\gamma}^2$
$\alpha\beta\gamma_{ijk}$	$C_{\alpha\beta\gamma} - C_{\alpha\beta} - C_{\alpha\gamma} - C_{\beta\gamma}$	$(p-1)(q-1)$	$\sigma^2 + n\sigma_{\alpha\beta\gamma}^2$
	$+C_\alpha + C_\beta + C_\gamma - C$	$(r-1)$	
ε_{ijka}	$C_\varepsilon - C_{\alpha\beta\gamma}$	$pqr(n-1)$	σ^2
Total	$C_\varepsilon - C$	$pqrn-1$	

donde: $p' = 1 - \frac{p}{P}$ $q' = 1 - \frac{q}{Q}$ $r' = 1 - \frac{r}{R}$

Experimentos a 3 Factores Cruzados

Tabla ANOVA

- ♦ Donde las Q siguen las fórmulas de cálculo:

$$C = pqrn \bar{y}_{\dots}^2$$

$$C_{\alpha} = qrn \sum \bar{y}_{i\cdot\cdot}^2$$

$$C_{\beta} = prn \sum \bar{y}_{\cdot j\cdot}^2$$

$$C_{\gamma} = pqn \sum \bar{y}_{\cdot\cdot k}^2$$

$$C_{\alpha\beta} = rn \sum \bar{y}_{ij\cdot}^2$$

$$C_{\alpha\gamma} = qn \sum \bar{y}_{i\cdot k}^2$$

$$C_{\beta\gamma} = pn \sum \bar{y}_{\cdot jk}^2$$

$$C_{\alpha\beta\gamma} = n \sum \bar{y}_{ijk}^2$$

$$C_{\varepsilon} = \sum y_{ijka}^2$$