

9. Análisis de varianza

9.1. Introducción. Conceptos básicos

Vamos a introducir en este tema las técnicas de análisis de la varianza. El análisis de varianza es una herramienta muy importante para comparar las medias distintas poblaciones. Su desarrollo, como se verá posteriormente, tiene una filosofía muy parecida a la de la regresión.

En esta sección introduciremos los conceptos básicos de análisis de la varianza. Empecemos con un ejemplo; así, consideremos el número de kilos de fruta que da un manzano. Tenemos entonces que la cantidad de fruta de un árbol concreto es una variable aleatoria, pues no puede predecirse *a priori* cuál será el resultado. Para esta variable aleatoria podemos realizar los estudios típicos de la inferencia estadística que hemos visto en los dos temas anteriores; por ejemplo, a partir de una m.a.s. podemos calcular un intervalo de confianza para la producción media de manzanas o cualquier otra característica de interés; podemos plantear el contraste de hipótesis para estudiar si la varianza de esta variable supera 4 o no; podemos plantear un contraste no paramétrico para ver si la producción sigue distribución normal; etc. Todos estos estudios nos darán información que puede ser muy valiosa para conocer el comportamiento *general* de la variable.

El análisis de la varianza plantea el problema desde otra perspectiva: ¿por qué no todos los árboles producen la misma cantidad de manzanas? O en otras palabras, ¿hay alguna característica que influye en la producción de manzanas? Por ejemplo, es posible que las distintas subespecies de manzano den lugar a diferencias en la producción de fruta; otra causa que podría influir en el resultado final de la variable producción es el tipo de abono utilizado; y así podemos obtener muchas

más características que puedan explicar las diferencias entre los distintos valores. Es decir, el análisis de la varianza se centra en estudiar las causas que originan las diferencias entre los resultados de los distintos individuos.

Este planteamiento es muy interesante en la práctica; así, si hay diferencias entre las distintas subespecies, podemos elegir la que dé una mayor producción; si hay dos abonos que potencian la producción de la misma manera, podemos elegir el que sea más barato, etc.; y en general, podemos buscar la mejor combinación entre subespecie y abono que maximice la producción. Esto hace que el análisis de varianza se aplique en muchos problemas médicos, industriales y económicos. Nótese la similitud con regresión, en la que tenemos una variable aleatoria y estudiamos si conocer el valor de otra u otras variables sobre un individuo nos dan información sobre el valor que toma la variable original sobre ese individuo.

Desde un punto de vista más formal, consideremos una v.a. X , en ocasiones llamada **variable dependiente**. Sobre esta variable aleatoria se cree que pueden influir una serie de características; cada una de estas características se llama **factor**; es usual denotar los factores por A, B , y así sucesivamente. Por ejemplo, en el caso anterior la variedad o el tipo de abono serían factores.

Cada factor puede tomar una serie de valores, *no necesariamente numéricos*. De hecho, esta será una de las diferencias con regresión, en la que las variables que influyen sobre la variable dependiente son numéricas y se busca una relación numérica que dé una aproximación de la variable dependiente a partir de los valores de las otras variables. Si estamos considerando el factor A , las distintas modalidades del factor se denotan por A_1, A_2, \dots y se llaman los **niveles** del factor A . Por ejemplo, en el caso anterior tenemos que las variedades fuji, reineta o golden son niveles del factor variedad. Nótese que las modalidades que puede tomar un factor pueden ser infinitas o finitas. Así, si consideramos la cantidad de lluvia como un factor, entonces tenemos infinitos niveles.

Cuando estamos estudiando el comportamiento de la variable X es posible que no podamos estudiar todos los niveles de un factor. Esto es debido a que el número de niveles es muy elevado (incluso infinito) o porque no podamos controlar la modalidad de la característica; por ejemplo, si estamos considerando la temperatura mínima, este es un

factor que no podemos controlar y tendremos que conformarnos con los valores que aparecen al realizar el estudio. En estas condiciones, diremos que el factor es de **efectos aleatorios**. Para un factor de efectos aleatorios se supone que los niveles del factor se escogen al azar, ya sea por el investigador, que los fija mediante un sorteo antes de comenzar el estudio, o por la naturaleza, que los determina al realizarse el experimento.

Por otra parte, es posible que sí seamos capaces a estudiar todas las modalidades del factor *que nos interesan*. Nótese que es posible que el factor tome muchas más modalidades, pero que no sean relevantes para nuestro estudio; por ejemplo, de entre todos los posibles abonos tal vez nos queramos centrar en los que no superan un nivel determinado de cierto componente muy contaminante; en este caso, los abonos que no están en estas condiciones dejan de ser niveles del factor abono. Si estudiamos todos los niveles que nos interesan diremos que el factor es de **efectos fijos**.

Aunque en el caso que estudiaremos nosotros, que es el modelo con un solo factor, los desarrollos matemáticos son iguales para un factor de efectos fijos que para un factor de efectos aleatorios, para modelos más generales sí existen diferencias en los cálculos a realizar.

En todo caso, sea un factor de efectos fijos o de efectos aleatorios, siempre tendremos en la práctica una cantidad finita de niveles. Denotaremos por a el número de niveles que se consideran para el factor A , por b el número de niveles que se consideran para el factor B , y así sucesivamente. Así, los niveles de A se denotan por A_1, \dots, A_a , los de B por B_1, \dots, B_b , y así sucesivamente.

Inicialmente, el objetivo del análisis de varianza es estudiar si verdaderamente existe una influencia de los distintos factores sobre el resultado de la variable X . Por otra parte, además de una influencia aislada del factor, es posible que existan interacciones entre los factores, de forma que la combinación de niveles tenga una influencia superior o inferior de la que tendría cada uno de los factores por separado y luego sumando todas esas influencias. Por ejemplo, es posible que una variedad de manzana con un determinado abono produzca muchas más manzanas que otra variedad con el mismo tipo de abono. Por ello, tenemos que estudiar todas las posibles combinaciones de factores. Si tenemos varios factores, una combinación de un nivel de cada factor se llama

un **tratamiento**. Así, si tenemos dos factores tenemos los tratamientos $A_1B_1, A_1B_2, \dots, A_aB_b$, en total $a \times b$ tratamientos. La tabla 9.1 se muestra el proceso para dos factores.

	B_1	B_2	...	B_b
A_1	A_1B_1	A_1B_2	...	A_1B_b
A_2	A_2B_1	A_2B_2	...	A_2B_b
\vdots	\vdots	\vdots	\vdots	\vdots
A_a	A_aB_1	A_aB_2	...	A_aB_b

Tabla 9.1. Construcción de los distintos tratamientos en el caso de dos factores.

En definitiva, lo que estamos haciendo es particionar una población, que viene dada por la variable X , en varias subpoblaciones, en cada una de las cuales los niveles de cada factor están fijados. Denotaremos por $X_{i_1, i_2, \dots}$ la variable que sigue la subpoblación en la que el primer factor toma el nivel i_1 , el segundo factor el nivel i_2 , etcétera. Si tenemos k factores y denotamos por r_i el número de niveles que estamos considerando en el factor i , entonces estamos comparando

$$\prod_{i=1}^k r_i$$

poblaciones distintas.

$$X \begin{cases} X_{11\dots} \\ \vdots \\ X_{r_1 r_2 \dots} \end{cases}$$

El objetivo entonces es estudiar si todas las variables $X_{i_1, i_2, \dots}$ siguen la misma distribución, lo que significaría que los distintos tratamientos influyen por igual para la variable que se está considerando o, por el contrario, existen diferencias entre ellos. Por lo tanto, el contraste que nos planteamos es:

$$\begin{cases} H_0 : X_{i_1, i_2, \dots} \text{ homogéneas (misma distribución)} \\ H_1 : \text{No } H_0 \end{cases} .$$

Esto es un contraste no paramétrico, pues no conocemos la distribución de cada una de esas variables. Sin embargo, nosotros supondremos las siguientes condiciones sobre las poblaciones:

- **Normalidad:** Todas las subpoblaciones siguen distribución normal. Es decir, supondremos que $X_{i_1, i_2, \dots} \sim \mathcal{N}(\mu_{i_1, i_2, \dots}, \sigma_{i_1, i_2, \dots})$.
- **Homocedasticidad:** Todas las subpoblaciones tienen la misma varianza. Es decir, $\sigma_{i_1, i_2, \dots} = \sigma, \forall i_1, i_2, \dots$
- **Independencia:** Todas las subpoblaciones son independientes entre sí.

Como consecuencia de estas hipótesis previas, $X_{i_1, i_2, \dots} \sim \mathcal{N}(\mu_{i_1, i_2, \dots}, \sigma)$, donde las distintas medias $\mu_{i_1, i_2, \dots}$ y la desviación típica común σ son desconocidas. El problema de comprobar su homogeneidad se reduce ahora a contrastar la igualdad de medias

$$\begin{cases} H_0 : \mu_{i_1, i_2, \dots} = \mu, \forall i_1, i_2, \dots \\ H_1 : \text{No } H_0 \end{cases},$$

donde μ es un valor común para todos los tratamientos y que coincide con la media de la variable X . Nótese además que por la hipótesis de independencia, el análisis de la varianza es una generalización del contraste de igualdad de medias para poblaciones normales independientes con varianzas iguales y desconocidas que se estudió en el capítulo anterior. Los distintos modelos de análisis de la varianza permiten resolver estos contrastes en distintas situaciones.

9.2. Modelo unifactorial

En esta sección veremos el modelo más sencillo de análisis de varianza. En el modelo unifactorial tenemos un único factor de interés. Denotaremos este factor por A y supondremos que se evalúan a niveles. Cada nivel del factor puede influir positivamente o negativamente en el valor de la variable.

Como hemos visto, el análisis de la varianza se basa en descomponer la variable X en varias variables, una por cada tratamiento. En este caso, como solo hay un factor, los tratamientos coinciden con los niveles del factor, y así los escribimos como X_1, \dots, X_a , donde estas variables corresponden a la variable X según los distintos niveles del factor. Como en el caso general, nosotros supondremos las siguientes condiciones sobre las poblaciones:

- **Normalidad:** Todas las subpoblaciones siguen distribución normal. Es decir, supondremos que $X_i \sim \mathcal{N}(\mu_i, \sigma_i)$.
- **Homocedasticidad:** Todas las subpoblaciones tienen la misma varianza. Es decir, $\sigma_i = \sigma, \forall i = 1, \dots, a$.
- **Independencia:** Todas las subpoblaciones son independientes.

Entonces, por estas condiciones previas de análisis de la varianza, tenemos que $X_i \sim \mathcal{N}(\mu_i, \sigma)$, donde σ es un valor desconocido de la desviación típica, común a todas las variables, y $\mu_i, i = 1, \dots, a$ son los valores de las medias (también desconocidas). El contraste que queremos resolver es si hay diferencias entre las distintas variables X_1, \dots, X_a . Por las condiciones previas, la única diferencia posible es en la media, lo que se traduce en el contraste

$$\begin{cases} H_0 : \mu_1 = \mu_2 = \dots = \mu_a \\ H_1 : \text{No } H_0 \end{cases} .$$

Para resolver este contraste, tenemos que estudiar el comportamiento de cada una de las variables. Para ello, tenemos a muestras aleatorias simples, de tamaños n_1, \dots, n_a , respectivamente, una para cada una de las variables (una para cada nivel del factor). Estas muestras vienen dadas por

$$(X_{11}, \dots, X_{1n_1}), \dots, (X_{a1}, \dots, X_{an_a}).$$

De esta manera los datos de todas esas muestras pueden escribirse todos juntos en una tabla tal y como se da en la tabla 9.2.

Tenemos entonces un número de datos dado por:

$$n_1 + \dots + n_a = n.$$

Nivel 1	Nivel 2	...	Nivel a
x_{11}	x_{21}	...	x_{a1}
x_{12}	x_{22}	...	x_{a2}
\vdots	\vdots	\vdots	\vdots
x_{1n_1}	x_{2n_2}	...	x_{an_a}

Tabla 9.2. Datos muestrales para resolver un análisis de la varianza unifactorial.

Nótese que los tamaños de muestra pueden ser diferentes para las distintas muestras, lo mismo que pasaba para el contraste de medias de dos poblaciones normales que se vio en el capítulo anterior.

Como se trata de comparar las medias, lo lógico es que nuestras conclusiones dependan de las estimaciones de cada una de ellas. Para cada subpoblación podemos estimar la media teórica mediante la correspondiente media muestral; la media correspondiente al nivel i se denota por \bar{X}_i y viene dada por

$$\bar{X}_i := \frac{1}{n_i} \sum_{j=1}^{n_i} X_{ij},$$

que para unos datos concretos $(x_{i1}, \dots, x_{in_i})$ nos da el valor

$$\bar{x}_i := \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij}.$$

De manera análoga podemos estimar la media de la población global mediante la media muestral de todos los datos, ya que podemos considerar la unión de todas las muestras como una muestra de tamaño n de la variable X ; denotaremos esta media por \bar{X} ; este valor viene dado por

$$\bar{X} = \frac{1}{n} \sum_{i=1}^a \sum_{j=1}^{n_i} X_{ij}.$$

Es fácil comprobar que estos valores están relacionados mediante la fórmula

$$\bar{X} = \frac{n_1\bar{X}_1 + n_2\bar{X}_2 + \dots + n_a\bar{X}_a}{n},$$

y los valores muestrales por

$$\bar{x} = \frac{n_1\bar{x}_1 + n_2\bar{x}_2 + \dots + n_a\bar{x}_a}{n}.$$

Si la hipótesis nula es cierta, entonces $\mu_i = \mu_j, \forall i, j = 1, \dots, a$. Esto significa que nosotros esperamos que \bar{X}_i y \bar{X}_j están estimando la misma cantidad, por lo que es de esperar que \bar{x}_i y \bar{x}_j sean valores similares (nótese que no es de esperar igualdad porque estamos haciendo una estimación). Por otra parte, si la hipótesis nula es cierta, todas las medias teóricas μ_i coinciden con μ , la media (desconocida) de X , con lo que \bar{X}_i y \bar{X} están estimando la misma cantidad y es de esperar que \bar{x}_i y \bar{x} sean valores similares; en esto último es en lo que nos vamos a basar para resolver el contraste.

Si \bar{X}_i y \bar{X} son similares, entonces $(\bar{X}_i - \bar{X})^2$ será pequeño, y en consecuencia

$$\sum_{i=1}^a n_i(\bar{X}_i - \bar{X})^2$$

será pequeño. Aquí n_i aparece porque si los tamaños de muestra son grandes, esperamos que la estimación sea mejor, y entonces penalizaremos las diferencias grandes más que si el tamaño de muestra es pequeño.

Por el contrario, si la hipótesis nula es falsa, hay diferencias entre algún μ_i y μ_j ; entonces, \bar{X}_i y \bar{X}_j estiman cantidades diferentes y es de esperar que \bar{x}_i y \bar{x}_j se diferencien mucho. Procediendo como antes, es de esperar también que \bar{x}_i y \bar{x} se diferencien y así $(\bar{x}_i - \bar{x})^2$ debería ser grande. Entonces, si la hipótesis alternativa es cierta, esperamos que

$$\sum_{i=1}^a n_i(\bar{X}_i - \bar{X})^2$$

sea grande. Entonces, nuestra región crítica debería ser de la forma

$$R.C. = \{(x_{11}, \dots, x_{1n}, \dots, x_{a1}, \dots, x_{an}) t.q. \sum_{i=1}^a n_i(\bar{x}_i - \bar{x})^2 > k\}.$$

Tenemos entonces que calcular el valor de la constante k que marca la frontera entre la región crítica y la región de aceptación. Supongamos que hemos obtenido $\sum_{i=1}^a n_i(\bar{x}_i - \bar{x})^2 = 4,3$. ¿Es este valor suficientemente grande para concluir la hipótesis alternativa? En otras palabras, ¿cuál es el valor de k para un nivel de significación α concreto?

Nótese en primer lugar que, si la hipótesis nula es cierta, entonces las diferencias entre \bar{X}_i y \bar{X} dependerán de la varianza σ^2 , valor que es desconocido. De esta manera, 4.3 será grande para varianzas poblacionales muy pequeñas, pero será un valor que podemos achacar a la aleatoriedad de la muestra si la varianza poblacional es muy grande. Lo mismo sucede si tenemos un error de medición entre dos distancias de un metro. Si es sobre una distancia de varios kilómetros, el error que se está cometiendo es muy pequeño, pero si es sobre una distancia de cinco metros, el error de medición es muy grande. ¿Qué hacer entonces?

Para resolver el problema volveremos al inicio de nuestro estudio, en el que nos planteamos por qué no son iguales todas las observaciones, y procederemos de forma similar a como se hizo para regresión lineal. Para ello tenemos que pensar en que las posibles diferencias en el valor de la variable X para dos individuos de la población pueden ser debidas a dos causas:

- Si los individuos están en niveles del factor diferentes, como la medida de cada individuo es un valor cercano a su media más o menos una cantidad debida a la variación σ de la población, entonces si las medias para distintos niveles son diferente, es de esperar que haya una diferencia entre los valores obtenidos para esos dos individuos. Es decir, que variar el nivel puede producir diferencias. Diremos que esta posible variación es debida al factor.
- Por otra parte, individuos que están en las mismas condiciones para el factor, es decir, están bajo el mismo nivel, no tendrán el mismo valor. Eso es debido a que en este caso, si tenemos el nivel A_i , tenemos dos individuos con distribución $\mathcal{N}(\mu_i, \sigma)$, y debido a que tenemos una variable aleatoria (en otras palabras, σ no es 0), es de esperar una diferencia entre los valores para los dos individuos. Diremos que esta diferencia es debida a un *error aleatorio*.

Una medida de la variación de todas las observaciones viene dada por

$$\sum_{i=1}^a \sum_{j=1}^{n_i} (X_{ij} - \bar{X})^2 = (n - 1)S^2 = nV(X).$$

Este término se llama la **suma de cuadrados del total** y se denota *SCT*. Puede demostrarse que esta suma de cuadrados puede descomponerse en dos sumandos mediante

$$\sum_{i=1}^a \sum_{j=1}^{n_i} (X_{ij} - \bar{X})^2 = \sum_{i=1}^a \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2 + \sum_{i=1}^a \sum_{j=1}^{n_i} (\bar{X}_i - \bar{X})^2.$$

El primer término se llama **suma de cuadrados del error** y se denota por *SCE*; el segundo término se llama **suma de cuadrados del factor** y se denota por *SCA*. Así, la variación total de la muestra se ha desglosado en dos términos:

$$SCT = SCE + SCA.$$

Nótese que se tiene lo siguiente:

$$SCE = \sum_{i=1}^a \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2 = \sum_{i=1}^a (n_i - 1)S_i^2 = \sum_{i=1}^a n_i V(X_i),$$

$$SCA = \sum_{i=1}^a \sum_{j=1}^{n_i} (\bar{X}_i - \bar{X})^2 = \sum_{i=1}^a n_i (\bar{X}_i - \bar{X})^2.$$

Veamos estos dos términos por separado. La suma de cuadrados del error es una suma de las varianzas muestrales en cada muestra multiplicadas por n_i . Como dentro de cada muestra la población tiene la misma distribución ($\mathcal{N}(\mu_i, \sigma)$), este término será más o menos grande si σ es más o menos grande. En otras palabras, este término es debido a errores aleatorios en las muestras y no depende de las diferencias entre los niveles.

En el caso de *SCA*, ya hemos visto que es el término que depende tanto del error aleatorio (que hace que no coincidan las medias muestrales con las teóricas) como de posibles diferencias entre los distintos

niveles (pues al estimar cantidades distintas hay una tendencia a que aumente el valor absoluto de sus diferencias).

Entonces tenemos que comparar SCA con SCE. El problema ahora es que estas cantidades no están en las mismas «unidades». En concreto, se puede demostrar que

$$E(SCE) = (n - a)\sigma^2, \quad E(SCA) = (a - 1)\sigma^2 + c,$$

donde c es una constante positiva que mide las posibles diferencias entre los niveles y que vale 0 si no hay diferencias entre ellos.

Definimos entonces el **cuadrado medio del factor** como

$$CMA = \frac{SCA}{(a - 1)}.$$

Análogamente, definimos el **cuadrado medio del error** como

$$CME = \frac{SCE}{(n - a)}.$$

Ahora sí podemos realizar la comparación: Al dividir CMA entre CME , si no hay diferencias entre los niveles, entonces están estimando la misma cantidad y esperamos que nos queden valores similares, con lo que el cociente será cercano a 1. Si hay diferencias entre los niveles, entonces c es positivo y estamos estimando cantidades diferentes y esperamos que ese cociente sea mayor que 1. Nuevamente tenemos que encontrar un valor para el que consideremos que es un valor suficientemente grande como para concluir que hay diferencias entre los niveles. La ventaja que tenemos ahora es que este cociente es adimensional. Se puede demostrar que

$$\frac{CMA}{CME} \sim F_{a-1, (n-a)}.$$

En definitiva, rechazaremos para valores grandes de $\frac{CMA}{CME}$. Es decir, que nuestra región crítica para un nivel de significación α será

$$R.C. = \{(x_{11}, \dots, x_{1n_1}, \dots, x_{a1}, \dots, x_{an_a}) t.q. \frac{cma}{cme} > F_{a-1, n-a; \alpha}\}.$$

En general, todos los modelos de análisis de varianza se resuelven rellenando una tabla. En el caso del modelo unifactorial la tabla es la que se puede ver en la tabla 9.3.

F. Variación	SC	g.l.	CM	F
Factor A	sca	a-1	cma	$\frac{cma}{cme}$
Error	sce	n-a	cme	
Total	sct	n-1		

Tabla 9.3. Tabla ANOVA para el modelo unifactorial.

Nótese también que se tiene el mismo estimador para resolver el modelo de efectos fijos y el modelo de efectos aleatorios. Como hemos dicho anteriormente, esto es algo que no ocurre en general. Sí cambia la interpretación de los resultados:

- En un modelo de efectos fijos, si concluimos H_0 , entonces afirmamos que no hay diferencias entre los niveles del factor. Si concluimos H_1 , entonces existen diferencias.
- En un modelo de efectos aleatorios, si concluimos H_0 , entonces afirmamos que no hay diferencias entre los niveles del factor, *incluyendo los niveles que no han sido estudiados*. Si concluimos H_1 , entonces existen diferencias.

Ejemplo 136.

Se están comparando cuatro métodos de estudio. Para ello, se prueba cada método y se observa la puntuación obtenida en una prueba. El profesor supone que, aparte de la variación usual en la nota para alumnos que han utilizado el mismo método, puede existir una variación significativa de la nota debido a que la eficacia de los distintos métodos sea distinta. Para investigar esto, selecciona 16 alumnos al azar que divide en cuatro grupos de cuatro alumnos cada uno, y evalúa a cada alumno en una prueba. Los datos recogidos aparecen en la tabla 9.4.

Suponiendo que se cumplen las condiciones de ANOVA, ¿podemos concluir que hay diferencias entre los métodos de estudio?

Método / Alumno.	1	2	3	4	Media
1	8	7	9	6	7,5
2	1	0	3	2	1,5
3	6	5	7	5	5,75
4	5	6	9	8	7

Tabla 9.4. Calificaciones correspondientes al ejemplo 136.

Aquí la variable X es la nota que se obtiene y el factor es el método (con cuatro niveles); es un factor de efectos fijos, pues solo estamos interesados en comparar esos cuatro métodos.

En este problema, se tiene que la media global es $\bar{x} = 5,4375$. Aplicando un análisis de la varianza se obtiene la tabla 9.5.

F. Variación	SC	g.l.	CM	F
Método	89.19	3	29.73	15.68
Error	22.75	12	1.9	
Total	111.94	15		

Tabla 9.5. Tabla ANOVA para el ejemplo 136.

El valor 15.68 corresponde a un p -valor inferior a 0.05, puesto que $F_{3,12;0,05} = 3,49$, por lo que rechazamos la hipótesis nula y concluimos que existen diferencias entre los métodos de aprendizaje.