

La hipótesis nula y la significación práctica

Frías, M. D.; Pascual, J. y García, J.F.
Universidad de Valencia

Durante la década de los noventa del siglo XX se han publicado debates y opiniones (se puede consultar una revisión exhaustiva en Nickerson, 2000) que han cuestionado, en ocasiones duramente, la aplicación del procedimiento de significación de la hipótesis nula (Null Hypothesis Significance Testing, NHST) como instrumento para el progreso del conocimiento científico. Las críticas al procedimiento incluyen desde las falsas concepciones sobre la información que facilita el procedimiento, provocando interpretaciones erróneas de sus resultados, hasta su uso inadecuado como medio de obtención de datos de significación práctica, sustantiva o clínica. El objetivo de nuestro estudio es abordar la problemática no resuelta aún sobre la utilización apropiada de NHST dentro del área aplicada de la comprobación de efectos del tratamiento y del cambio clínico así como analizar la dificultad de la especificación a priori del valor del tamaño del efecto y analizar las diferentes propuestas que actualmente se proponen al procedimiento. Entre dichas propuestas destaca la estimación del tamaño del efecto y el planteamiento de tamaños de efecto mínimo y datos normativos que enriquecen el diseño de investigación tanto en relación a su validez como en su aportación al progreso científico de la disciplina objeto de estudio. Los resultados de los trabajos de meta-análisis son especialmente útiles para la planificación de los valores del tamaño del efecto y para los análisis a priori de potencia de la prueba estadística.

During the decade of ninety in the XX century debates and opinions have been published (it can consult himself an exhaustive revision in Nickerson, 2000) that have questioned, in occasions difficultly, the application of the procedure of significance of the null hypothesis (Null Hypothesis Significance Testing, NHST) like instrument for the progress of the scientific knowledge. The critics to the procedure include from the false conceptions about the information that facilitates the procedure, causing erroneous interpretations of their results, until their inadequate use as half of obtaining of data of practical significance, sustantiva or clinic. The objective of our study is to not approach the problem still solved on the appropriate use of NHST inside the applied area of the confirmation of effects of the treatment and of the clinical change as well as to analyze the difficulty of the specification a priori of the value of the size of the effect and to analyze the different proposals that at the moment they intend to the procedure. Among this proposals the estimate of the size of the effect and the position of sizes of minimum effect and normative data that enrich the so much investigation design in relation to their validity highlights like in their contribution to the scientific progress of the discipline study object. The results of the goal-analysis works are specially useful for the planning of the values of the size of the effect and for the analyses a priori of power of the statistical test.

Dentro del área de la psicología aplicada interesa comprobar los efectos de los tratamientos donde lo importante no es únicamente mostrar que los tratamientos difieren o lo que es lo mismo si el tratamiento tuvo *algún* efecto sino que el interés se encuentra en señalar los beneficios de un tratamiento sobre otro, es decir, se desea conocer si el tratamiento tiene el efecto que el investigador plantea en su hipótesis científica (Fowler, 1985), apoyando la significación práctica o importante del cambio siguiendo el juicio del investigador. En general, en los estudios se plantea que el impacto del tratamiento producirá mayores cambios en la respuesta de los sujetos respecto al grupo de control no sometido a la intervención. El impacto del tratamiento es el denominado tamaño del efecto. Cuando se trata de analizar el denominado cambio clínico el efecto formulado en la hipótesis científica recoge el cambio hacia los valores de normalidad dentro del área psicológica sometida a intervención. En estos casos puede suceder que el cambio estadísticamente significativo no indique el verdadero valor terapéutico, prevaleciendo la importancia de la significación sustantiva, práctica o la denominada significación clínica entendida como la magnitud del cambio atribuida al tratamiento terapéutico

(Kendall, Flannery-Schroeder y Ford, 1999) que permite que el funcionamiento del sujeto pueda ser considerado normal. Ya sea como estimación de ciertos efectos o como cambio clínico, estamos hablando de analizar determinados efectos cuya cuantía está determinada por el modelo teórico y el campo psicológico en el que se trabaje.

De acuerdo con estos planteamientos de investigación, el procedimiento estadístico clásico basado en la hipótesis de nulidad de efectos con valor cero (hipótesis *nil* o vacía en términos de Cohen, 1994) no puede dar respuesta a las necesidades planteadas por hipótesis de nulidad con efectos distintos de cero (*non-nil* hipótesis) como pueden plantearse por ejemplo dentro del área de la valoración de programas o en el contexto del cambio clínico. Cuando se planifican diseños de investigación que están especialmente interesados en la comprobación de efectos del tratamiento, el cálculo del tamaño del efecto ayuda a la comprensión sustantiva de los resultados, complementando (sustituyendo, según algunos autores) el procedimiento tradicional de la comprobación de la significación estadística. No sólo por la necesidad de valorar el efecto del tratamiento desde el punto de vista del análisis del cambio práctico o clínico sino también debido a la interpretación incorrecta y uso inadecuado de las pruebas de significación estadística, la integración del análisis sustantivo de las diferencias detectadas entre los grupos debe ser incorporada al diseño de investigación. Analicemos brevemente los tres problemas más destacados que están vinculados con el uso e interpretación de los procedimientos clásicos de significación estadística.

Correspondencia: María Dolores Frías Navarro
Departamento de Metodología de las CC. del Comportamiento. Facultad de Psicología.
Universidad de Valencia
Avda. Blasco Ibáñez 21. 46010 Valencia.
E-mail: friasnav@uv.es

¿Qué significa utilizar el procedimiento de significación de la hipótesis nula? Implica poner en marcha el proceso de rechazar o mantener una hipótesis nula concreta basada en una serie de consideraciones teóricas planteadas a priori y en los valores p de probabilidad vinculados a la prueba estadística empleada con el objetivo de obtener conclusiones respecto a una hipótesis alternativa. Y desde este mismo planteamiento surge una de las críticas más difíciles de responder:

el procedimiento no especifica realizar predicciones sobre la hipótesis de investigación o científica sino que las predicciones se realizan sobre la hipótesis estadística de nulidad y será su rechazo la que dará crédito a la hipótesis de investigación que ni es especificada ni es comprobada con dicho procedimiento estadístico.

1. Interpretación incorrecta del valor de probabilidad p .

Se interpretan incorrectamente las probabilidades condicionales de resultados empíricos, $P(D/H_0)$, asociado al valor p de probabilidad, como probabilidades condicionales de hipótesis, $P(D/H_0)$ (Cohen, 1994), buscando respuestas que la técnica no puede ofrecer.

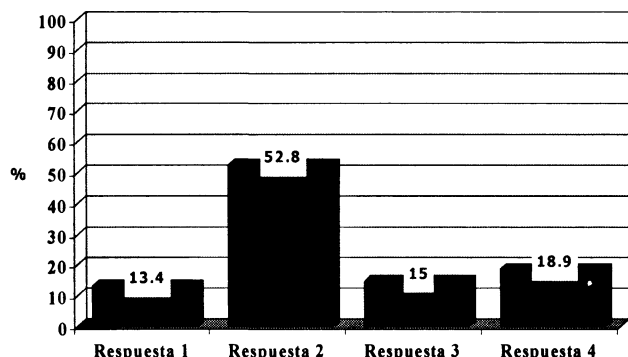
2. El valor p no informa del efecto del tratamiento.

La técnica no ofrece la probabilidad de que la hipótesis nula sea falsa, sólo facilita la probabilidad que tiene un resultado de ocurrir bajo la hipótesis nula pero no indica si dicho resultado ocurre o no ocurre bajo dicha hipótesis. El valor de p no tiene relación directa con la presencia o con la magnitud de efectos del tratamiento (Bracey, 1991; Cohen, 1994; Rosenthal, 1993). Además, la prueba estadística es función del tamaño del efecto y del tamaño de la muestra.

3. El valor p no informa de la significación sustantiva.

Las pruebas de significación estadística tampoco facilitan información sobre la magnitud o la importancia de un resultado (Grant, 1962; Rosenthal, 1993; Shaver, 1993), aspecto conocido como *significación práctica* (Kirk, 1996) o lo que Levy (1967) denominó *significación sustantiva* y que dentro del área de la psicología clínica se conoce como *significación clínica* (Jacobson, Follette, & Revenstorf, 1984) ¿Cuántos tipos de significación existen? Para realizar interpretaciones de efecto importante, práctico o clínico hay que recurrir al juicio del investigador sobre costos y beneficios tal y como se supone plantea en su hipótesis de investigación. Un juicio orientativo de significación práctica o cambio importante lo encontramos en los valores de Cohen (1990) sobre tamaño del efecto *pequeño* ($d=0.2$), *mediano* ($d=0.5$) y *grande* ($d=0.8$).

A un grupo de 132 alumnos de Psicología formados en las técnicas de significación estadística se les planteó un diseño experimental con dos condiciones experimentales de 150 sujetos cada una. Se les preguntó “si el valor de p obtenido hubiese sido de 0.04 ó quizás de 0.06 ó fuese de 0.0005 podríamos concluir respectivamente...” Como se observa en la gráfica, después de recibir formación de la técnica estadística durante dos años en diferentes asignaturas de análisis de datos, psicometría y diseños de investigación casi el 50% de los alumnos no interpretan correctamente un valor de probabilidad respecto al mágico 0.05 que por consenso establece el límite de la decisión dicotómica estadísticamente significativo/estadísticamente no significativo.



1. Cercano a la significación estadística ($p=0.04$), no del todo estadísticamente significativo ($p=0.06$) y muy significativo estadísticamente ($p=0.0005$).
2. Estadísticamente significativo los valores de $p=0.04$ y $p=0.0005$ y estadísticamente no significativo el valor de $p=0.06$.
3. Marginalmente estadísticamente significativo los valores de $p=0.04$ y $p=0.06$ y estadísticamente significativo el valor de $p=0.0005$.
4. No se puede responder a la pregunta.

Análisis del efecto de tratamiento

Diferentes métodos cuantitativos del cambio terapéutico sustantivo o de análisis de la eficacia del tratamiento permiten analizar si el valor obtenido cumple con el criterio estipulado por la hipótesis de investigación o si el nivel de funcionamiento del sujeto después de la terapia se encuentra dentro de los límites ‘normales’. Dentro de la primera orientación destacamos la técnica del meta-análisis donde la estimación del tamaño del efecto medio proporciona el índice de referencia (Glass, 1976; Hedges y Olkin, 1985; Hunter, Schmidt y Jackson, 1982; Rosenthal y Rubin, 1978) y el test de efectos mínimos (Murphy y Myers, 1999). Como técnicas representativas de puntuaciones dentro de la normalidad nos encontramos con el método de comparaciones normativas (Kendall y Grove, 1988) y el Índice del Cambio Fiable (Reliable Change Index, RCI, Jacobson, Follette & Revenstorf, 1984; Jacobson y Truax, 1991; Jacobson, Roberts, Berns y McGlinchey, 1999; Ogles, Lambert y Masters, 1996). En este trabajo nos centramos en las dos primeras técnicas.

Tamaño del efecto y estudios de meta-análisis

El cálculo del tamaño del efecto es una alternativa de análisis cuya información complementa los resultados del procedimiento de significación estadística (American Psychological Association, 1994, 2001; Abelson, 1995; Carver, 1978; 1993; Friedman, 1968; Loftus, 1993; Schafer, 1993; Thompson, 1994, 1996). Detallar el tamaño del efecto en los estudios individuales permite el desarrollo válido de trabajos de meta-análisis y además favorece que el investigador formule hipótesis cada vez más precisas. Sin embargo, conviene recordar que un tamaño del efecto grande no garantiza la importancia teórica o práctica del hallazgo que sólo puede estar respaldada por las hipótesis de investigación y la valoración o juicio del investigador. El problema es decidir en la fase de planificación del diseño qué valor del tamaño del efecto es el apropiado o qué valor se considera mínimo según el área concreta de estudio. Los estudios de meta-análisis ayudan al proceso de planificación proporcionando un tamaño del efecto medio. Dentro del campo clínico el problema es conocer qué valores determinan la normalidad de un constructo psicológico.

El uso del cálculo del tamaño del efecto va aumentando poco a poco especialmente provocado por la política editorial de las revistas que siguen las indicaciones de la A.P.A. pero aún no ha llegado a ser integrado en la interpretación sustantiva de los resultados de investigación, prevaleciendo los datos de significación estadística (Kirk, 1996; Thompson & Snyder, 1997, 1998; Vacha-Haase, Nilsson, Reetz, Lance & Thompson, 2000). Por ello en la quinta edición del manual del A.P.A. (2001) se señala que será *necesario* que siempre se incluya un índice de tamaño del efecto o de magnitud de la relación en la sección de *Resultados*.

“For the reader to fully understand the importance of your findings, it is almost always necessary to include some index of effect size or strength of relationship in your Results section. You can estimate the magnitude of effect or the strength of the relationship with a number of common effect size estimates... The general principle to be followed... is to provide the reader not only with information about statistical significance but also with enough information to assess the magnitude of the observed effect or relationship”. (pp. 25-26)

Robinson y Levin (1997) recomiendan utilizar el siguiente criterio. Primero se determina con NHST si el efecto es estadísticamente significativo. Si lo es, se informa del tamaño del efecto. Por su parte Fowler (1985) señala que se puede proceder a la inversa. Primero se comprueba si el tamaño del efecto tiene el valor mínimo que se plantea

en la hipótesis y si es así se pasa a contrastar su significación estadística.

Comprobación estadística de efectos mínimos

Se basa en determinar primeramente el tamaño del efecto mínimo que se considera sustantivo y comprobar posteriormente con NHST la hipótesis de que el efecto especificado es (o no es) estadísticamente significativo, eliminando la interpretación de efectos que no cumplen con el criterio mínimo sustantivo especificado en el planteamiento teórico del estudio dado que son considerados como efectos triviales o sin importancia práctica. Desde este planteamiento se ofrece una respuesta analítica a las críticas relacionadas con la especificación de tamaño de efecto cero, diferencia de medias igual a cero o correlación entre variables igual a cero que el modelo tradicional de contraste de hipótesis plantea con sus *nil hipótesis*. El problema con el que se encuentra el investigador es decidir a priori qué entiende por efecto mínimo cuantificando su valor y así poder determinar el tamaño del efecto que se considerará trivial en el análisis estadístico. Una posible solución puede ser utilizar los valores estándares de Cohen o mejor aún, utilizar los resultados de los trabajos de meta-análisis que ofrecen el tamaño del efecto medio dentro del área concreta objeto de estudio.

El planteamiento que orienta la perspectiva del test de efectos mínimos parte de la comprobación estadística de que la intervención tiene un efecto igual a o menor que un valor mínimo fijado a priori. Si, por ejemplo, el investigador decide que la psicoterapia que explica el 1% de la varianza tiene un efecto excesivamente pequeño como para que sea útil o práctico su uso, entonces su tarea consistirá en desarrollar un test estadístico que determine cuán grande ha de ser un tamaño del efecto en una muestra particular para que explique menos del 1% de la varianza. Ese test será una razón F , determinando el valor necesario para rechazar la hipótesis que el efecto en la población es igual o menor que el fijado.

El test de efectos mínimos se basa en la distribución F no centrada y para desarrollar las tablas con las que testar la hipótesis nula de efectos mínimos, el investigador deberá a) desarrollar la definición operacional de tamaño mínimo; b) calcular el tamaño del parámetro de no centralidad para tal efecto mínimo y c) tabular la correspondiente distribución de F no centrada. En este planteamiento se nos presentan dos problemas. Primero *cómo* definir de manera sensata el “efecto mínimo”. Segundo, *cómo* generar las tablas de F no centrada.

La definición de *qué es el efecto mínimo* requiere un juicio de valor que escapa al análisis estadístico y es previo. Es un problema dependiente de la teoría, de la praxis o del consenso. Otras veces la determinación del tamaño de efecto mínimo puede que exija un análisis de utilidad o de costos y de relevancia social o clínica o cualquier otro criterio que se estime oportuno. Siempre y en cualquier caso, habrá que justificarlo y definirlo y no dar por supuesto como un “prejuicio” dogmáticamente establecido que ese valor ha de ser siempre el valor cero. Por lo tanto, dependerá del área a investigar y de otras concomitancias y también de criterios estadísticos porque conforme se defina un tamaño del efecto mínimo mayor entonces decrecerá la potencia del modelo general lineal para detectar tal efecto, es decir aumenta por tanto la probabilidad de cometer error Tipo II (no detectar un efecto que realmente sí existe). Por lo tanto, antes de llevar a cabo la investigación habrá que delimitar cuál es la relación de riesgo que se acepta respecto de los errores estadísticos de Tipo I (detectar un efecto que realmente no existe) y de Tipo II (no detectar un efecto que realmente sí existe). La relación correcta se puede suponer que es de uno a cuatro como afirma Cohen (1988) o se puede suponer mayor o menor según las circunstancias concretas como por ejemplo el área de investigación.

La generación de las tablas de F no centrada es cada día más fácil gracias a las estimaciones que se pueden realizar con las funciones que aparecen en programas como el SPSS o el Excel. También se pueden utilizar programas específicos que computan las probabilidades (Narula & Weistroffer, 1968). Una aproximación basada en la distribución central de la prueba F (Horton, 1978; Murphy & Myers, 1999; Patnaik,

1949; Tiku & Yip, 1978) se puede realizar multiplicando el valor crítico de la F central, con grados de libertad g y v_2 , por el valor k . De este modo,

$$Pr(F_{v_1, v_2, \lambda} > c \cdot k) \text{ es aproximadamente igual a } Pr(F_{g, v_2} > c)$$

siendo Pr la probabilidad, $F_{v_1, v_2, \lambda}$ es un valor de la distribución F no central, F_{g, v_2} es un valor de la distribución F central y c es una constante. El cómputo de g y k es el siguiente:

$$g = \frac{(v_1 + \lambda)^2}{(v_1 + 2\lambda)} \quad k = \frac{(v_1 + \lambda)}{v_1}$$

El valor de lambda puede ser estimado por medio del porcentaje de varianza explicada, PV, y los grados de libertad del error, v_2 :

$$\lambda = \frac{(v_2 \cdot PV)}{1 - PV}$$

Por ejemplo en las tablas 1 y 2 se detallan los valores de la F no central obtenidos con el paquete estadístico SPSS y el Excel. Las tablas proporcionan los valores críticos necesarios para rechazar una hipótesis nula de efectos mínimos de varianza explicada del 1% o menos que en términos del estadístico d de Cohen (1988) estamos hablando de tamaños del efecto *pequeño*. Del mismo modo podemos calcular los valores para efectos mínimos de varianza explicada del 5% vinculándose con valores de d de efectos entre *pequeños* y *medianos*. Y así para cualquier valor que especifique el criterio de la hipótesis de investigación.

Conclusiones

Un análisis de las aportaciones de los investigadores dedicados al tema de la significación práctica o sustantiva conduce a una conclusión aceptada unánimemente: *se ha abusado de las pruebas de significación estadística aplicándose a todas las situaciones*. Muchas creencias estadísticas deben ser modificadas. El debate acerca de NHST está sin resolver, pero de lo que no cabe duda es que el juicio del investigador es indispensable en la interpretación de los resultados impidiendo un uso mecánico de NHST.

Mientras la polémica sigue su curso durante más de 70 años, el procedimiento de significación de la hipótesis nula sigue siendo la técnica más utilizada por los investigadores y como ya apuntaba Carver (1978), todas las críticas hacia las pruebas de significación estadística han tenido poco efecto. Como consecuencia, es necesario realizar una reflexión sobre su uso adecuado y sobre la formación académica que el futuro profesional recibe en metodología.

Actualmente la cuantificación del impacto del tratamiento, entendido como tamaño del efecto que provoca cambio práctico o formulado como cambio clínico, ha cobrado un papel destacado entre las nuevas formulaciones del análisis estadístico. La técnica del análisis de datos debe adaptarse a las formulaciones sustantivas de cambio práctico o clínico del diseño de investigación. La estimación del tamaño del efecto junto con los datos de significación estadística permiten una comprensión sustantiva de los resultados obtenidos, evitando interpretaciones difíciles de replicar debido a la falta de consistencia del impacto del tratamiento.

La búsqueda del análisis válido del cambio requiere computar y analizar el cambio sustantivo, práctico o clínico. Poder contrastar hipótesis nulas con efecto distinto a cero (hipótesis “non-nil nulls”) enriquecerá nuestras teorías psicológicas, avanzado el conocimiento y eliminando ciertas polémicas sobre la trivialidad de testar hipótesis con efecto cero, evitando al mismo tiempo la interpretación de resultados estadísticamente significativos sin importancia práctica

Si la objetividad de su planteamiento y su facilidad de cálculo son algunas de las razones del por qué los investigadores siguen adheridos al procedimiento de significación de la hipótesis nula entonces el planteamiento del *análisis de los efectos mínimos* cuenta con esas mismas ventajas, además se evitaría que un resultado no

- Horton, R. L. (1978). *The general linear model: Data analysis in the social and behavioral sciences*. New York, NY: McGraw-Hill.
- Hunter, J. E., Schmidt, F. L. y Jackson, G. B. (1982). *Meta-analysis: cumulating findings across research*. Beverly Hills, CA: Sage.
- Jacobson, N.S. y Truax, P. (1991). Clinical significance: A statistical approach to defining meaningful change in psychotherapy research. *Journal of Consulting and Clinical Psychology, 59*, 12-19.
- Jacobson, N.S., Follette, W.C. y Revenstorf, D. (1984). Psychotherapy outcome research: Methods for reporting variability and evaluating clinical significance. *Journal of Consulting and Clinical Psychology, 15*, 336-352.
- Jacobson, N.S., Roberts, L.J., Berns, S.B. y McGlinchey, J.B. (1999). Methods for defining and determining the clinical significance of treatment effects: Description, application and alternatives. *Journal of Consulting and Clinical Psychology, 67*(3), 300-307.
- Kendall, P.C. y Grove, W.M. (1988). Normative comparisons in therapy outcome. *Behavioral Assessment, 10*, 147-158.
- Kendall, P.C., Flannery-Schroeder, E.C. Ford, J.D. (1999). Therapy outcome research methods. En P.C. Kendall, J.N. Butcher y G.N. Holmbeck (eds.). *Handbook of research methods in clinical psychology*. New York: Wiley and Sons.
- Kirk, R. E. (1996). Practical significance: a concept whose time has come. *Educational and Psychological Measurement, 56*, 746-759.
- Levy, P. (1967). Substantive significance of significant differences between two groups. *Psychological Bulletin, 67*, 37-40.
- Loftus, G.R. (1993). Editorial comment. *Memory & Cognition, 21*, 1-3.
- Murphy, K. R. & Myers, B. (1999). Testing the hypothesis that treatments have negligible effects: Minimum-effect tests in the general linear model. *Journal of Applied Psychology, 84*, 234-248.
- Narula, S., & Weistroffer, H. (1968). Computation of probability and noncentrality parameter of a noncentral F distribution. *Communications in Statistics B, 15*, 871-878.
- Ogles, B.M., Lambert, M.J. y Masters, K.S. (1996). *Assessing outcome in clinical practice*. Boston: Allyn and Bacon.
- Patnaik, P.B. (1949). The non-central χ^2 and F-distributions and their applications. *Biometrika, 36*, 202-232.
- Robinson, D., y Levin, J. (1997). Reflections on statistical and substantive significance with a slice of replication. *Educational Researcher, 26*, 21-26.
- Rosenthal, R y Rubin, D. (1978). Interpersonal expectancy effects: the first 345 studies. *Statistical Science, 3*, 120-125.
- Rosenthal, R. (1993). Cumulative evidence. En G. Keren & C. Lewis (Eds.). *A handbook for data analysis in the behavioral sciences: Methodological issues*. Hillsdale, NJ: Erlbaum.
- Shaver, W.D. (1993). Interpreting statistical significance and nonsignificance. *Journal of Experimental Education, 61*, 383-387.
- Thompson, B. & Snyder, P.A. (1997). Statistical significance testing practices in the Journal of Experimental Education. *Journal of Experimental Education, 66*, 75-83.
- Thompson, B. & Snyder, P.A. (1998). Statistical significance and reliability analyses in recent JCD research articles. *Journal of Counseling and Development, 76*, 436-441.
- Thompson, B. (1996). AERA editorial policies regarding statistical significance testing: Three suggested reforms. *Educational Researcher, 26*, 29-32.
- Thompson, B. (1996). Guidelines for authors. *Educational and Psychological Measurement, 54*, 837-847.
- Tiku, M.L. & Yip, D.Y.N. (1978). A four-moment approximation based on the F distribution. *Australian Journal of Statistics, 20*, 257-261.
- Vacha-Haase, T., Nilsson, J.E., Reetz, D.R., Lance, T.S. & Thompson, B. (2000). Reporting practices and APA editorial policies regarding statistical significance and effect size. *Theory & Psychology, 10*, 413-425.