

Bogotá D.C  
Noviembre 2024  
ISSN: 2590 - 4663  
Publicación Trimestral

# SABER AL DETALLE

**EDICIÓN 16**

**¿Qué se entiende por confiabilidad  
y validez en el contexto de  
la medición con instrumentos?**



**Presidente de la República**  
Gustavo Francisco Petro Urrego

**Ministro de Educación Nacional**  
Jose Daniel Rojas Medellín

**Elaboración del documento**  
Nila Fernanda Amaya Melo  
Mishell Marcela Ramos de la Hoz  
Esteban Nicolás Arias Cubillos

**Revisado por**  
Ana María Cruz Pacheco  
Joan Gabriel Bofill Barrera  
Ricardo Macías Bohórquez

**Diseño y diagramación**  
Andrea del Pilar López Pulido

Bogotá D.C., Noviembre 2024

**Todos los derechos de autor reservados ©.**

**Directora General**  
Elizabeth Blandón Bermúdez

**Secretaria General**  
Brahiam Daniel Montoya Zuleta (E)

**Director de Evaluación**  
Rafael Eduardo Benjumea Hoyos

**Subdirector de Estadísticas**  
Cristian Fabian Montaña Rincón

**Subdirectora de Análisis y Divulgación**  
Alejandra Neira Aroca

**Directora de Producción y Operaciones**  
Luz Patricia Loaiza Cruz

**Subdirector de Producción de Instrumentos**  
Gustavo Andrés Monsalve Londoño

**Director de Tecnología e Información**  
Luis Rodrigo Cadavid Durán

**Subdirector de información**  
Diego Mauricio Salas Ramírez

**Jefe Oficina Asesora de Comunicaciones y Mercadeo**  
Byron Andrés Vélez Valdés

**Jefa Oficina Asesora de Gestión de Proyectos de Investigación**  
Jennyffer Paola Guío Veloza



## Términos y condiciones de uso para las publicaciones y obras que son propiedad del Icfes

El Instituto Colombiano para la Evaluación de la Educación (Icfes) pone a disposición de la comunidad educativa, y del público en general, de forma gratuita y libre de cualquier cargo, un conjunto de publicaciones disponibles en su portal web [www.icfes.gov.co](http://www.icfes.gov.co). Estos materiales y documentos están normados por la presente política y se encuentran protegidos por derechos de propiedad intelectual y derechos de autor a favor del Icfes. Si tiene conocimiento de alguna utilización contraria a lo establecido en estas condiciones de uso, por favor infórmenos al correo [prensaicfes@icfes.gov.co](mailto:prensaicfes@icfes.gov.co).

Queda prohibido el uso o publicación total o parcial de este material con fines de lucro. Únicamente está autorizado su uso para fines académicos e investigativos. Ninguna persona natural o jurídica, nacional o internacional, podrá vender, distribuir, alquilar reproducir, transformar<sup>1</sup>, promocionar o realizar acción alguna con la cual se lucre directa o indirectamente con este material. Esta publicación cuenta con el registro ISBN (International Standard Book Number, o Número Normalizado Internacional para Libros), que facilita la identificación, no sólo de cada título, sino, también de la autoría, la edición, el editor y el país dónde se edita.

En todo caso cuando se haga uso parcial o total de los contenidos de esta publicación, el usuario deberá consignar o hacer referencia a los créditos institucionales del Icfes, respetando los derechos de cita.

En otras palabras, se podrá hacer uso de esta publicación si dicho uso se contempla en los fines aquí previstos. Es posible, entonces, transcribir pasajes del texto si se cita siempre la fuente de autor. Por supuesto, estas citas no deberían ser excesivas ni frecuentes para que, así, no se considere una reproducción simulada y sustancial que redunde en perjuicio del Icfes.

Asimismo, los logotipos institucionales son marcas registradas y de propiedad exclusiva del Icfes. Por tanto, cuando su uso pueda causar confusión, los terceros no podrán usar las marcas de propiedad del Icfes con signos idénticos o similares respecto a cualquier producto o servicio prestado por esta entidad. En todo caso queda prohibido su uso sin previa autorización expresa por parte del Icfes. La infracción de estos derechos se perseguirá civil y penalmente (en caso de que sea necesario) de acuerdo con las leyes nacionales y tratados internacionales aplicables.

***El Icfes realizará cambios o revisiones periódicas a los presentes términos de uso y los actualizará en esta publicación.***

---

<sup>1</sup> La transformación es la modificación de la obra a través de la creación de adaptaciones, traducciones, compilaciones, actualizaciones, revisiones y, en general, cualquier modificación que se pueda realizar, haciendo que la nueva obra resultante se constituya en una obra derivada protegida por el derecho de autor, con la única diferencia, respecto de las obras originales, que aquellas requieren, para su realización, de la autorización expresa del autor o propietario, para adaptar, traducir, compilar, etc. En este caso, el Icfes prohíbe la transformación de esta publicación. Términos y condiciones de uso para las publicaciones y obras que son propiedad del Icfes.

## ¿Qué se entiende por confiabilidad y validez en el contexto de la medición con instrumentos?

En todas las pruebas que desarrolla el Icfes se emplean procedimientos para explorar sus propiedades psicométricas. Las propiedades psicométricas de un instrumento brindan información respecto a la calidad de la medición de los constructos que evalúa y su consistencia. Estos elementos se asocian con la confiabilidad de las pruebas y con la validez de las interpretaciones que se realizan a partir de ellas. Es importante explorar estas propiedades, considerando que, a partir de los resultados en la prueba, se toman decisiones que inciden en la calidad de la educación del país. En este sentido, el objetivo de este boletín de Saber al Detalle es presentar, de forma introductoria y general, los conceptos asociados a validez y confiabilidad<sup>1</sup>. Además, se espera mostrar las diferencias a nivel teórico de estos dos conceptos.

### ¿Qué es la confiabilidad?

La confiabilidad hace referencia a la precisión o consistencia de un instrumento en el proceso de medición. Como señalan Manterola et al., “un instrumento es confiable, preciso o reproducible, cuando las mediciones realizadas con él, generan los mismos resultados en diferentes momentos, escenarios y poblaciones si se aplican en las mismas condiciones” (2018, p. 680). Para ilustrar lo anterior, se puede tener en cuenta el ejemplo de una báscula que se usa para conocer el peso de una manzana en tres momentos diferentes. Si en cada ocasión se obtiene una medida distinta de la misma manzana, se puede considerar que la báscula presenta problemas de confiabilidad dado que esta medida no es consistente y, por tanto, presenta error.

En la medición con instrumentos se pueden presentar distintas clases de error; sin embargo, al hablar

de confiabilidad, cobra especial relevancia el error aleatorio. Este tipo de error hace referencia a eventos o situaciones que los evaluadores desconocen que pueden ocurrir y afectan la medición (Geisinger, 2013, p. 21; APA, AERA y NCME, 2018). En este caso, no se obtiene una misma medición cada vez que se aplica el instrumento, como se ilustra en el ejemplo de la báscula.

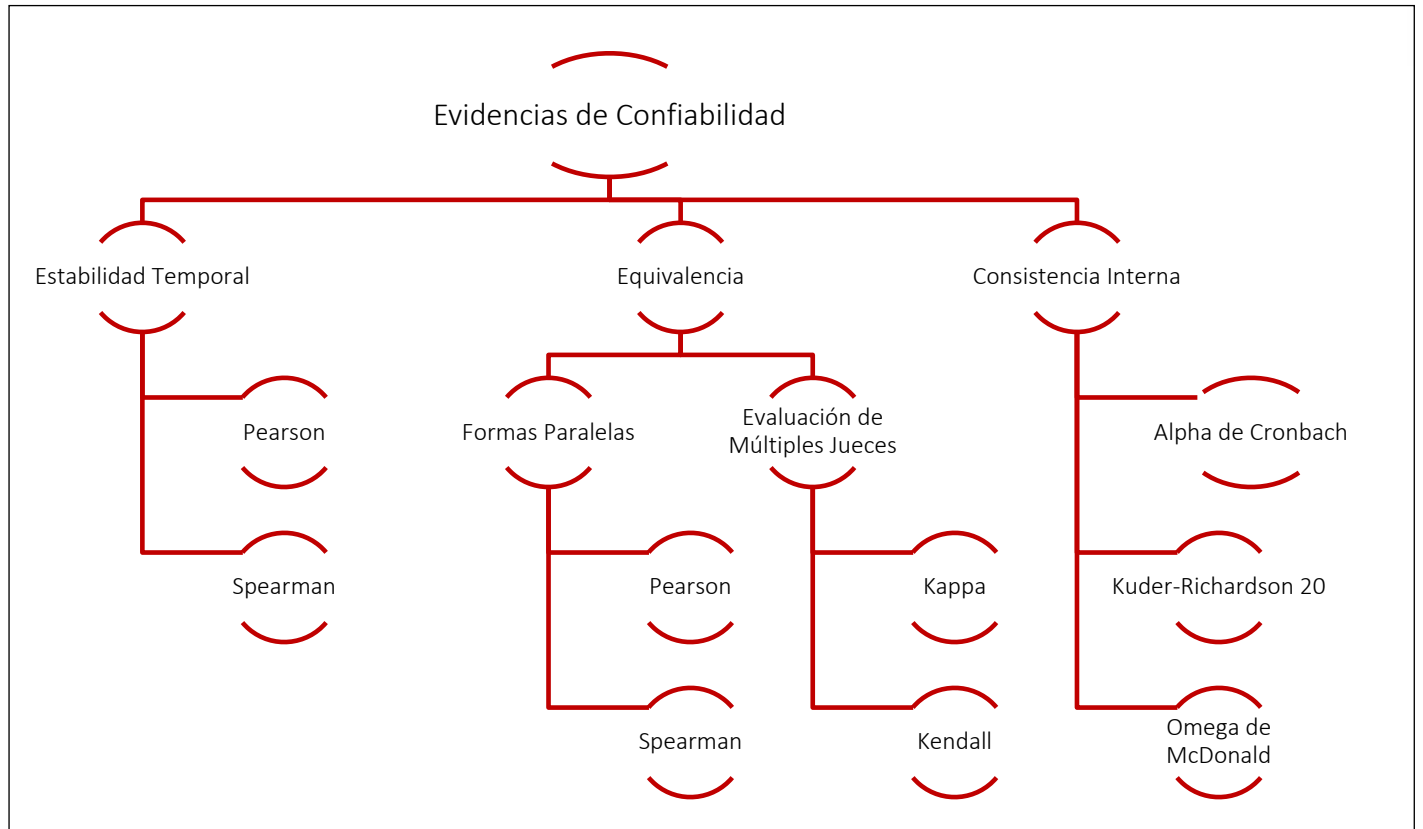
En resumen, la confiabilidad de un instrumento de medición, o test, es esencial para asegurar que las mediciones sean precisas y consistentes, permitiendo que los resultados sean reproducibles en diferentes momentos y contextos bajo las mismas condiciones. Al evaluar la confiabilidad, se busca minimizar los errores que puedan afectar las mediciones, garantizando así que los instrumentos utilizados proporcionen datos fiables. De este modo, se asegura la calidad y la precisión en la medición de los constructos evaluados.

### ¿Cómo se puede estimar la confiabilidad?

La implementación de procedimientos para estimar la confiabilidad depende de la fuente de error que pueda ser relevante según la naturaleza del objeto de medición y las características del instrumento. Se puede estimar la confiabilidad para conocer la estabilidad temporal de una prueba, la equivalencia entre las mediciones o la consistencia interna de los test. En la Figura 1 se presenta un esquema de los procedimientos que se pueden emplear para estimar la confiabilidad dependiendo de los métodos pertinentes para mitigar la fuente de error que interesa para la medición.

1. Para profundizar en los temas de confiabilidad y validez, le invitamos a consultar las siguientes referencias: American Educational Research Association -AERA, American Psychological Association - APA, & National Council on Measurement in Education -NCME (2018). Estándares para pruebas educativas y psicológicas. American Educational Research Association. Muñiz, J. (2018). Introducción a la Psicometría: Teoría Clásica y TRI. Pirámide. Martínez, M., Hernández, V., y Hernández, M. (2014). Psicometría. Alianza Editorial.

**Figura 1**  
**Mapa de decisiones de los diferentes índices de confiabilidad.**



Nota: Elaboración propia.

Frente a la estimación de la confiabilidad para examinar la estabilidad temporal de un instrumento, se busca cuantificar el grado en que varían las puntuaciones de las personas al aplicar la misma prueba en dos momentos diferentes (Argibay, 2006). La pertinencia de emplear este coeficiente va a depender de la naturaleza del objeto de evaluación, pues en constructos como la actitud es esperable una amplia variabilidad en periodos cortos, por lo que no sería apropiado estimar la confiabilidad mediante procedimientos relacionados con la estabilidad temporal. No obstante, estos procedimientos pueden ser apropiados en mediciones donde se esperen cambios leves a través del tiempo respecto al constructo evaluado, como por ejemplo en pruebas de personalidad.

La estabilidad temporal se mide a través de méto-

dos de test-retest, que buscan examinar la asociación entre los puntajes obtenidos en cada aplicación. Ejemplos de estos coeficientes son: la Correlación de Pearson, recomendada para la correlación de variables continuas con distribución normal bivariada, y la Correlación de Spearman, una medida no paramétrica que utiliza los rangos o posiciones de los datos en lugar de los datos que resulta conveniente en los casos en que no se puede cumplir con estos supuestos distribucionales (Ortiz Pinilla y Ortiz, 2021).

Respecto a la equivalencia, se suele hablar de dos tipos principalmente: la confiabilidad de la medición asociada a dos formas paralelas del instrumento y la asociada a la evaluación de múltiples jueces. La equivalencia de formas implica comparar dos versiones distintas del mismo instrumento aplicadas simultáneamente a los mismos sujetos, variando el orden de presentación para controlar posibles

sesgos de orden (Argibay, 2006). En este caso, se busca establecer si ambas formas proporcionan evaluaciones consistentes, lo que se evidencia mediante correlaciones elevadas de los puntajes obtenidos, al emplear índices como los mencionados en el párrafo anterior.

Por otro lado, la equivalencia entre evaluadores u observadores se aplica en instrumentos donde múltiples jueces deben calificar o puntuar la conducta o rendimiento de la persona (Argibay, 2006). En este caso, se utilizan índices de concordancia, como el coeficiente Kappa (Cohen, 1960) o el Kendall (Kendall, 1938).

Respecto a los coeficientes asociados con la estimación de la consistencia interna, estos se basan en las interacciones entre puntajes, producto de los ítems individuales o subconjuntos de ellos (Argibay, 2006). Uno de los índices más utilizados para estimar la consistencia interna en pruebas es el Alpha de Cronbach (Cronbach, 1951). Por ejemplo, en el Icfes se utiliza este índice para las pruebas que presentan ítems politómicos. Este coeficiente se basa en las varianzas y covarianzas de las respuestas a los ítems dadas por las personas evaluadas (Vila, 2010) y está dado por la siguiente fórmula:

$$KR_{20} = \frac{n}{n-1} \left( 1 - \frac{\sum_{j=1}^n p_j q_j}{S_x^2} \right)$$

Donde:

$n$  representa el número de ítems  
 $\sum S_j^2$  representa la suma de las varianzas de los ítems que conforman el test ( $n$  ítems), y  
 $S_x^2$  varianza de las puntuaciones en el test

El alfa de Cronbach resume la correlación entre los diferentes ítems en una escala entre cero (0) y uno (1), donde valores cercanos a uno (1) indican que los ítems incluidos en la prueba son consistentes en la medición, mientras que un valor cercano a 0 sugiere que no están relacionados de manera coherente.

Es importante destacar que este coeficiente no solo ofrece una medida de la confiabilidad de la prueba en su conjunto, sino que también permite identificar la contribución de cada ítem a la consistencia interna global. De esta manera, se pueden identificar ítems problemáticos que podrían afectar la fiabilidad de la prueba en su totalidad, facilitando así la mejora y refinamiento del instrumento de medición.

Otro índice que también es ampliamente utilizado es el coeficiente Kuder-Richardson (KR-20), que es una variación del coeficiente Alpha de Cronbach (Kuder y Richardson, 1937). El coeficiente KR-20 mantiene el mismo procedimiento del Alpha de Cronbach, pero se aplica para ítems dicotómicos, los cuales tienen una única respuesta correcta que se puntúa con un uno (1) cuando se presentan aciertos y con cero (0) cuando se selecciona una opción distinta a la correcta.

Para calcular este indicador, se analiza la correlación entre cada ítem y la puntuación total de la prueba, así como la correlación entre cada par de ítems, de acuerdo con la siguiente ecuación (Muñiz, 2018):

$$KR_{20} = \frac{n}{n-1} \left( 1 - \frac{\sum_{j=1}^n p_j q_j}{S_x^2} \right)$$

Donde:

$n$  representa el número de ítems  
 $p_j$  representa la proporción de respuestas correctas en el ítem  
 $q_j$  representa la proporción de respuestas incorrectas en el ítem, y  
 $S_x^2$  representa la varianza de las puntuaciones globales

El resultado final en el coeficiente KR-20 es un valor entre cero (0) y uno (1) que refleja cuán interrelacionados están los ítems y su interpretación es equivalente a la descrita para el Alfa de Cronbach.

Existen otros métodos menos conocidos para estimar la consistencia interna de los instrumentos como, por ejemplo, el índice Omega de McDonald

(McDonald, 1999), que puede ser empleado para el análisis de instrumentos conformados tanto por ítems politómicos como ítems dicotómicos. Teniendo en cuenta las diferentes formas en las que se puede estimar la confiabilidad, se espera que la información presentada en este apartado del boletín brinde herramientas que contribuyan a la exploración de estos métodos, teniendo en cuenta las características del constructo a evaluar.

## ¿Qué es validez?

De acuerdo con APA, AERA y NCME “la validez se refiere al grado en que la evidencia y la teoría respaldan las interpretaciones de los puntajes de una prueba para usos propuestos de las pruebas” (p. 11, 2018). Esta definición resalta la relevancia de considerar el contexto en el que se usan los instrumentos para garantizar la validez de las interpretaciones realizadas. Esto implica que un instrumento puede contar con evidencias de validez que respalden el constructo que pretende evaluar; sin embargo, si se utiliza para un propósito diferente al previsto, estas evidencias de validez pueden perder su significado.

Las evidencias de validez de una prueba pueden verse afectadas por el error sistemático, el cual se refiere a la obtención de resultados inexactos que se presentan de forma consistente. El error sistemático representa factores irrelevantes para la medición del constructo que comprometen la validez del instrumento (APA, AERA y NCME, 2018). Por ejemplo, si se empleara una báscula sin tarar para la medición del peso de una manzana, se obtendrían resultados sistemáticamente mayores o menores que el peso real.

Dado lo anterior, es importante contar con evidencias que soporten las interpretaciones realizadas a partir de los resultados obtenidos en una medición, ya que

así se puede garantizar que las decisiones tomadas a partir de las mediciones sean coherentes con el propósito de la prueba.

## ¿Qué fuentes de evidencias de validez existen?

De acuerdo con la versión actual de los estándares de la AERA, APA y NCME (2018), las evidencias de validez se categorizan teniendo en cuenta la fuente de información de la que se derivan. Estas se pueden clasificar en los siguientes tipos: basadas en el contenido de la prueba, los procesos de respuesta, la estructura interna, la relación con otras variables y las consecuencias del uso de la prueba.

Las evidencias centradas en el contenido de la prueba tienen como objetivo respaldar su representatividad con respecto al dominio que esta busca medir (AERA, APA y NCME, 2018). Estas evidencias pueden originarse tanto de análisis lógicos, como el juicio de expertos cuya efectividad se puede medir a partir de índices estadísticos como el Índice de Validez de Contenido (Lawshe, 1975), el coeficiente V de Aiken (Aiken, 1985), el Índice AC1 de Gwet (Gwet, 2014) o el Alpha de Krippendorff (Krippendorff, 2019); o ser de naturaleza empírica, al evaluar el funcionamiento de los ítems en la población<sup>2</sup>. Estos elementos ofrecen respaldo a aspectos cruciales, tales como: los temas evaluados por el instrumento, la redacción de los ítems, su formato, la metodología de administración y el procedimiento definido para la estimación de puntajes.

La evidencia derivada del proceso de respuesta evalúa la idoneidad de las respuestas de los individuos y el proceso cognitivo implicado en la resolución de los ítems (Padilla y Benítez, 2014). Estas evidencias suelen obtenerse mediante el análisis de las respuestas individuales de personas que forman parte

<sup>2</sup> Para mayor información respecto al proceso que realiza el Icfes en el análisis de ítems, puede consultar el boletín 9 de Saber al Detalle en el siguiente enlace:

<https://www.icfes.gov.co/documents/39286/2231027/Edici%C3%B3n+9+-+C%C3%B3mo+se+analizan+los+%C3%ADtems+de+las+pruebas+Saber.pdf/d08d9002-59d6-e71c-f066-e6cb41596b24?version=2.0&t=1685369830376>

de la población objetivo; de modo que se examinan las estrategias que emplean para llegar a las respuestas (AERA, APA y NCME, 2018).

La evidencia basada en la estructura interna evalúa en qué medida las relaciones entre los ítems de la prueba y sus componentes se alinean con el constructo subyacente que respalda las interpretaciones propuestas de los puntajes (AERA, APA y NCME, 2018). En esencia, estas evidencias respaldan la coherencia entre la interrelación de los ítems y la medida subyacente (Rios y Wells, 2014). Estas evidencias suelen obtenerse mediante técnicas como el análisis factorial<sup>3</sup>, que examina cómo los ítems se agrupan y relacionan dentro de la estructura global de la prueba.

Las evidencias basadas en la relación con otras variables analizan la correlación entre los puntajes de la prueba y variables externas para determinar en qué medida estas relaciones concuerdan con el constructo evaluado (AERA, APA y NCME, 2018). En este proceso, se busca identificar: a) relaciones convergentes, al evaluar la conexión entre las puntuaciones de la prueba y otras que miden el mismo constructo, b) relaciones discriminantes, al relacionar las puntuaciones de la prueba con medidas de constructos diferentes y c) relaciones prueba-criterio, al determinar la capacidad de la prueba para prever el comportamiento de un atributo que, aunque operativamente distinto, está teóricamente relacionado con el constructo que se pretende medir.

Las evidencias de validez asociadas a las relaciones convergentes se identifican cuando las correlaciones entre las mediciones de un rasgo, utilizando diferentes métodos, resultan significativas. Por otro lado, las evidencias de validez asociadas a las relaciones discriminantes se manifiestan cuando las correlaciones entre las mediciones del mismo rasgo que se obtienen utilizando diferentes métodos, son notablemente mayores que las correlaciones entre las mediciones de diferentes rasgos utilizando el mismo método. Las relaciones convergente y discriminante se pueden

extraer de los datos recopilados mediante la matriz multirrasgo-multimétodo, un conjunto de correlaciones que incluye múltiples rasgos evaluados por diversos métodos (Muñiz, 2018).

En cuanto a las relaciones prueba-criterio, según Muñiz (2018), tradicionalmente se ha utilizado la correlación como índice de validez, adaptando la fórmula según la fiabilidad del criterio y del instrumento y del número de variables a considerar. No obstante, en los casos en que el propósito del test es la predicción de un criterio, se han empleado otros tipos de estadísticos, como la regresión simple y la regresión múltiple; además de estadísticos que permiten evaluar el grado de certeza de las decisiones que el instrumento facilita, con índices como el de Sensibilidad, que cuantifica la precisión de la escala en la detección de personas con presencia del criterio, o el de Especificidad, que considera la capacidad de la escala para clasificar a las personas con ausencia del criterio.

Finalmente, la evidencia basada en las consecuencias de la medida evalúa los efectos prácticos y sociales de la prueba, considerando las implicaciones asociadas con su aplicación (Pan, 2009). Este tipo de evidencia se basa en la premisa de que el diseño de la prueba contempla su aplicación en condiciones específicas y con un propósito definido. En consecuencia, la recopilación de esta evidencia busca determinar si los beneficios esperados derivados de la interpretación y el uso de las puntuaciones generadas por el instrumento se han materializado completamente (AERA, APA y NCME, 2018).

## Finalmente, ¿En qué se diferencia la confiabilidad de la validez?

En el contexto de la medición de instrumentos se tiende a confundir los conceptos de confiabilidad y validez, sin embargo, estos presentan diferencias

<sup>3</sup> Para profundizar en el análisis factorial, puede consultar el Boletín 10 de Saber al Detalle en el siguiente enlace: <https://www.icfes.gov.co/documents/39286/2231027/Saber+al+detalle+Ed+10.pdf/d8d21b5c-fea9-2e19-9742-cbeba1e0601d?version=2.0&t=1691760390294>

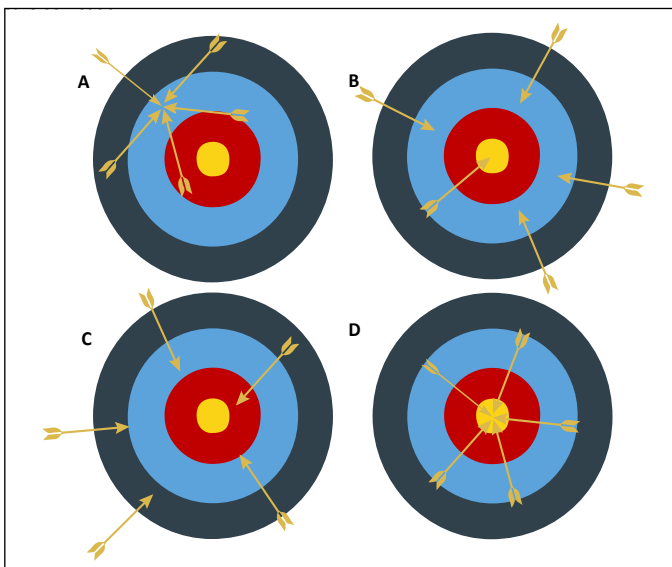
las cuales se deben tener presentes al momento de analizar la calidad de las pruebas. La confiabilidad, como se mencionó anteriormente, brinda información respecto a la precisión de los instrumentos, mientras que la validez se asocia con las evidencias que se recogen para respaldar la interpretación de los resultados de las pruebas que se aplican, teniendo en cuenta el propósito del instrumento.

En la Figura 2 se representan algunos posibles escenarios de validez que se pueden encontrar en una prueba. En el módulo A, todas las mediciones son consistentes, pero se alejan del valor real del nivel de habilidad, lo que representa una prueba con una alta fiabilidad, pero baja validez. En el módulo B, algunas mediciones aciertan en el nivel de habilidad de la persona, pero esto no sucede consistentemente, lo que representa una prueba con un grado de validez, pero baja fiabilidad. En el módulo C se observa que ninguna medición acierta sobre el nivel de habilidad de las personas, lo que representa una prueba con baja fiabilidad y validez. Por último, el módulo D representa una prueba que acierta consistentemente en el nivel de habilidad de las personas, lo que representa una prueba con alta fiabilidad y validez.

En conclusión, si bien la validez y la confiabilidad son conceptos diferentes que permiten conocer las propiedades de un instrumento, en conjunto logran garantizar la calidad de las mediciones realizadas. Esto es relevante ya que, a partir de los resultados y las interpretaciones de las mediciones, se suelen tomar decisiones importantes que impactan la vida de las personas evaluadas. Por ello, es relevante desarrollar procedimientos que permitan conocer y analizar estas propiedades para garantizar un instrumento riguroso y preciso.

Desde el Icfes siempre se busca que, tanto en las pruebas de Estado como en los proyectos, se implementen procesos que permitan recoger información sobre la confiabilidad y las evidencias de validez de las mediciones que se realizan. Para esto, se utiliza una amplia gama de estadísticas que se complementan con revisiones cualitativas de las pruebas para obtener un panorama amplio de las características psicométricas de los instrumentos.

**Figura 2**  
**Ilustración sobre conceptos de validez y confiabilidad.**



Nota: Figura tomada y adaptada de Manterola et al. 2018

## Referencias

- Aiken, L. (1985). Three coefficients for analyzing the reliability and validity of ratings. *Educational and Psychological Measurement*, 45(1), 131-142.
- American Educational Research Association, American Psychological Association y National Council on Measurement in Education. (2018). *Standards for educational and psychological testing*. American Educational Research Association.
- Argibay, J. C. (2006). Técnicas psicométricas: Cuestiones de validez y confiabilidad. *Subjetividad y Procesos Cognitivos*, 8, 15-33.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1), 37-46. <https://doi.org/10.1177/001316446002000104>
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, 16(3), 297-334.
- Geisinger, K. (2013). Reliability. En Geisinger, K. (Ed.) *APA Handbook of testing and assesment in psychology* (pp. 21-42). American Psychological Association.
- Gwet, K. L. (2014). *Handbook of inter-rater reliability: The definitive guide to measuring the extent of agreement among raters*. Advanced Analytics Press.
- Kendall, M. G. (1938). A new measure of rank correlation. *Biometrika*, 30(1-2), 81-93. <https://doi.org/10.1093/biomet/30.1-2.81>
- Krippendorff, K. (2019). *Content analysis: An introduction to its methodology* (4a ed.). Sage Publications.
- Kuder, G. F. y Richardson, M. W. (1937). The theory of the estimation of test reliability. *Psychometrika*, 2(3), 151-160.
- Lawshe, C. H. (1975). A quantitative approach to content validity. *Personnel Psychology*, 28(4), 563-575.
- Manterola, C., Grande, L., Otzen, T., García, N., Salazar, P., Quiroz, G., Manterola, C., Grande, L., Otzen, T., García, N., Salazar, P. y Quiroz, G. (2018). Confiabilidad, precisión o reproducibilidad de las mediciones. Métodos de valoración, utilidad y aplicaciones en la práctica clínica. *Revista chilena de infectología*, 35(6), 680-688. <https://doi.org/10.4067/S0716-10182018000600680>
- McDonald, R. P. (1999). *Test theory: A unified treatment*. Lawrence Erlbaum Associates, Inc.
- Muñiz, J. (2018). *Introducción a la Psicometría: Teoría clásica y TRI*. Ediciones Pirámide.
- Ortiz Pinilla, J. y Ortiz, F. (2021). ¿Pearson y Spearman, coeficientes intercambiables? *Comunicaciones en Estadística*, 14(1), 53-63. <https://doi.org/10.15332/23393076.6769>
- Padilla, J.-L. y Benítez, I. (2014). Validity evidence based on response processes. *Psicothema*, 26(1), 136-144. <https://doi.org/10.7334/psicothema2013.282>
- Pan, Y.-C. (2009). Evaluating the appropriateness and consequences of test use. *Colombian Applied Linguistics Journal*. <https://doi.org/10.14483/22487085.156>
- Rios, J. y Wells, C. (2014). Validity evidence based on internal structure. *Psicothema*, 26(1), 108-116. <https://doi.org/10.7334/psicothema2013.276>
- Vila, E. (2010). Principios básicos para la construcción de instrumentos de medición psicológica. En Barbero (Coord.), *Psicometría* (169-244).



**Educación**



**Icfes**



**Instituto Colombiano para la Evaluación de la Educación, Icfes**  
**Oficinas Calle 26 #69 - 76. Torre 2, pisos 16, 17, 18**  
**Edificio Elemento, Bogotá. Colombia**

 **@ICFEScol**

 **ICFES**

 **icfescol**

 **ICFES**