

## VALIDEZ Y FIABILIDAD DE LAS OBSERVACIONES SOCIOLOGICAS

*Margarita Latiesa*

### 1. Introducción

Una forma de abordar la problemática sobre la medida en sociología consiste en delimitar diferentes formas de medir y operacionalizar los conceptos por medio de las dimensiones, los indicadores y las escalas. Pedro Gonzalez Blasco<sup>1</sup>, en el capítulo anterior así lo hace y examina estos aspectos esenciales de toda observación científica.

Pero, una vez especificadas las mediciones surge en el presente capítulo la pregunta de hasta qué punto éstas ofrecen una información significativa y correcta. Así, la observación de cualquier realidad social implica una aproximación por medio de la medición y utilización de algún instrumento de medida o técnica de recogida de datos, que debe reunir las siguientes características (Seiffitz, 1976; Repetto, 1977):

1. **Relevancia.** La medida en términos de atributo presupone que la persona o el objeto pueden ser descritos adecuadamente.
2. **Fiabilidad.** Constancia de las observaciones que produce el instrumento de medida. Se deben ofrecer medidas fiables, de manera que se obtengan los mismos resultados al volver a medir el rasgo o aspecto, bajo condiciones similares del individuo u objeto en cuestión.
3. **Validez.** El instrumento de medida que se utiliza en una situación concreta y con propósito determinado debe realmente medir el rasgo que pretende medir. En otras palabras, la medida reproduce el atributo «verdadero».
4. **Sensibilidad.** Posibilidad de hacer suficientes distinciones con el instrumento de medida y permitir la especificidad y la exactitud de los atributos que se miden.

<sup>1</sup> La observación y la medida a las que nos vamos a referir en este capítulo se entienden, al igual que en el capítulo anterior por Pedro Gonzalez Blasco, en el sentido de correspondencia entre números y propiedades, desde la perspectiva de la sociología empírica, cuyo enfoque es similar al de las ciencias naturales y físicas.

Tipificación. El instrumento de medición debe poseer unas normas o patrones estandarizados con los que comparar las puntuaciones que el individuo ha obtenido al realizar la prueba.

Las técnicas de recogida de datos para ser útiles a la observación científica han de cumplir todos estos requisitos, pero muy especialmente los de validez (grado en el que la observación mide el concepto «verdadero» y posición de la persona) y fiabilidad (grado en el que la observación es consistente y constante).

La medición en sociología está sometida a error, lo que generalmente se denomina «problemas de la medida». Los objetivos de este capítulo son dos de estos problemas, la validez y la fiabilidad, que están implícitos en todo tipo de actividad científica que implique la operacionalización y la medición de conceptos y han sido tratados de estos conceptos y medidas por diversos investigadores. El estudio y esclarecimiento de estos conceptos permite depurar la metodología necesaria para llevar a cabo la medición en ciencias sociales, para contrastar las hipótesis y para acumular conocimientos científicos.

El actuar las pruebas de validez y fiabilidad parece no sólo recomendable sino imprescindible en las ciencias sociales. Sin embargo, hay que indicar, que tradicionalmente en la investigación sociológica no se suelen efectuar ni en el análisis de datos, ni en las encuestas, y mucho menos en otros tipos de observaciones más rigurosas. Más recientemente, esta escasez de análisis está cambiando, ya que el «concepto de observación desde los años treinta de encuestas sociales y de «medidas de validez social», ha permitido actualmente la existencia de numerosos bancos de datos y la acumulación de observaciones que han de ser analizadas. Esta sobrecarga hace que «volvamos a la mirada» sobre la validez y la forma en que los datos se obtienen, la fiabilidad de nuestras observaciones y de nuestras encuestas y métodos, a que seamos más meticulosos sobre el proceso de medición y los errores de la medida (Blalock, 1968, 96).

### medida y el error

Hay que asignar números para representar propiedades (Campbell, 1921). La relación correcta los conceptos «abstractos», las propiedades, con los «indicadores empíricos», los números (Blalock, 1968, 6; Riley, 1963, 23). Pero esta conexión que relaciona la medida de cualquier fenómeno resulta imposible de realizar sin errores. Errores diferentes no tendrían la misma percepción de un fenómeno. Esto quiere decir que ningún método de observación está libre de sesgos. El tipo de observación que interesa a la sociología (los conocimientos humanos, las opiniones, los valores, los comportamientos, las capacidades, etc.) impide que los instrumentos puedan observar estas observaciones sean completamente fiables y válidos. Pretender que los datos sean idénticos en todas las observaciones y que la medida refleje una realidad objetiva es imposible. Por otra parte, esto no debe limitar la obligación de los investigadores de ser rigurosos.

Los capítulos de este libro se ocupan también de estos temas, pero, desde una perspectiva muy amplia. La que se adopta aquí. El capítulo de Francisco Alvira expone la validez de los diseños de investigación y el capítulo de Eduardo López-Aranguren expone la validez y la fiabilidad en el análisis de datos. En ambos se trata de evaluaciones específicas aplicadas al diseño en el primer caso y al análisis de datos en el segundo. Sin embargo, en el presente capítulo se lleva a cabo un enfoque desde el punto de vista de la teoría de la medición.

que tiene el investigador de esmerarse en conseguir que sus observaciones sean precisas desde el punto de vista de la validez y de la fiabilidad.

Los resultados obtenidos por una medida no solamente dependen de las diferencias que existen entre los individuos respecto a una característica, sino que también dependen de factores transitorios y estables. Los factores transitorios o errores aleatorios determinan la fiabilidad y los factores estables, o errores sistemáticos, determinan la validez. La variación entre las medidas que produce un instrumento procede de la variación real entre los sujetos y de los errores de medición.

La observación y la medida tienen lugar en situaciones en las que numerosos factores contribuyen a ocasionar los errores (grado de atención o cansancio, grado de cooperación...) y el objetivo es incrementar la validez y la fiabilidad reduciendo tales errores.

En todas las ciencias, físicas naturales o sociales, el producto de una observación o de una medición surge de la interacción entre un OBSERVADOR el INSTRUMENTO DE MEDICIÓN y el sujeto u objeto OBSERVADO. Y las fuentes de error proceden de: a) cada uno de estos tres aspectos considerados aisladamente, b) la interacción entre los mismos, y c) el medio y situación física y/o social en la que se inscriben y desarrollan.

Parece indiscutible que las posibles fuentes de error son superiores en las ciencias sociales, donde los instrumentos de medición son más imprecisos, las situaciones donde se produce el proceso de aplicación más difícilmente reproducibles y los cambios reales de las características y hechos sociales muy probables.

### 2.1 Fuentes de error

Como consecuencia de lo expuesto, no es posible hacer una enumeración exhaustiva de todas las fuentes de error que existen en la investigación sociológica, pero sí al menos podemos mencionar aquellas fuentes más conocidas y citadas. Así, los tipos de errores que propician la no validez y la no fiabilidad son los siguientes (Webb y otros, 1966; Sellitz, 1976):

#### A) Errores del investigador (características y conducta de la persona que recoge los datos)

Una primera fuente de error emana de los propios supuestos y presunciones del investigador. El investigador, posee unas estructuras de orientación, prejuicios, etc. que pueden interferir en la interpretación de las conductas ajenas.

Todos estos valores y supuestos del investigador deben ser conocidos por el mismo, con el fin de reconocer los errores que como consecuencia se pueden presentar. En este punto no interesan cuáles sean los supuestos, tan sólo que el investigador sea consciente de ellos.

En el clínico, el experimentador o el observador participante tienen serias dificultades en detectar la desviación de su propia orientación, al mismo tiempo que recogen los datos. Sin embargo, el director de una encuesta encuentra mayores facilidades en aislar la problemática propia de su orientación y la problemática correspondiente a las desviaciones de los codificadores, entrevistadores, etc, pudiendo, así, descontar estos últimos efectos en la interpretación de los datos.

Es indudable que las características y la conducta de la persona que recoge los

influye en la calidad de los mismos e incluso las variaciones que se producen en la administración del instrumento. El entrevistador en la encuesta, el codificador en el análisis de contenido, o el observador de un grupo, pueden estar más o menos atentos, cansados, etc.; o por el contrario volverse más diestros con las aplicaciones.

#### *Errores del instrumento y del análisis*

Una segunda fuente de error se puede deber a factores de claridad del instrumento o de medida. Las personas pueden interpretar de modo distinto las preguntas incluidas en el instrumento de medida y las respuestas pueden reflejar diferencias de interpretación más que diferencias en el contenido de la pregunta.

También hay que tener en cuenta los factores debidos al formato del instrumento: el cuestionario recoge respuestas verbales, la observación recoge los aspectos actuales, ello hace que aunque se mida la misma característica los resultados sean diferentes. Igualmente, no es lo mismo recoger la respuesta en una escala de 10 ítems, que en una dicotómica de sí/no.

Por último hay que considerar los factores debidos al análisis. Errores en la interpretación, tabulación y análisis estadísticos y no estadísticos.

#### *Errores de los investigadores*

Conductas de los estudiados como sentirse amenazados, cooperar o por el contrario no cooperar con la investigación, y factores personales como humor, salud, distracción, etc., influyen en las respuestas y especialmente el grado de interés que tenga la persona en el tema que se le propone. Por ejemplo, alguien que no está interesado por las cuestiones políticas, ni por las acciones o palabras que pronuncian los miembros del gobierno, contestará con desinterés y desgana a las preguntas de este tipo. La misma medición hace que la observación en un instante y sobre unos temas afecte a los siguientes.

2. Elección de un rol que el estudiado considere adecuado para la situación, pero, que no coincida con su opinión (deseabilidad social).

3. La deseabilidad social consiste en la tendencia de las personas a dar una buena imagen de sí mismas, a dar respuestas de cierta respetabilidad y aceptación social.

Las personas se diferencian en el grado en el que manifiesten este aspecto y ello afecta a los resultados de la medición. Así, se ha confirmado la propensión que existe de dar respuestas deseables, de tal forma que las diferencias en las respuestas a determinadas preguntas pueden manifestar el grado en que una persona está dispuesta a admitir que mantiene conductas indeseables, más que diferencias respecto al contenido de la pregunta (Edwards 1957; Crowne y Marlowe, 1964).

3. La tendencia al asentimiento y a dar respuestas positivas o negativas, respuestas estereotipadas, etc (aquiescencia).

Algunos autores han identificado este aspecto como una propensión de las personas a expresar conformidad o disconformidad (Cronbach, 1950; Phillips, 1971;

Couch y Keniston, 1960). Sin embargo, Scott (1968) observa que no se trata de una tendencia de los individuos, sino que es una característica del instrumento de medida (Scott, 1968)

Los problemas de aquiescencia y deseabilidad social afectan a la validez de las mediciones y han sido muy estudiados en el campo de la psicología <sup>3</sup>.

#### D) Factores de la situación

1. Factores que influyen debido a la situación en que la medida tiene lugar: ambiente relajado, tenso, presencia o ausencia de determinadas personas, etc.

2. Factores mecánicos. Aspectos tan triviales como la falta de espacio para anotar, la rotura de la vestimenta o de un lapicero, pueden contribuir a distorsionar la medida.

#### E) Errores de muestreo

1. Limitaciones de la población a la que se puede generalizar los resultados.

2. La inestabilidad de las características de la población que cambian a través del tiempo hace que los datos de estudios longitudinales no sean comparables.

3. La heterogeneidad de la población a través de las zonas hace que un determinado procedimiento de recogida de datos sea más útil o efectivo en unas zonas que en otras.

4. Las limitaciones que impone el método de recogida de datos o de medición sobre los contenidos para los que resulta apropiado.

También hay que tener en cuenta las limitaciones del método de muestreo de ítems, ya que no es posible utilizar un universo entero con todos ellos.

5. La inestabilidad de los contenidos de la investigación a través del tiempo impide la comparabilidad.

6. La inestabilidad de los contenidos a través de las zonas geográficas.

#### 2.2. Cómo evitar las fuentes de error

El sesgo que introduce el *investigador* debe de reducirse en lo posible, por medio de dos procedimientos: la clarificación teórica y la aplicación de los métodos de investigación social.

La clarificación teórica es un requisito previo e indispensable para determinar qué se va a medir. En efecto, para asegurarnos la calidad de la investigación debemos saber qué cuestiones vamos a investigar y especificar tales cuestiones definiendo los conceptos.

La aplicación de los métodos en el establecimiento de procedimientos de medida implica la utilización de técnicas de recogida de datos y unas normas para la uti-

<sup>3</sup> Ante los problemas de la aquiescencia y la deseabilidad social, la American Psychological Association recomendó que los investigadores eliminaran estas influencias en la construcción de los tests.

de los mismos. En palabras de Sellitz (1976), los *procedimientos de medida* en «las definiciones de trabajo» de los conceptos y las *normas* «facilitan el uso de los datos».

Es de notar que es imposible evitar los errores que propician la falta de *validez* y *fiabilidad* si no se toman las debidas precauciones. Para ello, actualmente se está haciendo un gran interés a un enfoque que consiste en llevar a cabo el proyecto de investigación de manera que maximice la validez y fiabilidad, en lugar de evaluarlas a posteriori. En otras palabras, cuando se planea la investigación se hacen los *diseños* para disminuir los errores.

En otra parte, como ningún instrumento es totalmente válido y fiable, ya que por sí mismo no existe ninguno que refleje exclusivamente diferencias en la característica a medir, debido a que la conducta humana está influida por múltiples factores, se aconseja utilizar *varias medidas* en la misma investigación (Campbell 1959), y emplear *varias medidas* de un mismo concepto.

En un nivel más general, podemos hacer disminuir las fuentes de error durante el control riguroso de todos y cada una de las fases de la investigación. La mayoría de los tipos de errores pueden ser superados de esta manera. El control se refiere a los requisitos metodológicos y técnicos necesarios para el cumplimiento planificado de las observaciones. Es decir, utilizar correctamente los métodos adecuados en la aplicación de las técnicas de recogida de datos, en la preparación de entrevistadores o de observadores en general, en la depuración de observaciones y cuestionarios, así como la codificación.

El control estricto es especialmente necesario en sociología, por la peculiar relación que se establece entre observador y observado. Es más, aunque el investigador puede no poder dedicar el tiempo, esfuerzo y dinero necesarios para la compra de la fiabilidad y validez de su observación, análisis, cuestionario o entrevista, debe controlar cuidadosamente el trabajo de campo, las preguntas, las instrucciones dadas a los codificadores, entrevistadores, etc., para asegurarse de que se eviten los errores.

La literatura existente sobre los problemas de las medidas en ciencias sociales, concede la importancia que tiene para minimizar los errores, al control y al control exhaustivo en todos los momentos de la investigación. Esto quizás se puede hacer prácticamente imposible cuantificar la influencia positiva que tiene un investigador diligente y siempre pendiente de las fuentes de errores que se originan en el trabajo con otros temas preocupados y sensible a estos temas.

### Errores aleatorios y sistemáticos

Las fuentes de error que hemos enumerado propician que en la investigación se incorporen los *errores aleatorios* y los *errores sistemáticos*.

El error aleatorio obedece al azar, es debido a las situaciones *contingentes* que afectan a la persona, de la situación, del procedimiento de medida, etc y afecta a la fiabilidad.

Los errores aleatorios son aquellos factores azarosos que distorsionan la medida del fenómeno; son un conjunto de variables cuyos efectos nos son desconocidos y sus efectos múltiples y actúan de forma tal que no siguen ninguna ley que no sea, precisamente, la del azar. La cantidad de error aleatorio está directamente relacionada con el grado de fiabilidad de la medida del instrumento.

El error sistemático es el que se repite en todas las mediciones y puede ser precisamente el orden de la predicción y la medición (Kerlinger, 1975, 96).

Estos errores están siempre presentes en cualquier medición que llevemos a cabo (Stanley, 1971, 356) y, por tanto, en todos los tipos de investigación (encuesta, observación participante, análisis de contenido, simulaciones y experimentos). Es endémico a la investigación social y, también, en otras áreas como la física o la biología.

El error sistemático, por el contrario, influye *estructuralmente* en la característica objeto de medida y aunque muchas veces pasa inadvertido al investigador afecta a la validez.

Los errores sistemáticos siguen reglas fijas. Un ejemplo de error sistemático sería un termómetro que siempre registra dos grados por encima al tomar la temperatura.

Otro ejemplo donde podemos comparar los errores sistemáticos con los errores aleatorios sería el siguiente: si un rifle está mal calibrado y siempre que se dispara con él tiene un desvío hacia arriba y hacia la izquierda, todas las personas que lo utilicen cometerán un error al disparar sobre una diana. Este error es sistemático y afecta a todas las personas por igual. Ahora bien, si en el momento en el que dispara la persona «X» existe una ráfaga de aire y desvía ligeramente la bala y sin embargo esto no ocurre con la persona «Y», nos enfrentamos a un error aleatorio. También sería un error aleatorio si una misma persona en un disparo tiene una gran concentración y en otro disparo tiene menor concentración. En todos estos casos el error no afecta por igual a los disparos que se producen, ni a los individuos que los efectúan.

En resumen, el estudio de la fiabilidad se centra en la determinación de en qué grado las diferencias de las puntuaciones se deben a influencias aleatorias o casuales. Cuanto menos influyen los errores de azar mayor será la fiabilidad y la consistencia de nuestras mediciones. Pero, un instrumento puede ser fiable y tener errores sistemáticos y por tanto, ser no válido. El estudio de la validez se centra en la determinación de las influencias sistemáticas, que impiden que un instrumento mida realmente los comportamientos sociales.

Para algunos autores, como Sellitz, la validez incluye los errores aleatorios y sistemáticos, de tal forma que si se demuestra que un instrumento de medida es válido, no es necesario demostrar que es fiable. Otros autores, sin embargo, asocian los errores sistemáticos exclusivamente con la validez, y los errores aleatorios con la fiabilidad. A efectos prácticos, es lo mismo optar por un enfoque que por otro, ya que demuestran completamente la validez de un instrumento es imposible y en consecuencia, debemos estudiar siempre la fiabilidad y llevar a cabo un estudio detallado de los errores de tipo aleatorio que pueden influir.

### 3. Validez y fiabilidad

Este apartado constituye el núcleo central del capítulo y será, por tanto, en el que más nos vamos a extender. Se expone en primer lugar, la definición de los conceptos de validez y fiabilidad, en segundo, la clasificación de los diferentes tipos y estrategias que existen (criterios, definiciones, fórmulas, etc.), en tercero, las relaciones existentes entre la validez y la fiabilidad y por último, la evolución histórica que han tenido estos términos.

#### 3.1. Conceptos

La validez hace referencia a que el procedimiento utilizado mide lo que realmente pretende medir y la fiabilidad hace referencia a la propiedad del instrumento de medir de manera consistente y fiable.

lidez trata de determinar el modo en que la «realidad» queda reflejada en que hacemos de ella. Es decir, si tratamos de captar la «inteligencia», que tengamos de nuestra medida sea la «inteligencia» y no otra cosa. Podemos tener «realidades engañosas» y de esta forma nuestras mediciones dadas.

fine como el grado en que una medida mide «realmente» lo que pretende en otras palabras, el grado en que las diferencias de puntuación reflejan «las diferencias existentes entre los individuos en la característica estudiada». Definición tan sencilla de la validez implica serias dificultades, ya que no sabemos si la «realidad» ni la «verdad» puntuación del sujeto y por tanto no compararla con la puntuación obtenida por el instrumento de medida. Así lo, si queremos determinar cuál es la estructura de clases sociales en España medir la «estructura de clases real». Pero es evidente que esta adecuada la medida y la realidad no es fácil, porque cuando los sociólogos decimos «las clases sociales», ¿qué es lo que medimos? Desde luego, en su sentido absoluto, sino en sentido relativo, de acuerdo con definiciones apriorísticas.

«bilitad» hace referencia a la exactitud de las medidas, para ello, el instrumento de medida debe ser adecuado y estar bien «calibrado» y medir lo más exactamente lo que pretendemos captar. El problema de la fiabilidad es el de la determinación del grado en que las medidas con los instrumentos de medición, están libres de errores de tipo aleatorio. La puntuación observada (X) es una estimación de la realidad verdadera, también incluye un cierto error (E):

$$X = XV + E$$

na es algebraica, es decir, X puede sobreestimar o subestimar el valor de la puntuación del sujeto sea independiente o no del error.

ocedimiento de medida es fiable en el grado en que las medidas de un fenómeno resultados similares. El concepto de fiabilidad se ha aplicado para muchas operaciones y conceptos y se ha utilizado de forma muy diferente para su aplicación, no obstante, siempre supone algún tipo de medida de la constancia. Para algunos autores como Köning (1973) esta constancia de la medida el principio más general e importante del desarrollo de la observación. Los coeficientes de fiabilidad aportan una medida de la coincidencia o falta de fiabilidad de las respuestas, pero no indican cuáles son las razones por las que no coinciden. En otras palabras, los coeficientes de correlación que operan con la fiabilidad, cuantifican la diferencia pero no delimitan qué parte corresponde a una de las posibles fuentes de influencia. En el caso del test-retest las atribuyen a cambios en las características transitorias de las personas o a las condiciones de medición. En el caso de diferentes observadores con

— Los autores han propuesto el cambio del término «fiabilidad». Así, por ejemplo, Cronbach, indicó que sería preferible utilizar el término «generalización». Loewinger (1957) aconsejó separadamente el término «homogeneidad» del término «fiabilidad». Sin embargo, ninguna de las propuestas ha tenido éxito y los diversos autores se siguen refiriendo a la fiabilidad para los conceptos y aplicaciones.

un mismo instrumento las fuentes se atribuyen a diferencias entre los observadores. En el caso de formularios diferentes de un test las fuentes se atribuyen a diferencias de muestreo de los ítems.

### 3.2. Tipos o estrategias

Para medir la validez de un instrumento de medida se comparan los resultados que proporciona con otras «evidencias significativas» y se examina la exactitud o inexactitud de los mismos. Existen, no obstante, diversas formas de llevar a cabo estas comparaciones y, por tanto, diversos tipos de validez. Para hacernos una idea de la complejidad conceptual del término algunos autores como Brown (1970) indican que existen en la literatura aproximadamente 40 tipos de validez; posteriormente Brinberg y Mcgrath (1982) enumeran 10 términos distintos de uso frecuente y Messick (1980) encuentra 17 enfoques de validez.

Esta gran variabilidad se justifica en parte por la diversidad de los métodos de observación, ya que los instrumentos de medida difieren entre sí en los siguientes aspectos (Manheim):

- Unos son más estructurados que otros. El experimento de laboratorio sería el extremo de mayor rigidez, formalidad y se fija en ciertos comportamientos muy específicos; la observación participante sería el otro extremo de menor estructuración y mayor flexibilidad en la recogida de datos y en los aspectos que estudia.
- Con unos se pueden hacer más inferencias y generalizaciones que con otros, es decir, permiten en mayor o menor medida la abstracción y el alejamiento de los datos inmediatos tal y como los capta el investigador.
- También difieren en la medida en que los sujetos observados son o no conscientes de este hecho, ya que si son conscientes existe la posibilidad de modificación de la conducta.

Estas diferencias en los métodos de investigación implican una cierta diversidad en las características estudiadas e influyen notablemente en los tipos de validez y fiabilidad que podemos emplear.

A pesar de que numerosos autores se siguen refiriendo a los tipos de validez conviene hacer notar que actuamente, a raíz del artículo de Lawshe (1985), se habla más de tipos de estrategias para verificar el significado de la variable y usos puntuales del instrumento. En estas páginas utilizaremos, no obstante, las dos denominaciones, tanto en lo que se refiere a la validez como a la fiabilidad.

Los medios que se emplean para determinar la validez de un instrumento son diferentes según el tipo de validez a la que se refiera. Así por ejemplo, la validez aparente o la validez de contenido se refieren al instrumento de medida y son independientes de cómo se interpretan las respuestas. Sin embargo, en la validez de construcción, no se trata de una característica del instrumento, sino de una propiedad de las inferencias que pueden hacerse a partir de las puntuaciones con él obtenidas.

Y las mismas observaciones que estamos efectuando con la validez, son aplicables a la fiabilidad. Dadas las grandes diferencias que existen entre los numerosos métodos de observación existentes, cada uno de ellos tiene su propia forma de calificarla. Para llevar a cabo la comprobación de la fiabilidad desde el punto de vista de las observaciones se pueden efectuar las siguientes acciones:

Repetición de los actos de observación por la misma persona.  
Repetición de los actos de observación por personas diferentes, bien en el momento o en distintos momentos.

Contrastación con otros materiales, como son resultados de observaciones que han sido registradas por escrito. Siempre que sea posible en una intención, deben de contrastarse las informaciones con informes estadísticos, aun-  
tralmente estos también pueden tener errores.<sup>6</sup>

La clasificación de los tipos de fiabilidad es más conocida desde el punto de vista de la aplicabilidad del instrumento. De forma resumida podemos identificar los tipos:

1) **Fiabilidad de procedimiento.** El procedimiento consiste en efectuar varias aplicaciones del mismo instrumento de medida con el fin de comparar los comportamientos de los individuos en diferentes momentos del tiempo.  
2) **Fiabilidad de aplicabilidad.** El procedimiento consiste en efectuar varias aplicaciones con distintos instrumentos con el fin de comparar los resultados.  
3) **Fiabilidad de consistencia interna.** El procedimiento consiste en una sola aplicación del mismo instrumento, con el fin de medir el grado en que las respuestas de los individuos a los distintos ítems, u otros componentes de una medida, son consistentes.

Los aspectos implican diferentes formas de contrastar la fiabilidad y, por tanto, los tipos de fiabilidad.  
La fiabilidad en el método de observación se suele calcular por medio del acuerdo de jueces, ya que de la observación de los jueces proviene la mayor parte de los datos (Jueces, 1981). Se puede calcular el acuerdo para toda la sesión, o bien por ítems o por categorías. También se puede calcular por medio de la estabilidad de la observación en momentos diferentes, pero es menos frecuente.

Almente se emplea el método test-retest en las encuestas sociológicas, es la fiabilidad entendida como estabilidad, mientras que en los test psicológicos se emplean los métodos de las dos miradas o fiabilidad entendida como equivalencia interna.

### Tipos de validez

Los tipos de validez se refieren a la clasificación más aceptada:

1) **Validez de contenido.**  
2) **Validez de criterio.**  
3) **Validez de constructo.**

La clasificación hay que efectuar dos observaciones: en primer lugar que incluye la validez aparente en el apartado de validez de contenido, consiste en requerir una mención especial porque es fundamental en el campo social por el tipo de variables y observaciones en las que se interesa; en segundo

<sup>6</sup> Es aplicable cuando en vez de buscar el conocimiento exacto de hechos objetivos, el interés es saber que es lo que la gente cree que son los hechos.

lugar, que no nos vamos a referir a la validez interna y externa porque esta tipología se centra en los diseños de investigación y ha sido expuesta previamente en este libro. Podemos, no obstante, recordar de forma muy simplificada que la validez interna se basa en comparaciones entre las características diferentes de ciertos grupos y la validez externa se ocupa de las generalizaciones (Alvira, 1976, 1991). O en otros términos, la validez interna responde a la pregunta ¿se obtendrían resultados diferentes si se hubieran utilizado procedimientos diferentes?, mientras que la validez externa se pregunta ¿cuán generalizable es el procedimiento utilizado? (García Ferrando, 1985, 34).

### Validez de contenido

En este tipo de validez se utilizan los juicios con frecuencia consensuados, para determinar si el contenido de los ítems es apropiado o no lo es.

La validez de contenido se pregunta si el instrumento de medida sirve para medir el comportamiento que pretende. Para ello, debemos especificar el universo de los comportamientos y la variedad de formas en que podemos medirlos; y en el caso de los test identificar todos los posibles ítems útiles para la medición. Se asume que todas las mediciones sirven para medir el concepto objeto de estudio, pero no podemos asegurar que se están midiendo todas las dimensiones del concepto.

La validez de contenido se refiere al grado en que los ítems de una escala o test representan un determinado universo temático. Se fija en una propiedad del instrumento de medición (formulación de los ítems) y no en las inferencias que pueden hacerse con las puntuaciones obtenidas, por ello, algunos autores cuestionan que pueda hablarse con propiedad de validez.

La validez de contenido trata de determinar la relevancia o representatividad de los ítems en cuanto muestra adecuada de un dominio previamente especificado. La adecuación del contenido de los ítems con un plan previo, responde a la necesidad de construir un instrumento y nada tiene que ver con el significado de las puntuaciones. Por esta razón, algunos autores han llegado a cuestionar que la validez de contenido sea un tipo de validez y a afirmar que no es apropiada para determinar si se «mide lo que se pretende medir» (Messick, 1980, 1981; Tenopir, 1977; Cronbach y Meehl, 1955). Otros autores, sin embargo, como Yalow, Popham (1983) y Ebel (1983) aclaran que si el contenido es relevante y representativo esto significa que existe una inferencia implícita, porque, aunque es una propiedad del instrumento, una vez que se obtienen las respuestas, la validez de contenido hace posible su interpretación. Así, el contenido permite y justifica interpretar resultados.

Algunos autores también denominan validez aparente a la validez de contenido, pero conviene que delimitemos sus características propias por la transcendencia que en el campo sociológico tiene este tipo de validez.

Un instrumento tiene validez aparente si mide lo que parece. Este tipo de validez se considera en muchos casos trivial y no se tiene en cuenta, pero sin lugar a dudas tiene su importancia.

La validez aparente suele emplearse en dos sentidos o matices:

- 1) Un instrumento tiene validez aparente si parece apropiado a la situación. Así por ejemplo, un test de aptitud numérica parece apropiado para predecir el rendimiento en matemáticas.
- 2) Un instrumento tiene validez aparente si la denominación coincide con la

rece a simple vista coherente con la formulación de los ítems y, por lo tanto, mide lo que se dice que mide.

validez aparente se ocupa de las medidas que se efectúan directamente. Así, la validez de un atleta se cronometra observando su comportamiento. Las preguntas cuestionario como el sexo, los metros cuadrados que ocupa la vivienda, los rendimientos, etc., son ejemplos de observaciones directas, donde la validez rumental de medida se presenta en su propia apariencia. El problema que se plantea es que la apariencia puede ser momentánea y no representar la realidad íntima.

mbien se denomina validez facial o análisis lógico, porque la validez de los parece a simple vista utilizando el sentido común. Se aplica cuando los datos asados en la observación directa y no se necesitan inferencias. La mayor parte preguntas que se efectúan en los cuestionarios en Sociología pertenecen a este tipo, edad, etc.

#### de criterio (concurrente y predictor)

validez de criterio se mide directamente por medio de los coeficientes de correlación que se establecen con otras medidas o criterios. Suele dividirse en concurrentes y predictoras. Las medidas concurrentes se obtienen simultáneamente y predictoras (la medida en el futuro).

validez concurrente hace referencia a alguna característica por la que los individuos difieren en el presente. La validez predictoras hace referencia a alguna característica por la que los individuos se diferenciarán en el futuro. Lo importante es que ambos casos se trata de una validez pragmática porque la validación se efectúa con otro procedimiento de medida empírico. No interesa aquí si se mide algún concepto teórico determinado, lo que interesa es saber si el instrumento, si queremos predecir el rendimiento académico de los alumnos y en otros determinadas características como potenciadoras del mismo, no nos interesa por qué esas características son buenas predictoras, sino que el hecho de empujarse que efectivamente lo son, ya es suficiente.

enomina validez pragmática porque comparamos los resultados con algún indicador, como por ejemplo, contrastar las fuentes de datos secundarias fuentes de determinadas preguntas del cuestionario elaborado por el investigador. El criterio también puede ser alguna predicción basada en los resultados obtenidos validando, o puede basarse en grupos que se sabe ocupan posiciones en la variable que se está midiendo.

#### teórica o de constructo

Los tipos de validez están relacionados: la validación pragmática y la validez teórica no son excluyentes, ya que un instrumento de medida puede tener objetivos y por otra parte, no debemos quedar satisfechos si tan sólo hemos validado la validez pragmática y no hemos respondido al por qué de su efectividad. La validez teórica no se trata de demostrar un «tipo» de validez y quedar satisfechos,

hay que validar desde diferentes aspectos y difícilmente podemos concluir de manera definitiva que nuestras observaciones son válidas.

La validez no queda demostrada por unos coeficientes de correlación aislados (validez concurrente o predictiva), ni por juicios sobre la relevancia del contenido, sino que queda integrada en todo un proceso de verificación de hipótesis teóricas, donde entra tanto lo empírico como los juicios racionales. Como consecuencia, podemos afirmar que la validez de constructo subsume las nociones tradicionales de validez de contenido y de validez de criterio (Morales, 1988, 371). Así, la cuestión del contenido tiene que ver con la relevancia de criterio y la validez del criterio con la utilidad del instrumento, pero, ni el contenido, ni las relaciones con otras variables, pueden separarse de la interpretación del constructo. En este sentido la validez de constructo integra todos los tipos de validez.

La validez de constructo (o teórica) es un proceso mediante el cual acumulamos evidencias, pruebas empíricas, sobre relaciones teóricamente importantes que apoyan —o demuestran— en cierto grado una determinada inferencia o interpretación de las puntuaciones. De las relaciones observadas entre lo que medimos y otras variables, deducimos la naturaleza de las respuestas.

Por una parte hay que considerar el papel que tiene el instrumento de medida, ya que Fiske (1975) ha demostrado que las intercorrelaciones entre un conjunto de variables que miden la personalidad varían según sea el instrumento de medida, incluso utilizando una misma muestra de personas. Pero, la validez de constructo implica la validez no solamente del instrumento sino, también, de la teoría. No se interesa por la predicción de comportamientos, sino por la medición de alguna característica construida a nivel teórico (que no se puede observar ni medir directamente) y que se supone que nuestro instrumento de medida cuantifica. La característica que se mide trata de reflejar el grado en el que realmente la posee un individuo, relacionando el propio instrumento de medición con la estructura teórica general. Así, la inteligencia, la actitud hacia el aborto, etc., son constructos teóricos y de su validación como tales se ocupa la validez de constructo, aunque normalmente donde más se emplea es para los rasgos de personalidad.

Para determinar la adecuación de la medida del concepto Campbell y Fiske (1959) proponen dos tipos de validación:

- \* Validación convergente. Diferentes medidas de un concepto proporcionan resultados semejantes. Para comprobarlo es necesario medir el concepto con varios procedimientos.
- \* Validación divergente. La medición del concepto se diferencia de otros conceptos. Para comprobarlo es necesario medir con el mismo procedimiento otros conceptos de los que se supone se debe diferenciar el de la investigación.

El método de Campbell y Fiske se basa en la matriz multi-rasgo multi-método y trata de verificar los dos tipos de hipótesis (convergente y divergente). En la validación convergente se comprueba la correlación entre métodos distintos que presumiblemente miden el mismo rasgo. Así, cada instrumento de medición se concibe como un método-rasgo. En la validación divergente rasgos distintos medidos con el mismo método o métodos semejantes no están relacionados. Las correlaciones entre métodos distintos que midan el mismo rasgo, deben ser mayores que las correlaciones

<sup>7</sup> La validez de constructo ha generado muchos desacuerdos entre los sociólogos (Sjoberg, Nett, 1968).

re métodos semejantes que miden rasgos supuestamente distintos. Luego, el de estos autores «trata de clarificar el significado utilizando diversos enfo- métodos para medir lo mismo (convergenencia) y se distingue de conceptos afi- regencia)» (Morales, 1988, 435).

ra bien, estas pruebas de validación de constructo no son universales y rales. Así, una actitud perfectamente medida y validada con un instrumento ser muy útil para una época determinada, pero no para otra. Las generaliza- probadas no son estables en las ciencias físicas y mucho menos en las ciencias . En consecuencia, tenemos que afirmar con Nunnally (1978), que las pre- validación de constructo más que probar la «verdad» de la teoría o construc- tivo en un instrumento de medición, muestran su grado de pertinencia y uti- ara investigar una realidad.

### Tipos de fiabilidad

este apartado vamos a presentar los tipos de fiabilidad y las formas de deter- s.

#### Fiabilidad como estabilidad

estabilidad de los resultados de un instrumento de medida se establece com- ) los mismos en aplicaciones repetidas. Las diferencias en los resultados y por a inestabilidad de las mediciones, puede ser debida a cambios reales ocurri- cambios debidos al azar y a factores extraños. Es difícil determinar qué parte nde a los cambios reales y qué parte corresponde a los errores, aunque al- uiores han expresado algoritmos para hallarlos, demostrando que es posible ir los cambios verdaderos de los aleatorios mediante tres aplicaciones o bien e los indicadores múltiples de los modelos path (Heise, 1969; Wiley, Wiley, 374).

ndo el instrumento de medida se basa en la observación, es necesario reali- levado número de aplicaciones; en la encuesta se utilizan normalmente dos ones; y el análisis de contenido es un caso intermedio entre los dos expues-

riodo más conocido de comprobación de la fiabilidad como estabilidad se con el nombre test-retest. En una investigación por encuesta, se aplica el cuestionario a las mismas personas en distintos momentos del tiempo, bajo ras equivalentes y posteriormente se comparan los resultados. Las concor- obtenidas entre las puntuaciones de las dos pruebas se establecen mediante lación. En consecuencia, el coeficiente de fiabilidad se define matemática- or medio del coeficiente de correlación.

l test-retest también se pueden distinguir las dos medidas que hemos indica- una parte el cambio real o su inversa, el coeficiente de estabilidad; y por cambio aleatorio o su inversa, el coeficiente de fiabilidad. A mayor cambio or será el coeficiente de estabilidad y a mayor cambio aleatorio, menor será iente de fiabilidad.

método plantea algunos inconvenientes que se derivan de su aplicación. En ugar, el mismo proceso de medida al repetirlo puede aumentar las diferen- ras como por ejemplo el interés y la motivación pueden ser menores en la

segunda aplicación debido a que el encuestado ya participó activamente en la pri- mera. También puede recordar las respuestas que emitió en la primera aplicación y no contestar espontáneamente en la segunda. Así, el cuestionario es el mismo pero la situación de la entrevista puede ser diferente.

También puede darse la posibilidad de cambios reales entre las dos aplicaciones como consecuencia de la primera aplicación. Así por ejemplo, hay encuestados a los que el descubrimiento de un tema por medio del cuestionario, hace que se interesen por el mismo y cambien su opinión o actitud realmente como consecuencia de la in- formación adicional.

Todos estos problemas y limitaciones se concretan según Carmines y Zeller (1979) en cuatro:

- 1) A menudo los investigadores sólo pueden hacer una aplicación, no sólo por el elevado coste que supone aplicar dos veces el mismo cuestionario, sino porque además hay veces que es imposible hacerlo, aunque se dispusiera de dinero suficiente.
- 2) El cambio verdadero <sup>8</sup> que ha podido suceder en determinadas característi- cas o actitudes, se interpreta en este método como inestabilidad y subestima la fiabilidad, haciendo que los coeficientes sean menores. Esto sólo se puede solucionar separando el efecto del cambio verdadero del cambio aleatorio o error. Heise (1969) ha demostrado que esto es posible con tres aplicaciones del mismo cuestionario.
- 3) La reflexividad en las ciencias sociales. Es decir, la medición de un fenóme- no puede inducir a modificaciones del fenómeno mismo. Al obtener infor- mación de una persona en un momento de tiempo puede hacer que la per- sona se sensibilice hacia el tema y cambie su respuesta en el momento dos. Consecuentemente, la fiabilidad puede ser menor debido a esta reflexividad.
- 4) También se puede sobrestimar la medida de fiabilidad debido a la memo- ria. Si el tiempo que media entre las dos aplicaciones es corto, la memoria en la segunda puede influir y los temas que se recuerdan pueden dar coefi- cientes de correlación más altos que aquellos otros temas que no se recuer- dan (Nunnally, 1964).

Para obviar estos problemas se aconseja que transcurra un tiempo suficiente en- tre las dos aplicaciones para borrar los efectos de la primera aplicación y no sobre- timar la fiabilidad. Ahora bien, tampoco puede ser demasiado tiempo, para que no se produzcan cambios verdaderos y no subestimemos la fiabilidad. Si el tiempo no es suficiente puede influir la memoria de lo que se contestó en la primera aplicación y dará lugar a coeficientes de fiabilidad más altos; pero si se distancian mucho las aplicaciones, pueden darse cambios de hecho en el sujeto y dar coeficientes más bajos de los reales. En general, el intervalo de tiempo puede ser más reducido cuan- to menor sea la edad de los sujetos y cuanto más susceptible sea el rasgo medido de experimentar cambios.

Especialmente en psicometría se recomienda abiertamente que no se utilice el test-retest, o que si se hace, «que sea con aquellos instrumentos en los que menos

<sup>8</sup> Sobre este tema de la diferencia entre cambio verdadero y cambio aleatorio, o lo que se denomi- na comúnmente como el problema de la estabilidad de la medida, se recomienda la lectura de los si- guientes autores: Carmines y Zeller (1979); Heise, (1969); Wiley and Wiley (1970, 1974); Erikson (1978)



la memoria, la práctica o el aprendizaje, como son todos aquellos que implican más o menos rutinarias: sin embargo, para aquellos elementos que exigen azonamiento, desarrollo de estrategias, ingeniosidad..., no es adecuado» (R., 1986, 196). Así por ejemplo, en los test utilizados en psicometría sobre resolución de problemas, puede ocurrir que si en la primera ocasión se ha resuelto, en una segunda se puede ir directamente a la solución, sin pasar por las etapas intermedias. Por el contrario, la mayor parte de las variables que se miden en sociología son esta problemática.

Comparativamente existe un mayor número de publicaciones dedicadas a la fiabilidad entendida como equivalencia o consistencia interna que a la entendida como estabilidad. Esto se debe a que tradicionalmente ha sido el campo de la Psicología la que en mayor medida se ha ocupado de estos problemas. Y para calcular la fiabilidad de sus test son apropiados los conceptos de «equivalencia» o «consistencia interna». Sin embargo, el concepto de «estabilidad» que se recoge por el método test-retest es más adecuado para la disciplina sociológica, que se ha ocupado con algún retraso al estudio de esta problemática.

#### *fiabilidad como equivalencia*

Se trata de disponer de dos o más formas paralelas de un mismo instrumento de medida, que permita obtener dos o más puntuaciones del mismo sujeto, aplicadas simultáneamente o dejando un lapso de tiempo entre ellas. Se supone que los sujetos deberían responder de forma equivalente a estas muestras de ítems. La diferencia en la respuesta será la medida del error y se llama coeficiente de equivalencia o fiabilidad. Matemáticamente se expresa por la correlación entre las dos acciones.

Una cierta medida este procedimiento es similar al método test-retest, porque se puede también aplicar dos test a las mismas personas, la diferencia es que no se aplica a segunda vez el mismo test, sino otro test alternativo, diferente, aunque se refiere que mide lo mismo que el primero.

Las características que deben cumplir los instrumentos de medida paralelas

Deben tener el mismo número de elementos intercambiables uno a uno.

La redacción y la estructura de cada elemento debe ser idéntica en los elementos paralelos de ambas formas.

El contenido y el objetivo apreciado, elemento a elemento, deben ser el mismo.

Los índices de dificultad de los elementos deben ser iguales.

Las instrucciones dadas para la realización de la prueba, el tiempo asignado, las condiciones en que se aplica, etc., han de ser las mismas.

También deben ser idénticos los aspectos externos: presentación, formato, etc.

Un problema más grave que comporta este procedimiento es, sin duda, la consistencia de las formas paralelas de manera que sean equivalentes. Así, las medias, varianzas típicas, etc., verdaderas, de ambas formas deben ser idénticas y las medias encontradas empíricamente, casuales o aleatorias.

Matemáticamente el modelo de las pruebas paralelas es muy estricto en sus supuestos. Dos pruebas son paralelas si cumplen las siguientes condiciones:

1. Tienen idénticas medias, varianzas y covarianzas.
2. Ambas correlacionan en idéntico grado con las puntuaciones verdaderas<sup>9</sup>.  
La correlación entre las dos pruebas paralelas indica la proporción de verdadera varianza, o varianza que corresponde a un único factor común (Linn, Werts, 1979). La unidimensionalidad, por tanto, está implícita en este modelo<sup>10</sup>.
3. La varianza de cada prueba paralela no explicable por las puntuaciones verdaderas se debe a errores de medición.

La variable no experimenta cambio real, o en otras palabras, las puntuaciones verdaderas son siempre las mismas para cada sujeto, las diferencias entre las puntuaciones observadas se deben a que los errores de medición varían; los componentes verdadero y de error son independientes, en consecuencia la correlación entre los mismos es cero.

Este método no tiene en cuenta la inestabilidad que proporciona el tiempo como fuente de no fiabilidad, ya que se administran a los mismos individuos en una misma sesión dos formularios supuestamente equivalentes del mismo test y aunque los dos formularios contienen ítems distintos, todos ellos tratan de medir la misma característica. La correlación entre las puntuaciones de los dos formularios representa la medida de hasta qué punto los ítems miden la misma característica y son consistentes.

Tampoco se tienen en cuenta las fluctuaciones azarosas de las personas o de las situaciones en las que transcurre la administración, debido a que las dos formas se administran a un tiempo en una sola sesión.

Para Carmines y Zeller (1979), este método presenta las siguientes ventajas y desventajas respecto al método test-retest:

- 1) Tiene la ventaja de que reduce los problemas de memoria, que afectan a la fiabilidad en el método test-retest.
- 2) Al igual que el test-retest tampoco distingue los cambios verdaderos de los errores aleatorios, por lo que la fiabilidad es menor.
- 3) La desventaja respecto al test-retest es, como ya hemos visto, la dificultad de construir formas alternativas que sean paralelas: «si ya es difícil construir un test mucho más difícil lo es construir dos que midan exactamente lo mismo» (Carmines, Zeller, 1979, 41). Así, los problemas que plantea su utilización son los siguientes:
  - 3.a) La falta de equivalencia de los ítems muestrados, ya que aunque se construyen con la intención de que sean equivalentes, la formulación de los mismos puede dar origen a diferencias en las puntuaciones.
  - 3.b) La diferente experiencia y conocimiento que un sujeto puede tener en uno o en otro de los universos de ítems aunque sean equivalentes, ya que el sujeto puede entender mejor, o conocer mejor, unos aspectos que se reflejan en una aplicación que los que se reflejan en otra.

El procedimiento de los formularios alternos se desarrolló sobre los test de inteligencia y aptitud, donde es más fácil medir una característica repetidamente por medio

<sup>9</sup> Esta condición es imposible de comprobarla.

<sup>10</sup> La unidimensionalidad se define en el sentido de que todos los ítems miden una sola y única dimensión.

formas paralelas. Pero la dificultad de conseguir construir dos formas equivalentes de una medida ha impedido una mayor utilización.

La psicometría se considera un buen procedimiento para los instrumentos que deben apreciar la velocidad de ejecución (número de elementos resueltos por 1 de tiempo). Se recomienda que entre las dos formas alternativas medie una acción temporal de dos semanas.

Pesar de la defensa que determinados psicometras hacen de este método, es el más inconveniente presente porque se parte de unos supuestos que es improbable que se den en la práctica, razón por la cual no se utiliza mucho esta técnica psicológica.

*Fiabilidad como consistencia*

Existen dos posibilidades: a) Se construye un instrumento de doble longitud, que de partir en dos mitades, y todos los ítems miden el mismo rasgo o característica. Lo que debe darse una coherencia o consistencia en las respuestas de los sujetos en los dos subconjuntos de la prueba; b) O bien se establece hipotéticamente el rasgo total de ítems que miden un rasgo o característica y se compara con la acción observada en una muestra de ítems; o bien se supone un número indefinido de pruebas paralelas y se compara con la muestra obtenida. En ambos casos, el coeficiente de puntuaciones que se obtienen en la misma prueba, obtenemos una medida de consistencia interna.

*Fiabilidad de dos mitades*

Se aplica una sola vez, pero en la aplicación de la prueba se obtienen dos puntuaciones para cada sujeto y nos permite calcular un coeficiente de fiabilidad, que se llama coeficiente de consistencia interna.

Para obtener dos puntuaciones se divide el conjunto de elementos que integran el instrumento de medida en dos mitades equivalentes respecto a la característica medida. En realidad, son dos subconjuntos, lo que puede hacer pensar que son dos formas paralelas.

El método de las dos mitades es similar al anterior de formas paralelas equivalente de hecho, ambos implican los mismos supuestos y restricciones (unidad de medida, etc.). La diferencia con el procedimiento anterior es que allí cada forma es una prueba que se aplica independiente de la otra, mientras que aquí se aplica una sola prueba que después se parte en dos mitades. Por lo demás, consiste en una modalidad del procedimiento de equivalencia, ya que en este caso la fiabilidad de los resultados de diferentes muestras de ítems se estiman mediante el mismo ítem interno de las respuestas a los ítems en un solo test. Es lo que se denomina «corrección de corte por la mitad» y que supone un caso especial de medidas paralelas o alternas.

La forma más usual de obtener dos mitades es eligiendo los elementos pares y impares, pero existen otras muchas, como incluir en la primera mitad los primeros y en la segunda mitad, los segundos. Sin embargo, esta partición no permite comparar los resultados, ya que la fatiga y tensión en la primera prueba afecta a la segunda.

También se procede para igualar ambas mitades a considerar la dificultad de los ítems. La dificultad se determina dividiendo para cada elemento el número total de sujetos que los resolvieron correctamente por el de aquellos que lo intentaron. El valor oscilará entre 0 (no resuelto por nadie, difícil) y 1 (resuelto por todos, fácil). De esta forma se distribuyen los elementos a una u otra mitad. También es conveniente, a la vez, procurar un equilibrio entre ambas mitades en el contenido y los objetivos apreciados por los elementos que las integran.

En este procedimiento se puede demostrar que existe relación entre la longitud (número de elementos de la prueba) y la fiabilidad. A mayor número de elementos, mayor fiabilidad o consistencia interna. Ahora bien, esta relación no se da con cualesquiera elementos o ítems que se añadan, sino cuando son equivalentes a la que integran la prueba, de forma que la estructura y composición de la misma queda inalterada.

Existen tres procedimientos de cálculo para la estimación de la fiabilidad de las pruebas de longitud doble.

- a) Ecuación de Spearman-Brown (Spearman, 1910; Brown, 1910), que fue descubierta al mismo tiempo de forma independiente por estos dos autores.

$$R_{xx} = \frac{2r_{xx}}{1 + r_{xx}}$$

$R_{xx}$  = Coeficiente de fiabilidad total.  
 $r_{xx}$  = Correlación entre las dos mitades (varía de 0 a 1)

- b) Ecuación de Rulon

$$r_{xx} = 1 - \frac{S_d^2}{S_t^2}$$

- c) Guttman

$$r_{xx} = 2 \left( 1 - \frac{S_{1a}^2 + S_{2a}^2}{S_t^2} \right)$$

- d) Posteriormente se ha desarrollado la fórmula de Kuder-Richardson (1937). Permite como las anteriores obtener un coeficiente de consistencia interna, pero además un coeficiente de homogeneidad entre los elementos, ya que puede haber instrumentos con mitades consistentes, pero heterogéneos en sus elementos.

Se utiliza para instrumentos calificados en dos categorías (acierto, error, sí, no).

$$r_{xx} = \frac{n}{n-1} \left( \frac{S_1 - \sum p q}{S_2} \right)$$

mero de elementos.  
-p. iente entre el número de sujetos que aciertan y el total.

inconveniente del método de dos mitades es que existen numerosas estimaciones de la fiabilidad. Así, en una escala de 10 ítems, existen 125 formas de dividirlo, obtener diferentes estimaciones de fiabilidad. La ventaja es que se aplica una vez a los mismos individuos (Carnines, Zeller, 1979).

dos mitades se consideran formas alternas del mismo test y el coeficiente mide indica la consistencia interna del test. Un elevado coeficiente supone situación del individuo no se halla determinada por el muestreo concreto de los que ha contestado en cualquiera de las dos mitades.

un principio, las aplicaciones que se llevaron a cabo de este procedimiento que las dos mitades debían ser equivalentes, y cada una de ellas representaría la mitad del test (Guilford, 1954). Se comparan los ítems pares con los impares y el coeficiente de correlación que es una estimación del coeficiente de equi- test, por medio de la fórmula de Spearman-Brown.

embargo, estudios más recientes estiman que si todos los ítems del test tratan de la misma característica, deberían ser comparables cualesquiera dos mitades calculen al azar y no dos que se consideren equivalentes. Los coeficientes que calculado basándose en esta nueva concepción son el coeficiente alfa y la 20 de Kuder-Richardson. Estos dos índices expresan la correlación media por la mitad para todos los casos posibles de división del test en dos partes. La tendencia de procedimiento anterior (pares, impares) estos coeficientes de equi- lisis en el concepto de homogeneidad entre todos los ítems de un test. La a es la siguiente: ¿hasta qué punto miden la misma característica todos los un test? <sup>11</sup>.

#### del Universo de ítems <sup>12</sup>

los indicado que una limitación importante del método de las dos mitades odemos hallar numerosos coeficientes de fiabilidad. Así, el coeficiente halla- nbién existen otros métodos de medición de la homogeneidad basados en la comparación de las puntuaciones obtenidas de una medición con la varianza que resultaría si todos los ítems no correlacionados y con la varianza que resultaría si todos los ítems estuvieran correlacionados (Scott, 1960).

recientemente Cronbach y colaboradores (1972) han ampliado el modelo original del uni- versum. En el modelo original se controla como fuente de error el contenido de los ítems. En el posterior de este modelo se amplía la interpretación del universo, que ya no es solamente de ítems, sino de situaciones, y en él se incluyen otras fuentes de error (ítems, contextos, etc.). En los modelos se controlan diversos tipos de errores y no hay un único coeficiente de fiabilidad, er tantos como fuentes de error se analicen. Un resumen de este modelo lo lleva a cabo (1976).

do para la primera y segunda mitad es diferente del hallado para los pares e impares y así sucesivamente. Sin embargo, en el modelo de universo de ítems se obtiene un solo coeficiente.

Por otra parte, en el método de las dos mitades existen unos supuestos muy restrictivos. Por el contrario, en el modelo de la muestra de ítems procedentes de un determinado dominio, no supone necesariamente unidimensionalidad <sup>13</sup>, y aun- que de hecho se interprete que la tiene, no se exige explícitamente.

La dicotomía de puntuación verdadera y error se sustituye por el universo de ítems cuya suma definitiva la puntuación verdadera. El coeficiente de fiabilidad indica si una puntuación observada (en una muestra de ítems) estima bien la pun- tuación verdadera (universo de ítems). En realidad es como si se supusiera un nú- mero indefinido de pruebas paralelas, que se forman aleatoriamente por ítems. La diferencia con las pruebas paralelas es que en ellas todos los ítems tienen idéntica varianza e idéntica correlación con el total del dominio de ítems, mientras que en las pruebas formadas aleatoriamente no se exigen estos requisitos. De esta forma y como afirma Morales (1988) el modelo de las pruebas paralelas viene a ser un caso particular del modelo del universo de ítems.

Los supuestos del modelo de universo de ítems son los siguientes (Nunnally, 1978):

- 1/ Se supone que existe una población de ítems que pertenecen al dominio o constructo que deseamos medir; esta población o universo de ítems podemos considerarlo como hipotéticamente infinito.
- 2/ Cualquier instrumento, escala, test, etc., está compuesto por una muestra de k ítems, tomados aleatoriamente de ese conjunto hipotético de ítems; esta muestra puede ser tan pequeña o tan grande como se quiera.
- 3/ La puntuación verdadera de un sujeto es la que obtendría si respondiera a todos los ítems de la población; la puntuación observada es la que proviene de sus respuestas a una muestra de ítems, y es sólo una estimación de su puntuación verdadera.
- 4/ Las puntuaciones observadas serán fiables en la medida que exista una correlación alta con la hipotética puntuación verdadera; es decir, si el instrumento compuesto por k ítems tiene una correlación alta con el instrumento compuesto por todos los ítems. La fuente de error que se considera aquí es la que proviene de un deficiente muestreo de ítems; la que tiene que ver con el contenido de los ítems; existen otras fuentes de error, pero esta es la verdaderamente importante que se considera.
- 5/ Un elemento básico de este modelo es suponer que existe una matriz de correlaciones con las intercorrelaciones de todos los ítems del hipotético universo de ítems. La correlación media de esa matriz indicaría en qué grado todos los hipotéticos ítems tienen algo en común, que es precisamente lo que se pretende medir con la muestra de ítems. Lo que no se supone necesariamente es la existencia de un único factor. El modelo supone que las puntuaciones de los ítems se han tipificado, con lo que las varianzas de los mismos son idénticas (e iguales a la unidad).
- 6/ Las correlaciones de cada ítem con la suma de todos los demás deben ser iguales.

<sup>13</sup> No se identifica varianza verdadera con un único factor común (Morales, 1988).

partir de estos presupuestos Nunnally (1978) ha comprobado que:

- 1) La correlación entre un test de  $k$  ítems (puntuaciones observadas) y el test compuesto por todos los ítems de la población (puntuación verdadera), es igual a la raíz cuadrada de la correlación media de una serie de test paralelos todos con idéntico número  $k$  de ítems procedentes de la misma población o dominio de ítems:

$$r_{wv} = \sqrt{r_{kp}}$$

Las fórmulas para hallar las puntuaciones típicas son:

$$r_{11} = (K/K-1) ((R-K) / R)$$

$$r_{11} = (K/ K-1) (1 - (K/R))$$

es la suma de todas las correlaciones en la matriz, incluyendo los valores de la diagonal que son las varianzas de los ítems.

$K$  es la suma de las correlaciones menos las varianzas de los ítems.

Las fórmulas para hallar las puntuaciones directas que en la práctica es lo normal son:

$$r_{11} = (K/ K-1) ((C - \Sigma \sigma_i^2) / C)$$

$$r_{11} = (K/ K-1) (1 - (\Sigma \sigma_i^2 / \sigma^2))$$

numero de ítems

matrónico de la varianza de los ítems

varianza total

es la fórmula más conocida como el coeficiente alpha de Crombach. Este se ha demostrado que su coeficiente «alta» equivale a la fiabilidad media que se obtiene dividiendo el test en todas sus posibles mitades y aplicando la fórmula de an-Brown en cada posible división del test en dos mitades.

El coeficiente alpha de Crombach es una generalización del coeficiente de Kuder-Richardson (1957) para estimar la fiabilidad de las escalas dicotómicas. Así, si de ítems dicotómicos utilizaremos la fórmula de Kuder-Richardson 20, que es particular de la fórmula más general de Crombach (1951):

$$KR20 = N(N-1) / (N - \Sigma p_i q_i / \sigma_i^2)$$

numero de ítems dicotómicos.

proporción de respuesta positiva a los ítems.

$p_i$

varianza total compuesta.

Tiene la misma interpretación que el coeficiente alpha: es una estimación de la correlación esperada entre un test y una forma alternativa hipotética que contiene el mismo número de ítems.

Actualmente, ya no se usa el método de las dos mitades (de Spearman-Brown a partir de la correlación entre dos formas paralelas del mismo test). Se ha impuesto el coeficiente alpha de Crombach (o Kuder-Richardson 20). Las ventajas son evidentes, como hemos expuesto con anterioridad.

El coeficiente alpha responde a un modelo conceptual muy claro y simple: la proporción de varianza verdadera (fiabilidad) es igual a la varianza compartida dividida por la varianza total. La varianza verdadera queda definida operativamente por la suma de las covarianzas: por lo que discriminan los ítems precisamente por estar relacionados unos con otros. Este «estar relacionados» es lo que con propiedad se llama también consistencia interna u homogeneidad. El problema es la confusión que se da entre estos términos y el de unidimensionalidad.

El coeficiente alpha indica en qué proporción discriminan los ítems precisamente por estar relacionados entre sí. La interpretación más clara del coeficiente es la que se limita a expresar lo que señala la misma fórmula: «es la proporción de covariación, que si es grande implica relaciones claras entre los ítems» (Morales, 1988, 248). En este caso el término consistencia interna parece adecuado. Sin embargo, no está tan claro que se pueda identificar consistencia interna y unidimensionalidad, a pesar de que es precisamente Crombach el que introduce la idea de un único factor.

El coeficiente de fiabilidad alpha se utiliza como criterio para evaluar hasta qué punto un test o escala está compuesto por ítems lo suficientemente homogéneos, como para justificar que su suma constituya una medida del constructo subyacente. Lo que se verifica no es la correlación con unas supuestas puntuaciones verdaderas (contenidas en un universo de ítems), sino cuánto hay de común en los ítems, ya que no todo es común.

Este coeficiente nos dice cuánto hay de interrelación, pero no cómo es esa interrelación. El cómo se relacionan las variables nos lo dirá la matriz de correlaciones y el análisis factorial, pero según Crombach, y en general según la teoría clásica, un test es poco fiable y sus ítems poco homogéneos, si una proporción apreciable de varianza se debe a que los ítems son distintos; lo que se busca es que las diferencias en las puntuaciones observadas provengan fundamentalmente de que los sujetos son distintos en aquello que se intenta medir, no de que los ítems sean distintos.

El coeficiente alpha tiene limitaciones si se interpreta como índice de unidimensionalidad. Por otra parte, no guarda una relación clara ni con la magnitud u homogeneidad de las correlaciones, ni con su estructura factorial, por eso se presta a interpretaciones equívocas. No puede interpretarse este coeficiente de forma automática como un índice de homogeneidad o de unidimensionalidad sin tener en cuenta algunas matizaciones. Así por ejemplo, una alta consistencia interna (fiabilidad) no implica unidimensionalidad, pero la unidimensionalidad sí implica una alta consistencia interna. Green y colaboradores (1977) han demostrado que pueden darse valores altos de alpha y una estructura pluridimensional. En efecto, un coeficiente de fiabilidad alto puede darse cuando una proporción grande de la varianza está determinada por varios factores comunes, pero en situaciones en las que la unidimensionalidad puede ser muy discutible: la varianza explicada por el primer factor puede ser mínima y los factores comunes pueden no ser comunes a todos los ítems; unos factores pueden tener relación grande con unos ítems y ninguna con otros, sin que haya realmente factores comunes en todos los ítems.

En conclusión, a pesar de que el coeficiente alpha tiene el atractivo de su claridad conceptual y es el método más utilizado para hallar la fiabilidad como consistencia interna, también tiene algunas limitaciones<sup>14</sup>.

- 1/ Dependencia del número de ítems. La fiabilidad aumenta al incrementar el número de ítems y las intercorrelaciones entre los mismos. Sin embargo, existen limitaciones porque no se puede incrementar el número de ítems indefinidamente (Magnusson, 1968). Esta dificultad la reconoció el mismo Cronbach (1951) y la expresó con la siguiente analogía: un galón de leche no es más homogéneo que un cuarto de galón.
- 2/ Los presupuestos en los que se basa no se dan en la práctica. Se presupone en la teoría clásica de la fiabilidad que todos los ítems miden lo mismo y con la misma intensidad. Esto expresado matemáticamente significa que la correlación inter-ítems es idéntica a la media de todas las correlaciones y las varianzas idénticas. Pero lo cierto es que aunque los ítems midieran lo mismo, no lo hacen con la misma intensidad. Frecuentemente miden más de un constructo o dimensión y, además, unos y otros ítems, de manera desigual.

El hecho de que los ítems no midan exactamente lo mismo, ni con la misma intensidad, lo pone de manifiesto el análisis factorial, al identificar lo que podemos nominar subconstructos o factores (la pluridimensionalidad de la escala o test) y peso o correlación de cada ítem en cada factor. Por estas razones los últimos años, tanto conceptuales como metodológicos, en el estudio de la fiabilidad están acionados con el análisis factorial.

Por último, nos vamos a referir al valor mínimo que debe de tomar el coeficiente alpha y sobre lo que hay que indicar que no hay normas, ni práctica común. Nunally (1978) propone 0,70, Guilford (1954) propone un coeficiente de 0,50 para investigaciones básicas. Pfeiffer y colaboradores (1976) dan orientaciones más matizadas teniendo en cuenta el uso del instrumento, y ponen el límite en 0,85 si se van a tomar decisiones sobre individuos y de 0,60 para otros usos.

### 1.3. Elección del tipo de coeficiente de fiabilidad

La elección del tipo de coeficiente de fiabilidad a utilizar está relacionado en primer lugar con la naturaleza del rasgo o característica a medir. Así por ejemplo, las pruebas de conocimiento funciona bastante bien hallar la fiabilidad por el procedimiento de dos mitades.

Los instrumentos destinados a resolver problemas, destrezas, encontrar una ley una serie, establecer analogías, completar series, etc, hacen poco conveniente el procedimiento de la repetición e incluso, de las formas equivalentes. Sin embargo, instrumentos que miden tareas rutinarias o de atención, admiten bien los procedimientos de repetición. En este mismo caso se encuentran las características estables, como edad, sexo o estado civil.

<sup>14</sup> Amor tras reconocer las limitaciones del coeficiente alpha propone un nuevo coeficiente «The» basado en el análisis factorial. Además de este coeficiente han surgido otros como el coeficiente *g* de Heise y Bohmstedt para suprir algunos inconvenientes.

También influye en la elección del tipo de coeficiente de fiabilidad la finalidad a la que se destina. Si nuestro objetivo es predecir puntuaciones en un tiempo posterior, procede emplear la técnica de repetición.

Otra fuente que influye en la elección es la modalidad de la prueba. Así por ejemplo, si las pruebas son de velocidad hay procedimientos que dan lugar a coeficientes exageradamente altos, como es el caso de las dos mitades o de Kuder-Richardson, por lo que deben ser evitados.

A estas recomendaciones realizadas por R. Juste (1986) para la elección del procedimiento hay que añadir al menos una más: la disciplina en la que trabajamos y el tipo de variables que se utilizan. Si en psicología las herramientas de trabajo son test y escalas (Gomez-Buono, 1991) y, por tanto, variables ordinales, en Sociología, son las encuestas que recogen variables nominales y, ocasionalmente, ordinales y de intervalo. Esto revierte ciertas peculiaridades que conviene tener en cuenta, ya que no son siempre los mismos problemas los que nos ocupan a los psicólogos y a los sociólogos.

Los psicólogos se han centrado fundamentalmente en el instrumento de medición porque construyen escalas, sobre las que han de garantizar la precisión y la fiabilidad. Las construcciones que se llevan a cabo en psicología implican que todos los ítems deben de medir lo mismo (inteligencia, autoritarismo,...). En contraposición, los sociólogos no se pueden centrar tan sólo en el instrumento de medición, ya que éste es un aspecto que influye al mismo nivel que otros como son los rasgos medidos, las características de los sujetos, las épocas del año, los estados transitorios, las condiciones ambientales, las personas que intervienen en la medición, etc. Las construcciones de los ítems que se llevan a cabo en sociología tienen un valor independiente. Así, las frases que empleamos para medir una actitud no tienen la exigencia de medir todas lo mismo, sino que cada frase mide literalmente lo que se dice.

### 3.3. Relaciones entre la validez y la fiabilidad

La fiabilidad es básicamente un asunto empírico, porque se ocupa del grado de consistencia de los resultados al repetir las medidas. La validez esta orientada más teóricamente, porque ha de contestar a la pregunta ¿válido para qué propósito?. Ahora bien, estas afirmaciones conviene matizarlas porque en ocasiones es muy difícil distinguir entre validez y fiabilidad, especialmente entre la validez teórica o de constructo y la fiabilidad entendida como homogeneidad de las medidas paralelas, y también en los contrastes de «consistencia externa», donde las mediciones y los conceptos se comparan con variables externas.

Una prueba puede ser fiable y no ser válida, pero no cabe pensar en la situación opuesta. Se puede efectuar una medida con gran precisión y estabilidad, pero si mide algo diferente de lo que se pretendía medir para nada nos servirá. Ahora bien, si la medición no es fiable, no podemos medir de forma satisfactoria.

Validez y fiabilidad están íntimamente relacionadas, de forma que las ganancias o las pérdidas en una de ellas repercuten en la otra. De ello no se desprende que la relación entre ambas sea homogénea ni se pueda representar en un continuo. Niveles altos de fiabilidad no llevan necesariamente a niveles altos de validez, ya que los resultados pueden tener errores constantes. Tampoco se puede asegurar que necesariamente niveles bajos de fiabilidad impliquen niveles bajos de validez.

En efecto, por una parte los niveles bajos de fiabilidad procedente de las inconsistencias entre los jueces u observadores reduce la validez de la medida; pero, por

parte, los niveles bajos de fiabilidad procedentes de la inestabilidad en los resultados del método test-retest puede o no reducir la validez. En este último supuesto, característica que se mide es bastante estable y persistente, entonces, si se cumple que a menor fiabilidad menor validez, pero si la característica puede cambiar el tiempo, lógicamente la inestabilidad de los resultados de una administración no implica necesariamente error y por tanto, puede no afectar a la validez. También Loewinger (1954), en el caso de la relación entre la fiabilidad entendida como consistencia interna y la validez, ha establecido que el aumento de la fiabilidad a partir de cierto límite puede hacer disminuir la validez.

#### *Evolución de los términos*

El término fiabilidad comenzó a utilizarse en los test de aptitud y se suponía que consistencia o inconsistencia en los resultados de repetidas mediciones, con el uso de instrumentos similares, reflejaba la fiabilidad y, por tanto, errores de medición. La validez se concebía en relación a un sólo criterio y «se daba poca atención tanto a la influencia de las características relativamente duras, pero extrañas, como a la posibilidad de que el propio proceso de medición era modificar la característica» (Sellitz, 1976, 257).

La medida que se fueron midiendo otras características menos estables como las actitudes o las opiniones, se observó que la variable que el instrumento medía podía variar y, por tanto, las inconsistencias no se debían a errores de medición. También se observó que el instrumento podía reflejar otras características diferentes de la que se pretendía medir y en consecuencia, la consistencia de los resultados no tenía que reflejar necesariamente la ausencia de error.

El concepto de validez, desde principios del siglo XX se está escribiendo sobre la fiabilidad, especialmente de la validez de contenido y predictiva. En los años cincuenta, la American Psychological Association estableció tres tipos de validez: de contenido; criterio (concurrente y predictiva) y de constructo (Matesanz, 1975). En los años sesenta se dio a la publicación de estas recomendaciones técnicas de la A.P.A., surgen numerosos trabajos centrados sobre todo en el concepto de validez de constructo (Nunnally, 1955).

El concepto de validez de constructo nace de la limitación de la validez de contenido cuando se aplica a la medición de rasgos, actitudes, etc, es decir, a conceptos que no son directamente medibles (Nunnally, 1978). La validez de constructo resultó ser la más adecuada y por aplicación para las ciencias sociales (Zeller y Carmine, 1980) y su uso se ha operado cambios en los conceptos no solamente de validez sino también de fiabilidad.

Adicionalmente se había considerado la validez como el grado en que un instrumento «realmente» lo que se proponía medir. Actualmente el concepto es más amplio y se considera validez el grado en que una observación o medición concuerda con su referente, lo que implica dar una mayor importancia a los resultados de la medición en detrimento del instrumento con el que se mide. Lo más característico de la validez tal y como se entiende actualmente es que no se refiere a una característica del instrumento, sino de las inferencias que se hacen de las puntuaciones obtenidas. En otras palabras, no es el instrumento lo que es válido, sino sus interpretaciones.

Las recomendaciones técnicas que recoge la A.P.A. en 1974, hacen referencia a modificaciones y agrupa los tipos de validez en dos grandes bloques: los que

tienen que ver con la naturaleza o significado del atributo (validez de constructo) y los que tienen que ver con el uso de los instrumentos como indicadores de otras variables (validez de criterio).

Posteriormente, se va asumiendo paulatinamente que es la validez de constructo la que engloba a los otros tipos de validez, que se ven como pruebas de la validez de constructo o como coeficientes de utilidad, o de validez de hipótesis específicas. Esto constituye un avance conceptual importante, ya que en la actualidad la validez de constructo es casi sinónima de validez en general<sup>15</sup>.

La situación en la que nos encontramos hoy en día es muy diferente entre la validez y la fiabilidad. Mientras en la validez no existe un método claro y unificado y no se ha impuesto ninguna estrategia en particular, en la fiabilidad sí existen los métodos y las fórmulas específicas que hemos descrito, lo que facilita la comparación de los resultados y de las interpretaciones.

#### 4. Valoración final

Los análisis de validez y fiabilidad, tal y como los hemos descrito, están relacionados en su surgimiento, evolución y aplicación con los análisis cuantitativos y concretamente con el análisis de la investigación por encuesta. Ha sido en la aplicación de esta técnica donde sistemáticamente se han estudiado y desarrollado los aspectos ligados a «los problemas de la medida», y donde ha existido una preocupación constante por el control y evaluación de los errores. Sin embargo, los análisis cualitativos (entrevistas en profundidad, grupos de discusión, estudios de campo,...) no han mostrado un gran interés por estos temas, ya que se asume el subjetivismo inherente a la investigación social y, por tanto, no existen procedimientos de análisis definidos y comparables<sup>16</sup>.

En un extremo, se han sucedido las críticas internas de los cuantitativistas en sus exigencias de control, supervisión y medición de los sesgos en un intento de objetividad y objetividad, que si ya es imposible para las ciencias físicas, mucho más lo es para las ciencias sociales.

En el otro extremo, también se han sucedido las críticas externas de las perspectivas cualitativas. Estas, tienden a exagerar la significación de los «errores y sesgos», considerando superficialmente el hecho de que los científicos sociales dependen universalmente de los datos que han sido recogidos por medio de informes orales o escritos, y que estas observaciones y mediciones, con independencia de la forma en que hayan sido recogidas, se hallan sujetas invariable y esencialmente a las mismas fuentes de error y desviación que las recogidas por los entrevistadores de una encuesta.

El clínico, el experimentador o el observador participante, se hallan en el mismo peligro de desviar las respuestas de sus sujetos que los encuestadores. La diferencia fundamental es que cuando los científicos sociales tienen que depender de los informes de entrevistadores a los que han preparado y seleccionado, se dan cuenta mejor de los peligros y dificultades que van implícitos.

Podemos concluir que si por una parte, las posibilidades de controlar los errores,

<sup>15</sup> Son numerosos los autores que defienden este punto central para la validez de constructo: Messick (1980), Tenopir (1977), Anastasi (1986).

<sup>16</sup> Ambas corrientes obedecen a dos tradiciones que en sociología se conocen como «comprensión» y «explicación» y que actualmente parece que convergen mediante el pluralismo metodológico (Latiesa, 1991).

dir con precisión las características y de demostrar leyes son muy modestas, no lo, y digámoslo una vez más, hemos de abandonar la idea de intentarlo. No es electivamente ante una ciencia exacta: hay que tomar decisiones valorativas que tener en cuenta las limitaciones de los métodos para interpretar los existen diversas alternativas metodológicas y los conceptos son imprecisos, precisamente por estas mismas razones, podemos y debemos comprobar que derivaciones obtenidas con los métodos de medición son razonablemente válidas y que se encuentran suficientemente respaldadas por construcciones y por construcciones empíricas.

### Referencias bibliográficas

- ... y otros
- «Los dos métodos en las ciencias sociales», CIS.
- «Diseños de investigación», en Latiesa, M. (ed.): *El pluralismo metodológico en la investigación social*. SP Universidad de Granada.
- «Theta reliability and factor scaling», en CORNER, H. L. (ed.): *Sociological methodology 1973-1974*, Jossey-Bass, San Francisco.
- H. M.
- «The measurement problem», en BLAUROCK, H. M.; BLAUROCK, A. (ed.): *Methodology in social research*, McGraw Hill, New York, pp. 5-27.
- «What is Science», Dover, New York.
- «Psychometric theory», en DUNNETT, M. E. (ed.): *Handbook of industrial and organizational psychology*, Ran Mc Nally, Chicago.
- D. T.; FISKE, D. W.
- «Convergent and discriminant validation by the multitrait-multimethod matrix», *Psychological Bulletin*, 56, 81-105.
- E. G.; ZELLER, R. A.
- Reliability and validity assessment series: quantitative applications in the social sciences*, Sage, papers 17, London.
- KENISTON, K.
- «Yasayers and naysayers: Agreeing response set as a personality variable», *Journal of Abnormal and Social Psychology*, 60, 151-174).
- L. J.
- «Further evidence on response sets and test design», *Educational and Psychological Measurement*, 10, 3-31.
- «Response sets and test validity», *Educational and Psychological Measurement*, 6, pp. 475-494.
- Essentials of psychological testing*, Harper, New York.
- «Coefficient alpha and internal structure of tests», *Psychometrika*, 16, pp. 297-334.
- «Beyond the two disciplines of scientific psychology», *American Psychologist*, 30, pp. 116-127.
- «What price simplicity?», *Educational measurement issues and practice*, 2, 2, pp. 11-12.
- «Balancing the qualitative and the quantitative in psychological research», *Psychological assessment*, 2, pp. 3-12.
- R. MARLOWE, D.
- The approval motive: Studies in evaluative dependence*, N. Y., Wiley.
- L. T.
- Techniques of attitude scale construction*, N. Y., Appleton-Century-Crofts.
- 1978: «Analyzing one variable-three wave panel data: a comparison of two models», *Political methodology*, 5, pp. 221-231.
- FISKE, D. W.
- 1973: «Can a personality construct be validated empirically?», *Psychological Bulletin*, 80, 89-92.
- GARCÍA FERRANDO, M.
- 1985: *Sociestadística*, Alianza.
- GÓMEZ BENO, C.
- 1991: «Escala: Problemática subyacente», en LATIESA (ed.): *El pluralismo metodológico en la investigación social. Ensayos típicos*, Centro de publicaciones de la Universidad de Granada, Universidad de Granada.
- GREEN, LISSITZ, MURAIK
- 1977: «Limitations of coefficient alpha as a index of text unidimensionality», *Educational and psychological measurement*, 37, pp. 827-838.
- GUILFORD, J. P.
- 1954: *Psychometric methods*, Nueva York, McGraw-Hill.
- HEISE, D. R.
- 1969: «Separating reliability and stability in test-retest correlations», *American Sociological Review*, 34, pp. 93-101.
- HYMAN, H.
- 1971: *Diseños y análisis de las encuestas sociales*, Amorrortu.
- HYMAN, H.
- 1972: *Secondary analysis of sample surveys: Principles, procedures and potentialities*, New York, Wiley).
- KERLINGER, F. N.
- 1975: *Investigación del comportamiento*, Interamericana.
- KUDER, G. F.; RICHARDSON, M. W.
- 1973: «The theory of the estimation of test reliability», *Psychometrika*, 2, pp. 151-160.
- LINN, R. L.; WERTS
- 1979: «Covariance structures and their analysis», en TRAVIS, R.: *New directions for testing and measurement: methodological developments*, 4, pp. 53-74, Jossey-Bass, San Francisco.
- LOEVINGER, J.
- 1954: «Effect of distortions of measurement on item selection», *Educational and Psychological Measurement*, 1954, 14, 441-448).
- MAGNUSON
- 1968: *Teorías de los tests*, Trillas, México.
- MANHEIM, H. L.
- 1986: *Investigación sociológica*, Ed. Ceac, Barcelona.
- MORALES, F.
- 1981: *Metodología y teoría de la Psicología*, col. Psicología, núm. 2, F. Filosofía y Letras, Madrid.
- MORALES, P.
- 1988: *Medición de actitudes en psicología y educación*, Ed. Icaria, San Sebastián.
- NUNNALLY, J. C.
- 1964: *Educational measurement and evaluation*, McGraw Hill, New York.
- NUNNALLY, J. C.
- 1978: *Psychometric theory*, McGraw Hill, New York.
- PREIFFER, HESLUN, JONES
- 1976: *Instrumentation in human relations training*, La Jolla, University of California.
- KOHNIG, R.
- 1973: *Tratado de sociología empírica*, Madrid, Tecnos.
- REPETTO, E., y otros
- 1977: *Pedagogía experimental*, Uned, Madrid.
- RILEY, M. W.
- 1963: *Sociological research: A case approach*, Harper Brace Jovanovich, New York.
- SCOTT, W. A.
- 1960: «Measures of test homogeneity», en *Educational and Psychological Measurement*, 1960, núm 20, 751-757).