



www.universitaria.cl

Dieudo LECLERCQ



Álvaro CABRERA MARAY



UNIVERSIDAD
DE CHILE



Directores de la publicación:

Dieudonné Leclercq
Universidad de Liège (ULg)

Álvaro Cabrera Maray
Universidad de Chile (UCH)

IDEAS e INNOVACIONES
Innovaciones en Dispositivos de Evaluación
de los Aprendizajes en la enseñanza Superior
2014

Se pueden bajar gratuitamente
desde <http://orbi.uliege.be>, después Leclercq D., o
desde www.evaluaraprendizajes.cl

- Los **resúmenes** de los 23 capítulos
del libro IDEAS <http://hdl.handle.net/2268/173543>
- El **índice** de este libro para buscar entre
entradas de 1500 conceptos y
400 de autores <http://hdl.handle.net/2268/180060>

Dieudonné Leclercq

Dr. en Educación (1975) en « La Metacognición vía la autoevaluación con grados de certeza » y con postdoctorales en las universidades de Pittsburgh y UCLA. Fue profesor en las Universidades de Namur (1975-1980) y de Liège (1980-2010). Es emérito desde 2010. Enseña como invitado en las Ues. de Liège y Paris 13. Recibió el título de *Honorary Member of the World Cultural Council* (México). Ha colaborado, en Chile, con la U de Chile (UCH -Santiago), la UMCE, la UCT (Temuco), la UC del Maule, la UNAB y la UCSC (Concepción). En Perú con la PUCP y el SINEACE (Lima), la UNSAAC (Cusco) y la UNTRM (Chachapoyas). En México con la U A Chapingo. En España con la U de Sevilla y la U de Deusto (Bilbao). d.leclercq@uliege.be

Álvaro Cabrera Maray

Licenciado en Artes mención Teoría de la Música, y Master en Pedagogía en Educación Superior de la U. de Liège (Bélgica). Ha sido profesor en la Facultad de Artes y en Cursos de formación General, trabajando en el Depto. Estudios de Pregrado de la U. de Chile a cargo del Área de Formación. Integró la Red nacional de Centros de Enseñanza-Aprendizaje y la de expertos SCT-Chile sobre sistema de créditos transferibles. Trabajaba en el Ministerio de Educación de Chile, coordinando los programas de la reforma educacional en Educación Superior. alvarocabreramaray@gmail.com

Contenidos del libro IDEAS:

ES: Calificación ; Evaluación ; Productos ; Meta-cognición ; Resolución de problemas ; Proyectos ; Trabajo de grupo ; Portafolio ; Vigilancia cognitiva ; Pruebas de Progreso ; Taxonomía de Bloom ; Auto-evaluación ; Grados de certeza ; Test de Concordancia de Script ; Retroinformación ; calidades ; validez

EN : Assessment ; Evaluation ; Outcomes ; OSCE ; MCQ ; PARMs ; Metacognition ; Problem solving ; Projects ; Group produced work ; Portfolio ; Cognitive vigilance ; Progress Tests ; Bloom's Taxonomy ; Self-assessment ; Confidence Degrees ; Concordance Script Test ; Feedbacks ; Edometrics ; Metacognitive Spectral Test ; ETIC PRAD ; quality ; validity

FR : Notation ; Evaluation ; Résultats ; ECOS ; QCM ; PARMs ; Métacognition ; Résolution de problèmes ; Projets ; Travail de groupe ; Portfolio ; Vigilance cognitive ; Tests de progression ; Taxonomie de Bloom ; Auto-évaluation ; Degrés de certitude ; Test de Concordance de Script ; Rétro-information ; Edumétrie ; Test Spectral Métacognitif ; qualités d'une évaluation ; validité d'une mesure

IDEAS = Innovaciones en Dispositivos de Evaluación de los Aprendizajes en la educación Superior

La lista de los capítulos y el resumen de cada uno

aparece a continuación después de este capítulo.

CAPÍTULO IV

ETICPRAD. Ocho criterios de validez de un Dispositivo de Evaluación de los Aprendizajes (DEA)

DIEUDONNÉ LECLERCO

Las secciones A y B de este texto fueron publicadas en Leclercq D. (2006) "L'évolution des QCM". En Figari, G. y Mottier-Lopez, L. *Recherches sur l'évaluation en Education*. Paris: L'Harmattan, 139-146.

Este capítulo tiene una doble intención. En primer lugar, presentar el resumen que consideramos más completo de los criterios de calidad de un Dispositivo de Evaluación de los Aprendizajes (DEA): el modelo ETICPRAD (Leclercq, 2006). Esta sigla reúne las ocho dimensiones de la validez de un DEA, según presentamos más adelante en la sección A. La segunda intención es hacer perder su inocencia¹⁹ al lector respecto a las Preguntas de Selección Múltiple (PSM): no existe una única forma de plantear PSM. Lo anterior será ilustrado al mostrar que la evolución de las PSM es una sucesión de modificaciones que han intentado satisfacer los criterios de calidad de un DEA (y lo siguen haciendo). A continuación se presentan los 8 criterios, enumerados, definidos, marcados con un asterisco (*) y con la primera letra en mayúscula (formando el acrónimo ETICPRAD).

A. Criterios de validez de un Dispositivo de Evaluación de los Aprendizajes

- A.1. La validez *Ecológica** (Brunswick, 1943) de un DEA, o "validez aparente" (en inglés *face validity*), es mayor en la medida en que la situación de evaluación corresponde a una situación de la vida real. Se supone que la evaluación es capaz de representar o predecir las condiciones de una situación real.
- A.2. La validez *Téorica** (Cronbach y Meehl, 1955) se descompone en validez de contenido (o de "cobertura": ¿se ha testeado todo lo que debía evaluarse?) y en validez de *constructo* (¿está el DEA basado en un modelo creíble, científicamente fundado, de, por ejemplo, los procesos mentales?).
- A.3. La validez *Informativa** (diagnóstica, en lo posible) es la multiplicidad de informaciones que resultan de la evaluación, sus especificidades (apuntar sobre una

¹⁹ En el sentido que B. Bloom (1972) daba a esta expresión: nunca más poder decir "no sabía". Una persona formada debe declararse culpable.

capacidad y no sobre la vecina), sus niveles de precisión (sensibilidad), su valor explicativo.

- A.4. La validez *Consecuencial** (Green, 1998) se aprecia en las consecuencias del DEA sobre las representaciones, sobre los actos (por ejemplo: revisar o no el contenido, cambiar o no su método de estudio) de los aprendices, de los docentes y/o del sistema.
- A.5. La validez *Predictiva** de las medidas obtenidas es su capacidad de predecir eficazmente (es decir, con precisión y exactitud) otras medidas, a menudo posteriores, como por ejemplo el éxito escolar o profesional, el nivel de logro en otra prueba, etc.
- A.6. La *Replicabilidad** o *Fiabilidad* (en inglés *reliability*, en francés *fiabilité*) de una medida es su estabilidad en el tiempo o entre varios correctores. Una fórmula (Ebel, 1969) precisa el número de preguntas de una prueba, y, con Preguntas de Selección Múltiple, el número de soluciones propuestas (correctas + distractores) necesario para obtener un nivel de replicabilidad dado de la medida (0,8 por ejemplo). Otras fórmulas responden a lo mismo de la manera inversa: ¿Cuál debe ser el coeficiente n de longitud (extensión) de la prueba para alcanzar una replicabilidad deseada (por ejemplo 0,80 o 0,90)?
- A.7. La *Aceptabilidad** o practicabilidad del DEA, para el docente, se refiere a la adhesión a los principios en los que se basa la evaluación y la aplicabilidad (viabilidad) de los métodos (duración, materiales y lugares requeridos, precauciones antifraude, etc.).

Para el estudiante, la *Aceptabilidad** concierne a la adhesión y/o a la familiaridad con el tipo de prueba. Ha sido demostrado (Leclercq, 1986) que mientras más familiarizado está el/la estudiante con los procesos de *testing* y con las escalas de calificación, más avezado/a es en los tests (en inglés *test wiseness*), y más altas son sus probabilidades de éxito, especialmente con las Preguntas de Selección Múltiple.

- A.8. La validez *Deontológica** (o ética) se presenta en varias formas, siendo la equidad la más conocida. Desde hace tiempo, la docimología²⁰ crítica (Piéron, 1963) ha demostrado que la calificación (por parte de jueces) de hojas de respuesta no solo son no-concordantes inter-jueces, sino que son no-concordantes intra-juez. Las pruebas con corrección objetiva evitan otros efectos deplorables (de contraste, de severidad de juez, de halo, el efecto Posthumus, etc.), pero presentan otras limitaciones. Los derechos de los estudiantes son cada vez más y más reconocidos, y los sistemas de evaluación deben garantizar más y más la transparencia del DEA, por ejemplo en términos de re-calculabilidad de los puntos desde la hoja "en bruto".

²⁰ Ciencia de los exámenes.

B. Algunos momentos claves en la historia de las Preguntas de Selección Múltiple interpretados en términos de ETICPRAD

B.1. La gloria de la consigna clásica desde su nacimiento

Frederick J. Kelly (1915) desarrolló la primera versión de un test constituido por Preguntas de Selección Múltiple (PSM): el Kansas Silent Reading Test. Apurados por seleccionar a los oficiales entre los reclutas para la guerra de 1914-1918, EE.UU. depositó su confianza en los *Army tests* concebidos por Otis. Esos tests de inteligencia²¹, administrados a 1,7 millones de personas, eran en forma de PSM operando con la consigna clásica: "Una sola de las soluciones propuestas es correcta y se puede dar una sola respuesta". El que EE.UU. resultara estar entre los vencedores de la guerra dio una gran credibilidad a este tipo de test (Validez Predictiva*). Durante los años que siguieron las modalidades de test sistemático aumentaron aún más la exigencia por eficiencia (razón costo/eficacia) tan apreciada en EE.UU. (validez de Aceptabilidad*). Después de la guerra de 1940-1945 todavía en EE.UU., la exigencia creciente por no discriminación racial (*Civil Rights*) en la calificación hizo que se valorizara la corrección neutral de lo que llamaron los "*objective tests*" (validez Deontológica*).

El acto de añadir criterios de análisis *a posteriori* de las respuestas, para calcular los índices de discriminación (por ejemplo, vía correlaciones *point biserials*, como Davis en 1946) dio a las PSM una validez de Replicabilidad* vía la psicometría. Esos cuatro tipos de validez explican, desde nuestro punto de vista, la mayor adhesión a las PSM que existe en EE.UU. por sobre Francia, siendo los franceses pioneros de los test de coeficiente intelectual (Binet y Simon, 1905). Aunque desde 1963, en su libro *La Docimología*, Piéron viene evidenciando las discordancias importantes que pueden existir entre las calificaciones de diferentes jueces frente a una misma hoja "redactada", e incluso la inestabilidad de la calificación de un mismo juez para la misma hoja, en momentos diferentes, las autoridades francesas, estando al tanto de este problema, han mantenido la notación subjetiva con la justificación de que es ciega (los jueces no conocen la identidad del estudiante, por lo tanto las injusticias se reparten al azar según un viejo principio, francés también, de igualdad). Puede ser que la práctica más y más corriente de interponer recursos judiciales frente a decisiones de calificación (moda originada también en EE.UU.) exija que esta situación sea reconsiderada.

B.2. Un ataque teórico sobre el azar resulta en la *correction for guessing* clásica

Tversky (1964) definió el "poder" de un test ("*the power of a test*") como "1 - la probabilidad de obtener el resultado perfecto gracias al azar". Ahora bien, en cada PSM clásica que comporta k soluciones el estudiante tiene una probabilidad de $1/k$ de dar la solución correcta por azar. Varias técnicas han sido desarrolladas para paliar este defecto.

²¹ Alfa tests para los que sabían leer y Beta tests para los iletrados.

Desde 1920 Mc Call recurrió a la *correction for guessing* clásica, que consiste en establecer las tarifas como sigue: la Tarifa en caso de Respuesta Correcta (TC) vale +1 punto, la Tarifa en caso de Omisión (TOM) vale 0 y la Tarifa en caso de Respuesta Incorrecta (TI) vale $-1/(k-1)$.

Este proceso fue criticado muy poco tiempo después de su creación (West, 1923). Es la posición del autor, aun hoy, que este proceso es inadecuado, en primer lugar porque se basa en un modelo teórico falso de la actividad mental de un estudiante al contestar una PSM clásica: el primero de los tres modelos descritos por Bruce Choppin²² (1975).

En este modelo 1, cuando el estudiante “sabe” elige la respuesta correcta, y cuando “no sabe” elige al azar entre las soluciones propuestas. De donde resulta la *correction for guessing* clásica.

El modelo 2 empieza como el primero, pero en lugar de contestar al azar cuando no sabe, el estudiante empieza eliminando las soluciones que sabe que son falsas y elige entre las que quedan. Este modelo 2, del cual el 1 es una variante extrema, reconoce la noción de conocimiento parcial promocionada por De Finetti (1965). Este modelo ha inspirado consignas de tipo PSM (Preguntas con Soluciones y Respuestas Múltiples), que consisten en invitar al estudiante a eliminar las soluciones incorrectas, resultando en una PSM clásica (una sola solución correcta), con notas variando de $-(k-1)$ a $+(k-1)$. Esto produce una medida más sutil, más Diagnóstica*.

El modelo 3 de Choppin ahonda en el concepto de conocimiento parcial y dice que cuando una persona se enfrenta a una pregunta (abierta o PSM), empieza imaginando soluciones, después las ordena según un orden de probabilidad decreciente y si la consigna le fuerza a dar solo una, elige la que tiene (a sus ojos) la probabilidad más alta. Este modelo inspira a recurrir a los grados de certeza porque, como lo decía De Finetti (1965), “Solo la probabilidad subjetiva puede dar una significación objetiva a cada respuesta y a cualquier método de medición y de puntaje” (p. 111).

Al tener como referencia los modelos de Choppin es cuando aparecen las debilidades de la *correction for guessing* clásica: (1) Le falta validez Deontológica* porque es injusta: penaliza ciegamente a las personas a quienes se ha prohibido expresar su grado de duda. Además (2) le falta validez Informativa* para los docentes, porque no les muestra nada más que el número de respuestas correctas. Por fin (3), y por las mismas razones, le falta validez Consecuencial* para los estudiantes, porque excepto “omitir más frecuentemente” no tiene efecto alguno sobre su comportamiento. Cross y Fray (1977) demostraron (ver detalles en Leclercq, 1986) por qué este proceso tiene poco impacto en disuadir la “adivinación” de la respuesta por parte de los estudiantes.

²² Ver Capítulo 13 sección C.

Quienes promueven el uso de los grados de certeza²³ sostienen que esta técnica entrega respuesta a las tres fallas señaladas antes.

B.3. Un vendaval de críticas teóricas sobre los procesos mentales evaluados y no evaluados

a) Las Preguntas de Selección Múltiple clásicas no evalúan el recuerdo (evocación) de memoria.

Es evidente que las PSM no pueden pretender evaluar la capacidad de evocar conocimiento, sino de reconocerlo²⁴. En efecto, desde hace mucho tiempo se sabe (Luh, 1922) que la performance de reconocimiento tiene una tasa de éxito más alta que la de evocación. Esta observación ha sido muchas veces confirmada en contextos tan diferentes como el aprendizaje de idiomas extranjeros (Bahrick, 1984) o de medicina (Schuwirth, 1998). Añadir la solución “Ninguna” (u “Otra”) a las soluciones posibles mejora la validez Teórica*.

b) Las Preguntas de Selección Múltiple clásicas invitan a un razonamiento al revés

Aun con la solución “Ninguna” u “Otra”, las PSM inducen un proceso mental que no corresponde a lo que los estudiantes tienen que practicar en la vida normal. Con una PSM clásica, el estudiante tiene la tendencia a considerar (y eliminar) las soluciones propuestas, y solo después a elegir una solución, en lugar de primero evocar una personal. Esto corresponde al Modelo 2 de la actividad mental descrita por Choppin. Ahora bien, lo que es deseable evaluar es la capacidad del estudiante de evocar / concebir la solución, y solo después confrontarla a soluciones propuestas. Es el principio de las PSN (Leclercq, 2005) o Preguntas con Soluciones Numerosas, que permiten utilizar hojas para la lectura óptica de las respuestas en forma de marcas: el estudiante recibe una lista de centenas de soluciones arregladas en orden alfabético (como el index de un libro) entre las cuales tiene que elegir. Cada solución posible tiene un número en el orden del listado (por ejemplo de 001 hasta 700) y es con este número de tres cifras (leible por el lector óptico de marcas -LOM) que el estudiante contesta. En esta modalidad se mantienen las ventajas de la automatización de la corrección (validez de Aceptabilidad*), y al mismo tiempo se aumenta validez Teórica* (de cobertura o de contenido) a la dimensión de los conocimientos. La automatización de la corrección permite administrar muchas preguntas de este tipo (por ejemplo, 100 en una hora).

Las PSN no se justifican cuando las respuestas pueden hacerse con un teclado de computador. Hay incluso sistemas donde basta introducir las primeras letras de la palabra y el sistema propone las soluciones numerosas. Schuwirth (1998) ha aplicado este principio en tests de diagnóstico médico: si el estudiante introduce “diab”, el sistema le propone “diabetes de tipo 1, diabetes de tipo 2”, etc. Ha

²³ Ver sección B.3d y Capítulo 16.

²⁴ Ver detalles en el Capítulo 15, sección B.2.

llamado a esto *Long Menu Questions*, un sistema que mejora la validez de Aceptabilidad* para el estudiante.

c) *Las Preguntas de Selección Múltiple clásicas refuerzan el currículum oculto de la escuela.*

El currículum oculto (de la escuela) es lo que nadie enseña pero que todo el mundo aprende. Por ejemplo, que cuando hay una pregunta es obligatorio contestarla. Ahora bien, algunas preguntas, porque son absurdas o excesivamente intrusivas, ¡no deben o no pueden recibir respuesta alguna! Se aprende entonces que cuando la autoridad hace una pregunta, esta es naturalmente pertinente y bien formulada. Se aprende así que toda pregunta tiene una respuesta (correcta) y que, si uno no la sabe, no se debe reflexionar para encontrarla con razonamiento. En suma, el currículum habitual (aunque afortunadamente hay más y más excepciones) no ejercita en la vigilancia cognitiva, en la detección de las anomalías, de las incoherencias, etc., especialmente con sus modalidades de *testing*, siendo las PSM clásicas una de las más representativas. ¡Esta es una grave laguna en la validez Teórica* de esta técnica! Por todas estas razones es que hemos desarrollado (Leclercq, 1986) las PSM con Soluciones Generales Implícitas o PSM SGI. Estas soluciones son cuatro: Ninguna, Todas, Faltan datos en el enunciado de la pregunta y Absurdo en el enunciado. Son Generales porque valen (son idénticas) para todas las preguntas de un test PSM SGI. Son implícitas porque son presentadas solo una vez (al comienzo del test) y no se repiten en cada pregunta: el estudiante debe pensar en ellas solo, por su cuenta. En consecuencia, este proceso impacta positivamente en la validez Informativa* (o diagnóstica) porque permite distinguir dos niveles de la taxonomía de Bloom: la comprensión (sin trampa) y el análisis (con trampas en el enunciado de la pregunta). Gilles (1999) ha mostrado, en el ámbito de asignaturas de medicina, que las PSM SGI en las cuales la respuesta correcta es una de las cuatro SGI tienen una validez Predictiva* superior a las preguntas en las cuales la respuesta correcta es una solución simple (visible).

d) *El encuentro de las Preguntas de Selección Múltiple y los grados de certeza.*

Recurrir a grados de certeza es independiente de las PSM: aplican también para preguntas “abiertas”, que exigen respuestas redactadas. Shuford *et al* (1966), Van Naerssen (1962) y De Finetti (1965) demostraron que la consigna para indicar la seguridad en la respuesta no debe ser verbal (“poco seguro”, “medianamente seguro”, “muy seguro”) sino probabilística (en porcentajes de probabilidad). Además, hemos demostrado (Leclercq, 1982, 1993) que una precisión –un grado de discriminación sobre la propia seguridad– más alta que 20% es ilusoria, resultando nuestra consigna en 6 grados: 0%, 20%, 40%, 60%, 80%, 100% seguro/a de la respuesta.

En la línea de los autores ya citados, pensamos que este proceso tiene una mayor validez Ecológica* que los tests habituales que impiden a los estudiantes expresar sus dudas. Choppin (1975) ha descrito este problema en sus modelos 1,

2 y 3. Denunció la visión maniquea (todo o nada) de frases como “Contesten solo cuando saben, y omitan si no saben”, pues con frecuencia estamos en un estado mental de conocimiento parcial (De Finetti, 1965), en particular en situaciones de aprendizaje. El grado de duda explica comportamientos de verificación (en el diccionario, por ejemplo) como lo hemos demostrado experimentalmente (Leclercq y Gilles, 1993, p. 45).

Con las PSM, los grados de certeza resuelven (pero es un efecto secundario, no el principal) el problema del *guessing* (*adivinanza*), resultando en una contribución a la validez de Aceptabilidad* (para los docentes) de las PSM.

Al final, los grados de certeza muestran su contribución más importante en la validez Informativa* de las PSM, cuando las soluciones erróneas son elegidas más frecuentemente y con un grado de certeza más alto que la (o las) solución(es) correcta(s), lo que es anormal. Esta situación es reveladora de concepciones erróneas (en inglés *misconceptions*).

Dejamos hasta aquí esta dialéctica entre los mejoramientos sucesivos de las PSM y las críticas que continúan, ambas contribuyendo a mejorar varios aspectos de validez de las medidas. La historia no está acabada. Invitamos a quienes lo deseen a escribir unas de estas páginas en la historia de las PSM.

C. Aplicación de los criterios ETICPRAD a un DEA

C.1. Razones para concebir Dispositivos de Evaluación de los Aprendizajes con varios componentes

- Cada tipo de evaluación no puede, en solitario, evaluar todos los objetivos a un nivel aceptable, de acuerdo con los criterios de calidad o tipos de validez de ETICPRAD. Entonces, un DEA a menudo está constituido por varios eventos de evaluación, que ocurren en varios momentos, con varias finalidades, midiendo varios objetivos, con varias técnicas e instrumentos. El resultado es que el conjunto de todo ello constituye un DEA que satisface a la gran mayoría de los criterios ETICPRAD.
- Utilizar varios métodos de evaluación tiene ventajas, siendo la más importante la *validez consecucional* del conjunto. Ha sido demostrado que variar los modos de evaluación es formativo. La tabla que sigue es un resumen de la tabla original del informe de Chouinard *et al.* (2005) sobre una investigación en 106 clases (1.844 estudiantes) en escuelas primarias y 61 clases (1.719 estudiantes) en escuelas secundarias. La evaluación unimodal era utilizada en 62% de las clases de primaria y en 72% de las clases de secundaria. Los signos “-” y “+” indican cuál de los enfoques obtiene mejores resultados por cada criterio (o si se obtienen resultados similares: “=”).

Tabla 1: Comparación de enfoques unimodales y multimodales. Adaptado de Chouinard et al. (2005)

	EVALUACIÓN EN 2003	
	UNIMODAL	MULTIMODAL
Interpretación de las exigencias de la tarea	-	+
Implicación general en el trabajo	=	=
Implicación en tareas específicas	-	+
Planificación	-	+
Control	-	+
Rectificación-Ajuste	-	+
Recuerdo de lo que ha sido leído	-	+

La razón de este impacto parece evidente: hay varias evaluaciones, cada una con sus propios objetivos, apuntando a varios procesos mentales; como los alumnos estudian según las evaluaciones (lo que se llama *test driven curriculum*), esos alumnos deben entender las diferencias entre los varios tipos de evaluación y prepararse diferentemente para cada uno. Eso ejerce (y consecuentemente desarrolla) la capacidad de manejar su aprendizaje.

C.2. Cómo aplicar ETICPRAD en un DEA

La siguiente pauta se presenta vacía. Una pauta completa se encontrará en el Capítulo 5 sección D.

Tabla 2: Pauta para evaluar un Dispositivo de Evaluación de los Aprendizajes (DEA) con los ocho criterios de calidad ETICPRAD

CRITERIOS	COMPONENTES DEL DEA		
	COMPONENTE 1	COMPONENTE 2	COMPONENTE 3
DESCRIPCIÓN:			
E			
T			
I			
C			
P			
R			
A			
D			

Este tipo de descripción se puede hacer antes de la aplicación del DEA (en la fase PRE), durante (fase PER) y/o después (fase POST). En los dos últimos casos la reacción de los estudiantes puede ser conocida.

Se puede aplicar a nivel de un curso y también a nivel de un programa.

Referencias

- BAHRICK, H.P. (1984). Semantic memory content in permastore: 50 years of memory for Spanish learned in school. *Journal of Experimental Psychology: General*, 120, 1-29.
- BLOOM, B.S. (1972). L'innocence en pédagogie. *Education - Tribune Libre*, 135, 14-20.
- BRUNSWICK, E. (1943). Organismic achievement and Environment Probability. *Psychological Review*, 50, 255-272.
- CHOPPIN, B.H. (1975). Guessing the answer on objective tests. *British Journal of Educational Psychology*, 45, 206-213.
- CHOUINARD, BOWEN, CARTIER, DESBIENS, LAURIER y PLANTE (2005). Effets de différentes approches évaluatives sur l'engagement et la persévérance scolaires dans le contexte du passage du primaire au secondaire. *FQRSC*, 2003-PS-4321.
- CRONBACH, L. y MEEHL, P. (1955). Construct validity in psychological tests. *Psychological Bulletin*, 52, 281-302.
- CROSS, L. y FRARY, R. (1977). An empirical test of Lord's theoretical results regarding formula scoring of multiple choice tests. *Journal of Educational Measurement*, vol. 14, 313-321.
- DAVIS, F.B. (1946). *Analyse des items*. Louvain: Nauwelaerts (1966).
- DE FINETTI, B. (1965). Methods for discriminating levels of partial knowledge concerning a test item. *British Journal of Mathematical and Statistical Psychology*, 18, 87-123.
- EBEL, R.L. (1969). Expected reliability as a function of choices per item. *Educational and Psychological Measurement*, 29, 565-570.
- GILLES, J.L. (1999). Apports des mesures métacognitives lors d'un test de compréhension d'un article scientifique, in C. Depover y B. Noël (Eds). *Approches plurielles de l'évaluation des compétences et des processus cognitifs*. Actes de la 12^e Conférence de l'ADMEE Mons: UMH-FUGAM, 19-30.
- GREEN, D.R. (1998). Consequential aspects of the validity of achievement tests: A publisher's point of view. *Educational Measurement*, 17, 16-19, 34.
- LECLERCQ, D. (1982). Confidence marking, its use in testing, in Postlethwaite y Choppin, *Evaluation in Education*, vol. 6, 161-287, Oxford: Pergamon Press.
- LECLERCQ, D. (1993). Validity, Reliability and Acuity of Self-Assessment in Educational Testing, in Leclercq D. y BRUNO J. (1993), *Item Banking: Interactive Testing and Self-Assessment*, NATO ASI Series, F 112, Berlin: Springer Verlag, 114-131.
- LECLERCQ, D. (1986). *La conception des qcm*. Bruxelles: Labor.
- LECLERCQ, D. y GILLES, J.L. (1993). Hypermedia: Teaching Through Assessment, in D. Leclercq y J. Bruno (1993). *Item Banking, Interactive Testing and Self Assessment*. NATO ASI Series F112. Heidelberg: Springer Verlag, 31-48.
- LECLERCQ, D. (Ed.) (2003). *Diagnostic cognitif et métacognitif au seuil de l'université. Le projet MOHICAN mené par les 9 universités de la Communauté française Wallonie Bruxelles*. Liège: Editions de l'Université de Liège.
- LECLERCQ, D. (2005). *Edumétrie et docimologie pour praticiens chercheurs*. Editions de l'université de Liège.
- LECLERCQ, D. (2006) *L'évolution des qcm*. In G. Figari et L. Mottier-Lopez. *Recherches sur l'évaluation en Education*. Paris : L'Harmattan, 139-146.
- LUH, C.W. (1922). The conditions of retention. *Psychol. Monograph*, 31, 142, 401-410.
- MC CALL, W.A. (1920). A new kind of school examination. *Journal of Educational Research*, 1, 33-46.

- PIÉRON, H. (1963). Examens et docimologie, Paris: Presses Universitaires de France.
- SCHUWIRTH, L. (1998). An approach to the assessment of medical problem solving: Computerised Case-based Testing, Ph. D., Rijksuniversiteit Maastricht: Datawyse Universitaire Press.
- SHUFORD, E., ALBERT, A. y MASSENGIL, N.E. (1966). Admissible probability measurement procedures, *Psychometrika*, 31, 125-145.
- TVERSKY, A. (1964). On the Optimal Number of Alternatives at a Choice Point'. *Journal of Mathematical Psychology* 1(2): 386-391.
- VAN NAERSEN, R.F (1962). A scale for the measurement of subjective probability. *Acta Psychologica*, 20, 2, 159-166.
- WEST, P.V. (1923). A critical study of the right minus wrong method. *Journal of Educational Research*, 8, 1-9.

IDEAS E INNOVACIONES Dispositivos de Evaluación de los Aprendizajes en la educación

Dieudonné LECLERCQ y Álvaro CABRERA MARAY 2014

Resumen de cada capítulo

Los editores y autores principales del libro

p. 11-13

Prologo

Álvaro Cabrera &
Dieudonné
Leclercq

Parte 1: Conceptos clave en educación

p. 15-20

1	ATOME (Alineamiento en un Tablero de Objetivos, Métodos y Evaluaciones. Da una visión panorámica de los tres pilares de un programa de formación: los objetivos (y sus 4 niveles de alcance), los Métodos (y sus 8 Eventos de Enseñanza-Aprendizaje), las evaluaciones (y sus 4 niveles de profundidad), insistiendo sobre la Triple Concordancia (u alineamiento) O-M-E y dando ejemplos de inconsistencia.	D.Leclercq & Álvaro Cabrera p. 23-34
2	Los componentes de un dispositivo de evaluación de los aprendizajes (DEA) Da una visión de los vínculos entre las finalidades (formativas o sancionantes) de la evaluación, las competencias que desarrollar y los recursos que dominar, las condiciones de un dispositivo, las herramientas y los criterios de calidad de cada componente de un DEA.	D. Leclercq p. 35-50
3	El prisma de las características de un Dispositivo de Evaluación de los Aprendizajes (DEA) Presenta las características y las condiciones de un DEA como las facetas de un prisma: Quien (los agentes) evalúa, cuando (de manera definitiva o mejorable), quienes (individuo o grupo), para quienes (pública o confidencial), como (objetivamente o subjetivamente; estandarizada o adaptativa), que modifican la medición o su interpretación.	D. Leclercq p. 51-82
4	ETIC PRAD: Ocho criterios de validez de un Dispositivo de Evaluación de los Aprendizajes (DEA) Presenta 8 tipos de validez de un componente de un DEA: Ecológica (cerca de la situación real), Teórica (razonamiento o teoría que lo funda), Informativa (o diagnóstica), Consecuencial (lo que resulta del componente), Predictiva (correlada con otras mediciones), Replicabilidad (o fiabilidad), Aceptabilidad (para los profesores, los estudiantes, el público), Deontológica (equitativo).	D. Leclercq p. 83-92
5	Autodescribir y evaluar el Dispositivo de Evaluación de los Aprendizajes (DEA) de un curso Propone una secuencia que puede seguir un profesor para definir un DEA para su curso, es decir sus objetivos, sus métodos y sus evaluaciones, presentándoles en una tabla de modo que aparecen los vínculos y las ausencias de vínculos.	D. Leclercq & Álvaro Cabrera p. 93-102

6	<p>La calificación subjetiva de los desempeños complejos: Criterios y rubricas Presenta la docimología y sus evidencias de los efectos de notación o de calificación subjetiva (ley de Posthumus, ausencia de concordancia intra y inter-jueces, efectos de halo, de secuencia, de estereotipo, de confirmación (o de inercia). Además de esta docimología “negativa”, presenta principios de una docimología positiva y varios tipos de escalas (ej: la de Mercali) y rubricas.</p>	<p>D. Leclercq & Álvaro Cabrera p. 103-128</p>
7	<p>Evaluar la capacidad de resolver problemas Explica la diferencia entre una pregunta y un problema, el cono de la experiencia (Dale), y las heurísticas de Polya para resolver problemas. Da varios ejemplos de evaluaciones apropiadas a medir la capacidad y detectar los procesos utilizados en la resolución de problemas: las cascadas convergentes y divergentes, las análisis fraccionadas de casos (AFC), la facilitación progresiva, la medición de la búsqueda de información (Shannon, Rimoldi). Da ejemplos de medición de la creatividad, de la capacidad de aproximación y una teoría de la auto-fijación de la dificultad, como de la perseverancia.</p>	<p>D. Leclercq, S. Delcomminette (HERS) & A. Cabrera p. 129-152</p>
8	<p>ECO: Exámenes Clínicos Objetivos y Estructurados Esta técnica consiste en una sucesión de estaciones en cada de cuales se juegan roles (simulaciones) donde el profesor juega el paciente (el estudiante jugando el del medico o de la enfermera) u el cliente (el estudiante jugando el del farmacéutico), o... para medir competencias, es decir capacidad de actuar en situación compleja. El sistema de notación incluye las actitudes, las destrezas, y la cognición. Las reacciones de los participantes como la predictividad de estas mediciones son presentadas.</p>	<p>G. Philippe (ULg), D. Leclercq & J-P. Bourguignon (ULg) p. 153-170</p>
9	<p>Meta cognición y Tests Espectrales Metacognitivos (TEMs) Para los docentes que quieren desarrollar y medir capacidades como la vigilancia cognitiva, el espíritu crítico, la auto-evaluación (y la meta cognición) y el desarrollo epistemológico es presentada el método “Test Espectrales Meta cognitivos” que combina PSM con SGI (cap. 13, 14 y 15), grados de certeza (cap. 15 y 16), debate y reflexión meta cognitiva. Presenta los aspectos técnicos como los resultados obtenidos en varios ámbitos (cognitivo, epistemológico, meta cognitivo).</p>	<p>D. Leclercq & Álvaro Cabrera p. 171-196</p>
10	<p>Evaluar los Aprendizajes en la Pedagogía Por Proyectos (PPP) La PPP permite de desarrollar y medir competencias complejas (incluido trabajar en equipo), con un enfoque sobre rubricas, tan como sus componentes (recursos) en términos de cognición, actitudes, destrezas. Se puede aplicar los principios de evaluación a 360° (por los pares, por su mismo, por los docentes, por el público). El capítulo plantea (y ilustra sobre un caso) el problema de la convergencia (o ausencia de congruencia) entre estas varias fuentes de evaluación, y el problema de la ponderación de los criterios.</p>	<p>Álvaro Cabrera p. 197-220</p>
11	<p>Evaluar la contribución de cada participante a un trabajo grupal Distingue colaboración y cooperación, presenta los elementos que deben ser parte de un contrato al inicio, y después presenta 6 métodos para evaluar el valor añadido de cada participante al trabajo de grupo. Ilustra el método 4 (declaraciones de participación) con un ejemplo, el de PARMs (Proyectos de Animación Reciproca Multimedia) y sus criterios DECLAR, el método 5 (observación continua con la simulación de actividad parlamentaria y el método 6 (observar la colaboración) con la pauta de Bales. .</p>	<p>D. Leclercq, P. Gillet (ULg), M. Erpicum (ULg) & A. Cabrera p. 221-242</p>
12	<p>Los Portfolios: Hacia una evaluación más integrada y coherente con el concepto de desempeño complejo Este principio (y método) de evaluación sirve no solo a evaluar desempeños complejos como estancias en terreno, sino de constituir una integración de varias evaluaciones. Es ilustrado en dos carreras de la universidad de Liège: Formasup o Master en Pedagogía Universitaria (con sus instrucciones o consignas de redacción del portfolio) y el Master en Logopedia (que permite de discutir de 4 niveles de calidad de evidencias).</p>	<p>M. Poumay (ULg) & Chr. Maillard (ULg) p. 243-260</p>

13	<p>Las Preguntas de Selección Múltiples (PSM): del currículo escondido a la vigilancia cognitiva Presenta los retos del currículo oculto y de la espontaneidad vs la limitación a respuestas sobre sollicitación. Explica como la vigilancia cognitiva se puede entrenar y medir con una consigna valida por las PRB (Preguntas a respuesta Breve) y las PSM (Preguntas a Selección Múltiple): las Soluciones Generales Implícitas (SGI) como “Ninguna, Todas, falta datos, Absurdo”. Da una definición muy precisa de PSM, sus formas de presentación, sus ventajas y desventajas y presenta los modelos mentales que cada de 8 consignas (instrucciones) favorece. Presenta la fórmula que vincula la fiabilidad de la nota final en la prueba, el número de PSM y el número de soluciones en ella.</p>	<p>D. Leclercq & Álvaro Cabrera p. 261-286</p>
14	<p>Reglas de redacción de las Preguntas de Selección Múltiples y la habilidad para responder pruebas Presenta 24 reglas (repartidas en 5 categorías) y los dispositivos experimentales (preguntas sobre contenidos ficticios) que permiten verificarlas, tan como los resultados de estas verificaciones en caso de transgresión de las reglas.</p>	<p>D. Leclercq p. 287-300</p>
15	<p>Evaluar procesos cognitivos según la Taxonomía de Bloom Presenta modalidades de evaluación apropiadas a cada de los 6 niveles de los procesos mentales descritos en la taxonomía de Bloom: la memoria (de re-cognición y de evocación), la comprensión (con la definición de Smedslund), la aplicación, el análisis (y las Preguntas PRIM-BIS para diferenciar entre análisis y comprensión, la síntesis y la creación (y los criterios de Torrance), el juicio(incluido la capacidad de aproximar).</p>	<p>D. Leclercq p. 301-328</p>
16	<p>Auto-evaluación con grados de certeza: un microscopio para la evaluación de los aprendizajes Presenta los retos del uso de grados de certeza: epistemológico (de definición de “dominio”), de medición en investigación (la necesidad de un microscopio del pensamiento), de caracterización practica (utilizable – inutilizable) de niveles de conocimiento) y de fijación de umbrales de éxito os resultados y de excelencia. Presenta las condiciones metodológicas de uso (3 principios), las distribuciones espectrales de calidad de les respuestas, las nociones de meta memoria y de meta comprensión (el JOC o juicio de comprensión).</p>	<p>D. Leclercq p. 329-356</p>
17	<p>Grados de certeza y docimología: como calificar Denuncia varios sistemas de cotejo inapropiados y la importancia (impredecible) de tener en cuanta el realismo de las respuestas acertadas por un estudiante en una prueba. Explica como verificar (con la ley binomial) la presunción de realismo, cálculo de un índice de calibración. Trata de la sobrestimación y de resolución (Discriminación y lucidez), tan como de una pauta innovadora de cotejo basada en ;los grados de certeza.</p>	<p>D. Leclercq p. 357-386</p>
18	<p>PdP: Pruebas de Progreso Presenta una modalidad de evaluación en cual la universidad de Maastricht se ha ilustrada como pionera: la Pruebas de Progreso que consisten en presentar el mismo día a todos los estudiantes de una carrera (que sean de primer o de ultimo año) una prueba sobre todos los contenidos de la carrera (centenas de preguntas), cuatro veces por año (con pruebas “paralelas”). Las ventajas y desventajas son revisitadas, como el modo de comunicar los resultados, original también. Estos principios son ilustrados por su aplicación en Maastricht desde cuarenta años.</p>	<p>D. Leclercq, A. Cabrera & C. Van der Vleuten (U. Maastricht) p. 387-408</p>
19	<p>TCS : El Test de concordancia de Script Esta técnica ha sido concebida para medir la capacidad clínica de tratar la información. Ha sido utilizada principalmente en medicina (revisión de opinión desde una información adicional). Es ilustrada con un ejemplo y resultados de su aplicación en la univ. de Liège.</p>	<p>V. Massart (ULg), A. Collard (ULg) D. Giet (ULg) p. 409-418</p>

20	<p>Concebir Dispositivos de Evaluación de los Aprendizajes (DEA) al nivel de un programa Presenta tres experiencias de desarrollo de un DEA al nivel de una facultad: la de Farmacia en Liège y las de medicina en Liège y en Maastricht.</p>	<p>D. Leclercq, C. Van der Vleuten & A. Cabrera p. 419-430</p>
21	<p>Retroinformaciones (Feedbacks) Empieza con el problema de la profundidad de penetración de una retroinformación, desde sobre los detalles de ejecución de la tarea hasta el <i>Self</i> (es porque son presentadas las teorías de William James sobre la auto-estima y la <i>FIT</i> o <i>Feedback Intervention Theory</i>). Un modelo integrador (llamado CAIRO) es presentado. Varios modos de presentación de las retroinformaciones después de una prueba son presentados. Una modalidad, utilizada en la UCH (Universidad de Chile) que se focaliza al esencial, es presentada con un ejemplo.</p>	<p>D. Leclercq, M. de la Fuente (UCH) & A. Cabrera p. 431-454</p>
22	<p>Los roles de un SMART: Servicio Metodológico de Apoyo a la Realización de Tests Un (SMART) ayuda docentes en la concepción y la realización de pruebas estandarizadas y en el procedimiento de las respuestas de los estudiantes (calculo de varios índices relativos a cada pregunta y cada solución de las PSM), como en las retroinformaciones automatizadas a los estudiantes. Un enfoque especial es dedicado al uso de cajas de voto a distancia (<i>clickers</i>).</p>	<p>D. Leclercq & P. Detroz (ULg) p. 455-476</p>
23	<p>Índices cuantitativos en Docimología Consiste en un catálogo de conceptos útiles para tratar cuantitativamente los datos resultando de evaluaciones estandarizadas como</p> <ul style="list-style-type: none"> -los tipos de categorías (nominales, ordinales, métricas). -los índices relativos a una distribución : índices de centración (Modo, Mediana, Media), de dispersión (rango, cuartiles, desviación estándar), de posiciones relativas o normativas (la nota z, los percentiles) de la forma de la distribución (asimetría o <i>skewness</i>). -las presentaciones gráficas de distribuciones. -índices de comparación o de progreso: la amplitud del efecto (AE), la ganancia relativa (GR). -la fiabilidad de la nota (<i>reliability</i>) al total de la prueba y el alfa de Cronbach. -el umbral de éxito, fijado a priori o a posteriori. -el índice de discriminación (correlación punto <i>biserial</i> o <i>rpbis</i>) de un modo de respuesta aplicado a cada de las soluciones de cada PSM -el análisis automática de una prueba -el valor heurístico de los nubes de puntos. 	<p>D. Leclercq, R. Roco (Chile) & A. Cabrera p. 477-543</p>
24	<p>Index de los autores 426 autores citados.</p>	<p>D. Leclercq & A. Cabrera p. 545-549</p>
25	<p>Index de los conceptos Se puede bajar gratuitamente via http://hdl.handle.net/2268/180060</p>	<p>D. Leclercq & A. Cabrera</p>