Bibliography

- American Educational Research Association (AERA), American Psychological Association (APA), & National Council on Measurement in Education (NCME). (1999). Standards for educational and psychological testing. Washington, DC: AERA.
- Angoff, W. H. (1971). Scales, norms and equivalent scores. In R.L. Thorndike (Ed.), *Educational measurement* (pp. 508–600). Washington, DC: American Council on Education.
- Bloom, B., Englehart, M., Furst, E., Hill, W., & Krathwohl, D. (1956). *Taxonomy of educational objectives: The classification of educational goals. Handbook I: Cognitive domain.* New York, Toronto: Longmans, Green.
- Center for Applied Special Technology (CAST). (2002). Universal design for learning (UDL) guidelines version 1.0. Retrieved from http://www.cast.org/library/UDLguidelines/version1.html.
- Cook, H. G. (2008, June). *Expansion of Webb's alignment protocol*. Paper presented at the 2008 Council of Chief State School Officers Student Assessment Conference, Orlando, FL.
- Crocker, L., & Algina, J. (2006). *Introduction to classical and modern test theory.*Belmont, CA: Wadsworth Publishing Company.
- Davies, S., O'Malley, K., & Wu, B. (2007, April). *Establishing measurement equivalence of transadapted reading and mathematics tests*. Paper presented at the 2007 annual meeting of the American Educational Research Association, Chicago, IL.
- Efron, B. (1979). Bootstrap methods: Another look at the Jackknife. *The Annals of Statistics*, 7(1), 1–26.
- Efron, B., & Tibshirani, R. J. (1993). *An introduction to the bootstrap*. New York, NY: Chapman and Hall.
- Ferrara, S., Lewis, D., Mercado, R., D'Brot, J., Barth, J., & Egan, K. (2011, April). *A method for setting benchmarked performance standards: Workshop procedures, panelist judgments, and empirical results*. Paper presented at the 2011 annual meetings of the National Council on Measurement in Education, New Orleans, LA.





- Flowers, C., Wakeman, S.Y., Browder, D. M., & Karvonen, M. (2007). Links for Academic Learning: An alignment protocol for alternate assessments based on alternate achievement standards.

 Charlotte, NC: National Alternate Assessment Center, University of North Carolina at Charlotte.
- Hambleton, R. K., & Plake, B. S. (1995). Using an extended Angoff procedure to set standards on complex performance assessments. *Applied Measurement in Education*, 8, 41–56.
- Kane, M.T. (1992). An argument-based approach to validity. *Psychological Bulletin*, 112, 527–535.
- Kane, M.T. (2006). Validation. In R.L. Brennan (Ed.), Educational measurement (pp. 17–64). Washington, DC: The National Council on Measurement in Education & the American Council on Education.
- Kolen, M. J., & Brennan, R. L. (2004). *Test equating, scaling, and linking: methods and practices*. New York, NY: Springer.
- Lewis, D. M., Green, D. R., Mitzel, H. C., Baum, K., & Patz, R. J. (1998, April). *The bookmark standard setting procedure: Methodology and recent implementations*. Paper presented at the 1998 annual meeting of the National Council on Measurement in Education, San Diego, CA.
- Linacre, J. M. (2001). A user's guide to WINSTEPS: Rasch-model computer program. Chicago, IL: MESA Press.
- Lohr, S. L. (1999). *Sampling: Design and Analysis*. Pacific Grove, CA: Duxbury Press.
- Masters, G. N. (1982). A Rasch model for partial credit scoring. *Psychometrika*, 47(2), 149–174.
- Messick, S. (1989). Meaning and values in test validation: The science and ethics of assessment. *Educational Researcher*, 18(2), 5–11.
- Morgan, D. L. (2003, November). The performance profile method: A unique method as applied to a unique population. Paper presented at the 2003 annual meeting of the California Educational Research Association, San Francisco, CA.

- National Alternate Assessment Center. (2005, June). Designing from the ground floor: Alternate assessment on alternate achievement standard. Paper presented at the 2005 Council of Chief State School Officers Large Scale Assessment Conference, San Antonio, TX.
- O'Malley, K., Keng, L., & Miles, J. (2012). Using validity evidence to set performance standards. In G.J. Cizek (Ed.), *Setting performance standards* (pp. 301–322). New York, NY: Routledge.
- Petersen, N. S. (1987, September 25). *DIF procedures for use in statistical analysis* [ETS internal memorandum].
- Petersen, N. S., Kolen, M. J., & Hoover, H. D. (1989). Scaling, Norming, and Equating. In R.L. Linn (Ed.), *Educational Measurement* (pp. 221–262). New York, NY: Macmillan.
- Phillips, G.W. (2012). The benchmark method of standard setting. In G. Cizek (Ed.), Setting performance standards (pp. 342–364). New York, NY: Routledge.
- Rasch, G. (1966). An individualistic approach to item analysis. In P. Lazarfeld & N. W. Henry (Eds.), *Readings in Mathematical Social Science* (pp. 89–107). Chicago, IL: Science Research Associates.
- Rudner, L. M. (2001). Computing the expected proportions of misclassified examinees. *Practical Assessment, Research & Evaluation*, 7(14).
- Rudner, L. M. (2005). Expected classification accuracy. *Practical Assessment, Research & Evaluation*, 10(13).
- Schafer, W. D., Wang, J., & Wang, V. (2009). Validity in action: State assessment validity evidence for compliance with NCLB. In R. W. Lissitz (Ed.) *The concept of validity* (pp. 173–193). Charlotte, NC: Information Age.
- Shepard, L. A. (1997). The centrality of test use and consequences for test validity. *Educational Measurement: Issues and practice*, 16(2), 5–8, 13, 24.
- Torgerson, W. S. (1958). *Theory and methods of scaling*. New York, NY: Wiley.
- Way, W. D., Davis, L. L., & Fitzpatrick, S. (2006, April). Score comparability of online and paper administrations of the Texas Assessment of Knowledge and Skills.

 Paper presented at the 2006 Annual Meeting of the National Council on Measurement in Education, San Francisco, CA.
- Webb, N. L. (1997). Criteria for Alignment and Expectations and Assessments in Mathematics and Mathematics Education. Council of Chief State School Officers and National Institute for Mathematics Education Research Monograph No. 6. Madison: University of Wisconsin, Wisconsin Center for Educational Research.





- Wingersky, M.S., & Lord, F.M. (1984). An investigation of methods for reducing sampling error in certain IRT procedures. *Applied Psychological Measurement*, 8(3), 347-364.
- Wright, B. D. (1977). Solving measurement problems with the Rasch model. Journal of Educational Measurement, 14, 97–116.
- Wright, B. D., & Masters, G. N. (1982). *Rating scale analysis*. Chicago: MESA Press.
- Wright, B.D., & Stone, M.H. (1979). Best test design. Chicago: MESA Press.
- Zieky, M. (1993). DIF statistics in test development. In P.W. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 337–347). Hillsdale, NJ: Erlbaum.