

ZUR LEISTUNGSFÄHIGKEIT VON RASCH MODELL UND KLASSISCHER TESTTHEORIE FÜR DIE INFERENZSTATISTISCHE EVALUATION VON TESTSCORES

Kempf, Wilhelm (Konstanz)

Entgegen der verbreiteten Auffassung von der universellen Anwendbarkeit der klassischen Testtheorie (KT) ist der von ihr angegebene Vertrauensbereich des True-Scores

$$(1) \text{KONF}_\gamma \{x_{vt} - z_{\text{krit}}\sigma(F.t) \leq \tau_{vt} \leq x_{vt} + z_{\text{krit}}\sigma(F.t)\}$$

an drei Voraussetzungen geknüpft: (A1) an jene empirischen Zusatzannahmen, welche für die Reliabilitätsschätzung erforderlich sind, aus der sich der Standardmeßfehler errechnet, (A2) an die näherungsweise Übereinstimmung des Standardmeßfehlers $\sigma(F.t)$ mit der tatsächlichen Standardabweichung des Meßfehlers der Person $\sigma(F_{vt})$ und (A3) an die Normalverteilung des Meßfehlers.

Untersucht man, wie der Testscore X_{vt} aus den einzelnen Antworten A_{vi} einer V_p auf die k Items eines Tests zusammengesetzt ist, so zeigt sich, daß bei binären Items und unter der verbreiteten Scoreformel $X_{vt} = \sum_i A_{vi}$, nur die letzte dieser Voraussetzungen als einigermaßen unproblematisch gelten kann. Wegen des zentralen Grenzwertsatzes kann Verteilung des Meßfehlers bei großem k und mittlerem τ_{vt} durch die $N[0, \sum_i p_{vi}(1-p_{vi})]$ approximiert werden. Annahme (A1) ist dagegen hoch restriktiv und stellt strenge Anforderungen an die Testkonstruktion, während Voraussetzung (A2) wegen der Abhängigkeit der Fehlervarianz $\sigma^2(F_{vt}) = \sum_i p_{vi}(1-p_{vi})$ vom True Score τ_{vt} jedenfalls verletzt ist. Folglich kann die KT die Konfidenzgrenzen des True-Scores nicht einmal im Item-Sampling-Paradigma (ISP) korrekt wiedergeben. Die Breite des Vertrauensbereiches wird von der KT im mittleren Scorebereich unterschätzt und für großes oder kleines τ_{vt} überschätzt. Hier nehmen die Konfidenzgrenzen zudem unsinnige Werte an, welche außerhalb des Wertebereichs der Scorevariable liegen. Wie schlecht die Ergebnisse ausfallen, hängt von der Homogenität der Eichstichprobe ab, aus welcher die Reliabilität bestimmt wurde: je heterogener sie ist, desto stärker wird $\sigma(F_{vt})$ im mittleren Scorebereich unterschätzt. Bei Anwendung des Rasch Modells (RM) ergibt sich der Vertrauensbereich des True Scores durch Einsetzen der Konfidenzgrenzen des Personenparameters in $E(X_{vt}) = \sum_i f_i(\zeta)$. Dabei werden die Konfidenzgrenzen im ISP durch das RM (mit $\sigma_i = 1$) ebenfalls nicht perfekt, aber doch mit sehr guter Näherung rekonstruiert.

Entsprechend führt auch der Signifikanztest für den Unterschied zweier Testscores mittels der Prüfgröße $z = (X_{vt} - X_{wt}) / [\sigma(F.t) / \sqrt{2}]$ bereits im ISP zu einer eklatanten Überbewertung der Scoredifferenzen durch die KT, während die Prüfgröße

$z = (\zeta_v - \zeta_w) / \sqrt{[I_t(\zeta_v)^{-1} + I(\zeta_w)^{-1}]}$ des RM die korrekten Verhältnisse deutlich besser wiederzugewinnen vermag.

Indem sie von unabhängigen Testscores und damit von unabhängigen Itemstichproben ausgehen, sind beide Methoden nur im ISP anwendbar und bewirken im Fixed-Test-Paradigma (FTP) eine noch krassere Überbewertung der Score Differenzen. Ein angemessenes Prüfverfahren stellt im FTP der Test von McNemar dar, dessen Anwendung zudem an keine Modellannahmen gebunden ist. Geltung des RM ist allerdings notwendig und hinreichend für die spezifische Objektivität des Tests, so daß die Verteilung der Prüfgröße ausschließlich vom Unterschied zwischen den Vpn abhängt.

Kempf, W.: Ein pragmatischer Ansatz in der statistischen Theorie psychologischer Testscores. Universität Konstanz (im Druck).