

مقایسه خطای استاندارد اندازه‌گیری شرطی در روش‌های غیرخطی تبدیل نمره‌های خام به نمره‌های مقیاس

مجتبی جهانی فر^{۱*}، ابراهیم خدایی^۲، جلیل یونسی^۳، سید امین موسوی^۴

تاریخ دریافت: ۹۵/۱۱/۲۸

تاریخ پذیرش: ۹۶/۰۷/۱۵

چکیده

برای تفسیرپذیری بهتر و مقایسه پذیر کردن نمره‌ی آزمون‌ها با همدیگر، نمره‌های خام به دست آمده از خرده‌آزمون‌ها را به مقیاس مشترکی تبدیل می‌کنند که به آن نمره مقیاس گفته می‌شود. روش‌های متفاوت خطی و غیرخطی برای تبدیل نمره‌های خام به نمره‌های مقیاس وجود دارد. روش‌های متداول غیرخطی تبدیل نمره‌های خام به نمره‌های مقیاس، نرمال‌سازی و تبدیل آرک‌سینوس هستند. در این پژوهش که با هدف مقایسه خطای استاندارد اندازه‌گیری تبدیل نرمال و تبدیل آرک‌سینوس انجام گرفته از ۱۰۰۰۰ داده شبیه‌سازی شده و ۱۰۰۰۰ داده واقعی آزمون سراسری سال ۱۳۹۵ ایران بهره برده‌ایم. به منظور مقایسه این دو روش از نمودارها و شاخص‌های آماری و همچنین ویژگی‌های اندازه‌گیری بر اساس نظریه نمره قوی حقیقی از جمله خطای استاندارد اندازه‌گیری شرطی استفاده شد. نتایج نشان داد که دو روش تبدیل از ویژگی‌های متفاوتی برخوردار هستند. گرچه نمره‌ها در هر دو روش دارای ضریب پایایی بالایی هستند، اما روش آرک‌سینوس ضمن کاهش نوسان خطا برای سطوح مختلف نمره‌ها، دارای میانگین خطای استاندارد اندازه‌گیری شرطی کمتری نسبت به روش نرمال‌سازی بوده است.

واژه‌های کلیدی: تبدیل آرک‌سینوس، خطای استاندارد اندازه‌گیری شرطی، نرمال‌سازی، نمره‌های

مقیاس

۱. * دانشجوی دکتری سنجش آموزش، دانشگاه تهران، تهران، ایران. m.jahanifar@ut.ac.ir

۲. دانشیار گروه روش‌ها و برنامه‌های آموزشی و درسی، دانشگاه تهران، تهران، ایران

۳. دانشیار گروه سنجش و اندازه‌گیری، دانشگاه علامه طباطبائی، تهران، ایران.

۴. عضو هیئت علمی دانشکده آموزش، دانشگاه ساسکاچوان، کانادا.

مقدمه

هر آزمون مرکب شامل چند خرده‌آزمون است که هر کدام نمره خاصی را به خود اختصاص می‌دهند، برای تفسیرپذیری بهتر و مقایسه‌پذیر کردن نمره‌ی آزمون‌ها، نمره‌های خام^۱ را به نمره‌های مقیاس^۲ تبدیل می‌کنند. در صورتی می‌توان نمره خرده‌آزمون‌های مختلف را باهم جمع کرد که روی یک مقیاس مشترک باشند. روش‌های مقیاس‌سازی می‌توانند موجب کاهش اطلاعات و یا ایجاد خطاهای تصادفی و نظام‌دار در نمره‌های مقیاس شوند و نمره‌های هر حوزه را تحت تأثیر قرار دهند. یکی از قدیمی‌ترین روش‌های پیوند زدن^۳ نمره‌ها مقیاس‌سازی نمره‌ها است. مقیاس‌سازی همان تبدیل نمره‌های خام به نمره‌های مقیاس است، در روان‌سنجی همواره به مقیاس‌های متفاوت نمره‌ها نیاز بوده به همین خاطر قدمت این روش به قدمت علم روان‌سنجی است. هدف از مقیاس‌سازی نمره‌ها تبدیل نمره‌ها از دو آزمون مختلف به مقیاس مشترک^۴ است. این روش پیوند غیرمستقیم بین نمره‌های دو آزمون ایجاد می‌کند (دورانز، پامریچ و هولاند، ۲۰۰۷). در این روش نمره‌های هر آزمون به‌طور جداگانه به مقیاس مشترک برده می‌شوند. هرگاه دو یا چند آزمون که سازه‌های متفاوتی را اندازه می‌گیرند با هم روی یک جامعه‌ی مشترک از شرکت‌کنندگان اجرا شوند، مقیاس‌سازی بر روی هریک از آزمون‌ها به‌منظور داشتن مقیاس مشترک انجام خواهد گرفت. کولن (۲۰۱۴) این آزمون‌ها را آزمون مرکب^۵ نامید. نکته مهم در اینجاست که بردن نمره‌ها به مقیاس مشترک به معنی هم‌تراز کردن آن‌ها نیست و نمی‌توان نمره‌ی آزمون با ساختار مشخص را با آزمون دیگر با ساختار متفاوت تبدیل کرد و مقیاس مشترک تنها آن‌ها را مقایسه‌پذیر می‌کند نه این که راهی را فراهم کند که بتوان آن‌ها را به‌جای هم به کار برد. برای هر آزمون مانند Y که بخواهد مقیاس‌سازی شود، تابع توزیع تراکمی نمره‌ها وجود دارد که روی جامعه مرجع P به‌صورت رابطه (۱) تعریف می‌شود.

1. raw scores
2. scale scores
3. linking
4. common scale
5. test battery

$$F_{YP}(y) = \Pr(Y \leq y | P) \quad (۱)$$

در رابطه (۱) تابع توزیع تراکمی و $\Pr(Y \leq y | P)$ تابع توزیع احتمال هستند.

نمره‌های آزمون Y به وسیله تبدیل S به نمره‌های مقیاس تبدیل می‌شوند.

$$S = S(F_{YP}(y)) \quad (۲)$$

در رابطه‌ی (۲)، S تابع تبدیل به مقیاس است که می‌تواند خطی و یا غیرخطی باشد. از مزیت‌های این روش، تبدیل نمره خرده‌آزمون‌های یک آزمون مرکب به مقیاسی است که تفسیرپذیری و مقایسه آن‌ها را با همدیگر آسان‌تر می‌کند. نمره بالا در مقیاس مشترک برای یک آزمودنی در یک خرده‌آزمون به معنی عملکرد بهتر وی در آن خرده‌آزمون نسبت به سایر خرده‌آزمون‌ها است (دورانز و همکاران، ۲۰۰۷). تبدیل‌های مقیاس خطی^۱، تبدیل مقیاس نرمال^۲ و تبدیل آرک‌سینوس^۳ از مشهورترین روش‌های تبدیل مقیاس است، دو روش آخر را روش‌های تبدیل غیرخطی نمره‌های خام به نمره‌های مقیاس می‌نامند.

در روش معمول نرمال‌سازی ابتدا با استفاده از توزیع فراوانی نمره‌های کسب‌شده توسط تمامی شرکت‌کنندگان در آزمون، فراوانی تراکمی و تراکمی نسبی هر نمره محاسبه شده و سپس رتبه درصدی برای هر نمره مشخص می‌گردد. با استفاده از رتبه درصدی، نمره z متناظر با رتبه درصدی را از معکوس تابع توزیع تراکمی نرمال مطابق با رابطه (۳) محاسبه می‌کنند.

$$\Phi(z) = \frac{\hat{Q}(y)}{100} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\omega^2/2} d\omega \quad (۳)$$

تابع رابطه (۳) توزیع تراکمی نرمال استاندارد است. ω متغیر انتگرال‌گیری است که دامنه آن از $-\infty$ تا z است و $\hat{Q}(y)$ همان رتبه درصدی است (آنگاف، ۱۹۷۱ و کولن، ۲۰۱۴). به‌عنوان نمونه آزمون SAT^۴ برای تبدیل نمره‌های خام به نمره‌های مقیاس از تبدیل

1. linear transformation (LT)
2. normalizing transformation (NT)
3. arcsine transformation (AT)
4. scholastic aptitude test (SAT)

نرمال بهره می‌برد. این آزمون که سه ساعت و چهل و پنج دقیقه به طول می‌انجامد شامل خرده‌آزمون‌های مهارت‌های خواندن، نوشتن و همچنین ریاضیات است. در این تبدیل نمره‌های فرمولی^۱ برحسب فراوانی تراکمی، به مقیاس نرمال برده می‌شوند و حاصل این مقیاس‌سازی جدول تبدیلی است که در آن هر نمره بر اساس رتبه و نمره درصدی به مقیاس آزمون SAT برده می‌شود (گزارش فنی SAT، ۲۰۱۵). آزمون IBTS که در دانشگاه آیووا به منظور اندازه‌گیری پیشرفت تحصیلی از سطح مهد کودک تا پایه دوازدهم طراحی شده است، دارای ۱۵ خرده‌آزمون است. البته این خرده‌آزمون‌ها برای پایه‌های نهم تا دوازدهم شامل نه خرده‌آزمون است. این خرده‌آزمون‌ها شامل مهارت‌های خواندن و نوشتن، مهارت‌های دستور زبان، مهارت‌های شنیداری، مهارت‌های محاسباتی و ریاضیات، علوم و مهارت‌های اجتماعی است؛ که هر کدام هم در تعداد سؤال‌های خرده‌آزمون و هم در زمان پاسخگویی متفاوت هستند. در این آزمون نیز از روش تبدیل نرمال برای ساختن مقیاس نمره‌ها استفاده شده است. (گزارش فنی IBTS، ۲۰۱۶). آزمون سراسری ورود به دانشگاه‌های ایران نیز آزمونی مرکب است که با توجه به رشته تحصیلی دارای خرده‌آزمون‌های متفاوت است، در این آزمون نیز پس از محاسبه نمره خام، نمره‌ها به روش تبدیل نرمال به مقیاس مشترک برده می‌شوند. نمره‌های مقیاس در آزمون سراسری ایران به نمره‌های تراز شهرت دارند (نقی زاده، ۱۳۹۴)

روش دیگر تبدیل خطی نمره‌های خام به نمره‌های مقیاس روش تبدیل آرک سینوس است. در تبدیل نمره‌های خام به نمره‌های مقیاس با استفاده از رابطه تبدیل آرک سینوس از تعداد پاسخ‌های درست به هر خرده‌آزمون و تعداد سؤال‌های خرده‌آزمون برای ساختن نمره‌ی مقیاس استفاده می‌شود (کولن، ۲۰۱۴).

$$S(X_i) = \frac{1}{2} \left\{ \sin^{-1} \sqrt{\frac{X_i}{k+1}} + \sin^{-1} \sqrt{\frac{X_i+1}{k+1}} \right\} \quad (۴)$$

در رابطه (۴)، X_i نمره‌ی خام (تعداد پاسخ‌های صحیح) و k تعداد سؤال‌های خرده‌آزمون هستند. این روش مقیاس‌سازی، نمرات خام با استفاده از تابع تبدیل نمرات

1. formula score

آرک‌سینوس به‌وسیله کولن و هانسون در سال ۱۹۸۹ پیشنهاد شده است. کولن و هانسون (۱۹۸۹) در این پژوهش نشان دادند که در صورت استفاده از این مقیاس خطای نمره‌ها با هم تا حدود زیادی همسان خواهند شد.

آزمون ACT^۱. از چهار خرده‌آزمون زبان و ادبیات انگلیسی، ریاضیات، مهارت‌های خواندن و علوم تشکیل شده است. همگی این خرده‌آزمون‌ها به‌صورت چندگزینه‌ای طراحی شده‌اند. مهارت نوشتن هم به‌عنوان یک خرده‌آزمون اختیاری و به‌صورت انشائی در این آزمون گنجانیده شده است. در این آزمون نمره خام هر یک از خرده‌آزمون‌ها به‌طور مستقیم به مقیاس نمره‌ای آرک‌سینوس بین ۱ تا ۳۶ تبدیل می‌شوند (راهنمای ACT، ۲۰۱۴).

برای ساختن نمره مقیاس از روی نمره‌های خام روش‌های متفاوتی وجود دارد به‌گونه‌ای که موسسه‌های مختلف تولید آزمون از مقیاس‌های مختلفی برای تبدیل نمره‌ها بهره می‌برند، اما به‌طور قطع و یقین نمی‌توان اظهار داشت که کدام مقیاس بهترین است. به همین خاطر انتخاب بهترین مقیاس کار دشواری است. مفید بودن هر مقیاس به ویژگی‌هایی است که آن مقیاس برای معنادار بودن نمره‌ها ایجاد کرده و همچنین تمهیداتی است که برای کاهش کج‌فهمی‌ها در تفسیر نمره‌ها ایجاد می‌کند (کولن، ۲۰۱۴).

شون ون چانگ^۲ سه روش تبدیل مقیاس یعنی خطی، نرمال‌سازی و تبدیل آرک‌سینوس را با هم مقایسه کرده است. در این مقایسه، سه روش مقیاس‌سازی در شاخصه‌هایی مانند ضریب پایایی نمره‌های مقیاس و تعداد تغییر نمره‌ها در اثر برش نمره‌ها^۳ و همچنین شکاف‌های^۴ ایجادشده در مقیاس با هم مقایسه شده‌اند. در نتیجه این گزارش چنین آمده است که هر روش دارای معایب و مزایای مربوط به خودش است و هیچ روش همه ویژگی‌های مطلوب را دارا نیست، تصمیم برای انتخاب مقیاس مناسب بستگی هم به خواص اندازه‌گیری و هم به سهولت تفسیر خواهد داشت (چانگ، ۲۰۰۶). در پژوهش چانگ ضریب پایایی خرده‌آزمون‌ها پس از انجام سه نوع تبدیل به هم نزدیک بودند. درحالی که در ایجاد

1. american college test
2. Shun-Wen Chang
3. truncation
4. gaps

شکاف‌های بین نمره‌ها روش تبدیل آرک‌سینوس از دو روش دیگر پیشی گرفته است ولی در نمودار خطای استاندارد شرطی این روش از دو روش دیگر دارای خطای کمتری است و روش نرمال‌سازی بیشترین شباهت را به توزیع نمره‌های خام نشان داده است؛ و تبدیل خطی نمره‌ها کمترین شکاف را در مقیاس نمره‌ها نشان داده است. با توجه به منابع و بررسی‌های متفاوت نمی‌توان قطعاً گفت که کدام روش نسبت به سایر روش‌ها ترجیح بیشتری دارد، هر کدام از تبدیل‌ها دارای ویژگی‌های مخصوص به خود هستند و مفید بودن هر روش مقیاس‌سازی بستگی به مفید بودن و تفسیرپذیری بهتر آن داشته و هر روش کاربردهای متفاوتی پیدا کرده است.

علاوه بر پژوهش چانگ (۲۰۰۶) پژوهش‌های کولن و هانسون (۱۹۸۹) کولن و برنان (۲۰۱۴) پترسن، کولن و هاور (۱۹۸۹) نشان دادند که به‌طورقطع و یقین نمی‌توان اظهار داشت که کدام مقیاس بهترین است. به همین خاطر انتخاب بهترین مقیاس کار دشواری است. مفید بودن هر مقیاس به ویژگی‌هایی است که آن مقیاس برای معنادار بودن نمره‌ها ایجاد کرده و همچنین تمهیداتی است که برای کاهش کج‌فهمی‌ها در تفسیر نمره‌ها ایجاد می‌کند بستگی دارد ضمن اینکه بررسی‌ها نشان دادند نمره‌های مقیاس تفاوت آشکاری در ضریب پایایی ندارند، به طوری که اگر بخواهیم جمع‌بندی از مقایسه نتایج پژوهش‌های مختلف داشته باشیم، این مقایسه را می‌توان چنین خلاصه کرد که این پژوهش‌ها به این نتیجه رسیده‌اند که هیچ برتری شاخصی در انتخاب روش‌های ساختن نمره مقیاس با معیار ضریب پایایی وجود نداشته و ضریب پایایی به شکلی که در این پژوهش‌ها محاسبه شده است نمی‌تواند ملاک مناسبی برای بررسی روش‌های ساختن نمره مقیاس باشد. لذا در این پژوهش به منظور مقایسه روش‌های متفاوت ساختن نمره مقیاس به دنبال رویکرد خطای استاندارد اندازه‌گیری شرطی رفته‌ایم.

این پژوهش با هدف بررسی ویژگی‌های اندازه‌گیری در روش‌های تبدیل مقیاس نرمال و تبدیل آرک‌سینوس برای تبدیل نمره‌های خام به نمره‌های مقیاس و همچنین مقایسه خطای استاندارد اندازه‌گیری شرطی نمره‌ها بین این دو روش انجام گرفته است و سؤال اصلی آن

این است که این دو روش تبدیل در مقدار خطای استاندارد اندازه‌گیری شرطی چه تفاوتی با هم دارند.

روش پژوهش

پژوهش حاضر پژوهشی کمی و از نوع توصیفی (غیرآزمایشی) است و به منظور توسعه دانش کاربردی در زمینه ساختن نمره مقیاس انجام گرفته که با این رویکرد، پژوهشی کاربردی محسوب می‌شود. جامعه مورد نظر در این پژوهش داوطلبان شرکت کننده در آزمون سراسری سال ۱۳۹۵ در گروه آزمایشی ریاضی و فنی هستند، در سال ۱۳۹۵ تعداد ۱۶۲۸۷۹ نفر در آزمون سراسری در رشته ریاضی و فنی شرکت کرده‌اند. در این پژوهش روش‌های مختلف مقیاس‌سازی و تحلیل‌ها بر روی ۱۰۰۰۰ داده شبیه‌سازی شده‌ی خرده‌آزمون‌های مختلف عمومی و اختصاصی آزمون سراسری ایران در گروه آزمایشی رشته ریاضی و فنی اجرا شده است. برای تولید داده‌های شبیه‌سازی شده از ویژگی هر کدام از خرده‌آزمون‌ها از جمله تعداد گزینه‌ها، تعداد سؤال‌ها و همچنین دشواری آن‌ها استفاده شده است. هدف از تولید داده‌های شبیه‌سازی در این پژوهش (۱) بررسی اولیه ویژگی‌های آزمون مانند خطای استاندارد اندازه‌گیری شرطی نمره‌های مقیاس، ضریب پایایی (۲) مبنایی برای بررسی صحت محاسبات و تحلیل داده‌های واقعی است. از آنجا که داده‌های واقعی جهت تعیین پارامترهای شبیه‌سازی استفاده شد (به منظور تولید داده‌هایی، هرچه بیشتر؛ مبتنی بر شرایط واقعی آزمون)، می‌توان گفت که مطالعه شبیه‌سازی بنیانی برای بررسی نتایج حاصل از تحلیل داده‌های واقعی فراهم می‌نماید که نشأت گرفته از شرایط واقعی آزمون است. مبنای شبیه‌سازی داده‌ها میانگین دشواری سؤال‌های آزمون است، به گونه‌ای که میانگین ضریب دشواری را به عنوان میانگین نسبی $\bar{P} = \frac{M_x}{K}$ (میانگین تقسیم بر تعداد سؤال‌ها) در نظر گرفته سپس بر اساس تعداد گزینه‌ها و تعداد سؤال‌های آزمون داده‌ها تولید شد (بروکس و یوهانسن، ۲۰۰۳). به عنوان مثال هرگاه بخواهیم ۳۰ سؤال دو ارزشی و چهار گزینه‌ای با ضریب دشواری ۰/۵۰

در سطح آزمون ایجاد کنیم، داده‌هایی تولید خواهد شد که دارای میانگین و واریانس تقریبی ۱۵ (حاصل ضرب ۰/۵ در ۳۰) و ۷/۵ (حاصل ضرب ۰/۲۵ در ۳۰) خواهند بود. به‌منظور تأیید نتایج شبیه‌سازی و همچنین مبنایی برای تولید داده‌های شبیه‌سازی از طریق نمونه‌گیری تصادفی، نمونه‌ای از داوطلبان آزمون سراسری سال ۹۵ برای بررسی پاسخ‌ها و خرده‌آزمون‌ها انتخاب شدند. با توجه به اینکه در این پژوهش با مقیاس بزرگ^۱ سروکار داریم از طریق قاعده سرانگشتی^۲ می‌توان حجم نمونه نزدیک به ۱۰۰۰۰ نفر را مناسب دانست. در آزمون سراسری و در گروه آزمایشی ریاضی و فنی چهار درس عمومی زبان و ادبیات فارسی (۲۵ سؤال)، زبان و ادبیات عربی (۲۵ سؤال)، معارف اسلامی (۲۵ سؤال) و زبان انگلیسی (۲۵ سؤال) و سه درس اختصاصی ریاضیات (۵۵ سؤال)، فیزیک (۴۵ سؤال) و شیمی (۳۵ سؤال)، مورد استفاده قرار می‌گیرد.

ابزار اصلی گردآوری داده در این پژوهش همان سؤال‌های آزمون سراسری است، داده‌های این آزمون در اختیار سازمان سنجش آموزش کشور بوده و به‌منظور تحلیل آزمون ۱۳۹۵ گروه آزمایشی ریاضی و فنی از آن‌ها بهره گرفته شده است. ضمناً به‌منظور انجام مطالعه شبیه‌سازی از داده‌های شبیه‌سازی شده نیز در کنار آن استفاده می‌شود. برای تولید داده‌های شبیه‌سازی شده از نرم‌افزار TAP^۳ نسخه ۴، ۷، ۱۴ استفاده شده است (بروکس و یوهانسن، ۲۰۱۴). برای برخی از روش‌ها و تبدیل‌هایی که در این پژوهش از آن‌ها استفاده می‌شود، نرم‌افزار تجاری به بازار عرضه نشده است به همین منظور در طی انجام پژوهش از نرم‌افزار کد نویسی ریاضی و آمار MATLAB^۴ برای نوشتن کدهای مربوط به برخی روش‌ها و تبدیل‌ها استفاده شد.

الف) تبدیل نمره‌های خام به نمره‌های مقیاس به روش نرمال‌سازی: برای به دست آوردن نمره مقیاس در این پژوهش به روش نرمال از رویکرد آنگگاف (۱۹۷۱) و کولن (۲۰۱۴)

1. large-scale assessment
2. thumbnail rule
3. test analysis program
4. MATLAB (matrix laboratory) is a multi-paradigm numerical computing environment and fourth-generation programming language.

استفاده شده است. در این روش پس از تبدیل نمره‌ها به نمره Z (رابطه ۳) از تبدیل خطی رابطه (۵) برای تبدیل خطی و محدود کردن دامنه نمره‌ها استفاده شد.

$$SC = \sigma(SC)Z + \mu(SC) \quad (5)$$

نمره‌های مقیاس^۱ SC هستند. در این پژوهش میانگین تبدیل خطی ۵۰۰۰ و انحراف استاندارد آن ۱۲۵۰ است. هدف از انتخاب این اعداد محدود کردن نمره‌های مقیاس بین اعداد ۰ تا ۱۰۰۰۰ است. این مقادیر اختیاری است و تنها نمره‌ها را در بازه مشخصی نگه می‌دارد، مقادیر میانگین و انحراف استاندارد در تبدیل خطی، نمره‌ها را در طول محور اعداد جابجا و یا در بازه مشخصی محدود می‌کنند. به‌طور مثال: $0 < SC < 10000 \Rightarrow 0 < 1250 \times Z + 5000 < 10000 \Rightarrow -4 < Z < +4$. پس از تبدیل خطی نمره‌ها به نزدیک‌ترین عدد صحیح گرد می‌شوند.

ب) تبدیل نمره‌های خام به نمره‌های مقیاس به روش آرک‌سینوس: برای تبدیل مقیاس به روش آرک‌سینوس از رابطه پیشنهادی کولن (۱۹۸۹) و کولن (۲۰۱۴) استفاده شد. در این پژوهش به‌منظور محدود شدن دامنه نمره‌های مقیاس آرک‌سینوس بین ۰ تا ۱۰۰۰۰ از تبدیل خطی $sc(y) = \alpha \times S(X_i) + \beta$ استفاده شد. در این تبدیل خطی $\alpha = 6690$ و $\beta = -67$ هستند. $S(X_i)$ تبدیل آرک‌سینوس رابطه (۴) است. این تبدیل دامنه نمره‌های مقیاس را در فاصله ۰ تا ۱۰۰۰۰ محدود می‌کند. در اینجا نیز نمره‌های مقیاس ساخته شده به نزدیک‌ترین عدد صحیح گرد خواهند شد.

روش محاسبه خطای استاندارد اندازه‌گیری شرطی. هم‌اکنون گزارش‌های فنی بیشتر خطای استاندارد اندازه‌گیری را به‌صورت خطای استاندارد اندازه‌گیری کلی^۲ ارائه می‌دهند ولی موسسه‌های AERA^۳ و NCEME^۴ و APA^۵ از سال ۱۹۸۵ برای استانداردهایی که به‌منظور تولید آزمون‌های آموزشی و روانی توصیه کرده‌اند، گزارش خطای استاندارد اندازه‌گیری

1. scale scores
2. overall standard error of measurement
3. american educational research association
4. national council on measurement in education
5. american psychological association

شرطی (استاندارد شماره ۲-۱۰ APA) را به همراه گزارش‌های فنی آزمون توصیه کرده‌اند (AERA، ۲۰۱۴). خطای استاندارد اندازه‌گیری شرطی قادر است میزان خطای استاندارد اندازه‌گیری را برای همه سطوح نمره‌ها برآورد کند. این شاخص آماری هم برای نمره‌های خام و هم برای نمره‌های مقیاس و هم برای نمره‌های مرکب قابل محاسبه است (وودروف و دیگران، ۲۰۱۳).

بررسی این شاخص آماری نشان می‌دهد که میزان خطای استاندارد اندازه‌گیری برای همه نمره‌ها برابر نیست و سطوح مختلف نمره‌ها دارای خطای استاندارد اندازه‌گیری شرطی متفاوت هستند. خطای استاندارد اندازه‌گیری شرطی تعریف مشابهی با خطای استاندارد اندازه‌گیری دارد، خطای استاندارد اندازه‌گیری شرطی، واریانس نمره‌ی مشاهده‌شده هر شرکت‌کننده در طی برگزاری آزمون‌های موازی در شرایط مشابه است، البته با این فرض که نمره حقیقی او ثابت بماند (هارتل، ۲۰۰۶).

رویکردهای متفاوتی توسط کولن (۱۹۹۲)، فلت و کوالس^۱ (۱۹۹۶) و برنان و لی^۲ (۱۹۹۹) برای محاسبه CSEM پیشنهاد شده است که به دلیل سهولت انجام محاسبه و عدم نیاز روش برنان و لی (۱۹۹۹) به محاسبه خطای استاندارد اندازه‌گیری نمره‌های خام از آن برای محاسبه خطای استاندارد اندازه‌گیری شرطی نمره‌های مقیاس در این پژوهش استفاده شده است. روش برنان و لی برای محاسبه خطای استاندارد اندازه‌گیری شرطی که به روش دوجمله‌ای نیز مشهور است از تعمیم دو روش یعنی نظریه نمره حقیقی قوی لرد (۱۹۶۵) و کولن (۱۹۹۲) استفاده می‌کند و رابطه‌ای را برای خطای استاندارد اندازه‌گیری شرطی ارائه می‌دهد (برنان و لی، ۱۹۹۹). بر اساس نظریه نمره حقیقی قوی^۳، احتمال شرطی اینکه شخصی از مجموع k سؤال در یک آزمون بتواند به γ تا از آن‌ها پاسخ صحیح بدهد از رابطه (۶) محاسبه می‌شود:

1. Feldt & Qualls
2. Brennan & Lee
3. strong true score theory

$$p(y|\pi, k) = \binom{k}{y} \pi^y (1-\pi)^{k-y} \quad (۶)$$

در رابطه (۶) پارامتر π نمره حقیقی نسبت پاسخ‌های صحیح برای هر شخص است و y متغیر تابع احتمال است. لرد (۱۹۶۵) مقدار خطای اندازه‌گیری را با توجه به توزیع دوجمله‌ای رابطه (۶) برای شخصی که می‌تواند به x سؤال پاسخ صحیح بدهد (x نمره‌ی خام است) طبق رابطه (۷) محاسبه کرده است:

$$\hat{\sigma}_{E(x)} = \sqrt{\frac{x(k-x)}{k-1}} = c_k \sqrt{\frac{x(k-x)}{k}} \quad (۷)$$

$c_k = \sqrt{\frac{k}{k-1}}$ همان عامل تصحیح سوگیری برآورد است چون در رابطه (۶) به منظور برآورد π از مقدار $\bar{x} = \frac{x}{k}$ استفاده می‌شود. عبارت $c_k = \sqrt{\frac{k}{k-1}}$ باعث می‌شود که $\hat{\sigma}_{E(x)}$ برآوردگر نااریبی از $\sigma_{E(x)}$ باشد. از همین ایده می‌توان برای محاسبه خطای استاندارد اندازه‌گیری شرطی نمره y بهره گرفت، در رابطه (۱۲) عبارت زیر رادیکال واریانس توزیع دوجمله‌ای است، پس می‌توان گفت که عامل تصحیح سوگیری برآورد یعنی c_k در واریانس شرطی توزیع دوجمله‌ای، یعنی در عبارت:

$$\sqrt{\frac{x(k-x)}{k}} = \sqrt{\frac{x}{k} k \left(1 - \frac{x}{k}\right)} = \sqrt{k\bar{x}(1-\bar{x})}$$

ضرب شده است. طبق آنچه در آمار مقدماتی موجود است می‌توان رابطه واریانس شرطی نمره‌های y را به صورت زیر نیز نوشت (مود و همکاران، ۲۰۰۸):

$$\sigma^2(Y|X) = E((Y - E(Y|X))^2|X) \quad (۸)$$

اگر توزیع شرطی نمره‌ها به صورت $p(Y|X)$ تعریف شده باشد رابطه (۸) به صورت رابطه (۹) در خواهد آمد:

$$\sigma^2(Y|X) = \sum_{y=-}^k ((Y - \mu(Y|X))^2 p(Y|X)) \quad (9)$$

رابطه (۹) ساده‌تر هم می‌شود به گونه‌ای که به میانگین شرطی نیازی نباشد:

$$\sigma^2(Y|X) = \sum_{y=-}^k Y^2 p(Y|X) - \left(\sum_{y=-}^k Y p(Y|X) \right)^2 \quad (10)$$

خطای استاندارد اندازه‌گیری شرطی نیز واریانس شرطی خطاها به شرط هر نمره هستند

پس به کمک رابطه (۱۱) خطای استاندارد اندازه‌گیری شرطی قابل محاسبه است:

$$\sigma_{E^2}(y|x) = \frac{k}{k-1} \left(\sum_{y=-}^k y^2 p(y|\pi, k) - \left(\sum_{y=-}^k y p(y|\pi, k) \right)^2 \right) \quad (11)$$

برای محاسبه خطای استاندارد اندازه‌گیری شرطی نمره‌های مقیاس طبق نظریه نمره

حقیقی قوی می‌توان در رابطه (۱۱) به جای مقادیر نمره خام y از تبدیل‌شده‌ی غیرخطی آن‌ها

(مثلاً نمره‌های نرمال) استفاده کرد:

$$\sigma_{E^2}(s(x)|x) = \frac{k}{k-1} \left(\sum_{y=-}^k f(y)^2 p(y|\pi, k) - \left(\sum_{y=-}^k f(y) p(y|\pi, k) \right)^2 \right) \quad (12)$$

برای مقایسه روش‌های مختلف مقیاس‌سازی از میانگین خطای استاندارد اندازه‌گیری

شرطی که به شکل رابطه (۱۳) استفاده می‌شود.

$$\bar{\sigma}_{E(s(x)|x)}^2 = f(x) \hat{\sigma}_{E(s(x)|x)}^2 \quad (13)$$

در رابطه (۱۳)، $f(x)$ فراوانی نسبی نمره x است $s(x)$ همان نمره‌های مقیاس و

$\hat{\sigma}_{E(s(x)|x)}$ خطای استاندارد اندازه‌گیری شرطی برای هر نمره مقیاس است (کولن، ۲۰۰۴).

یافته‌ها

در جدول ۱ نتایج تحلیل داده‌های شبیه‌سازی شده را برای نمونه تصادفی به حجم ۱۰۰۰۰ نفر مشاهده می‌کنید. این داده‌ها که به وسیله نرم‌افزار TAP طراحی شده‌اند به صورت امتیاز یک برای پاسخ صحیح و امتیاز صفر برای پاسخ‌های غلط و بی‌پاسخ تهیه شده‌اند و سپس با جمع امتیازها نمره‌های خام به دست آمده است. این ویژگی‌ها شامل شاخص‌های آماری مانند میانگین حسابی، واریانس، چولگی و کشیدگی و شاخص‌های اندازه‌گیری مانند خطای استاندارد اندازه‌گیری و ضریب پایایی کوادر ریچاردسون ۱۲۰ است. مقادیر داخل پرانتز میانگین نسبی هستند. خطای استاندارد اندازه‌گیری در جدول ۱ به وسیله ضریب پایایی کوادر ریچاردسون محاسبه شده و منظور از میانگین نسبی نمره‌های خام در این جدول حاصل تقسیم میانگین نمره‌های خام به تعداد سؤال‌های آن آزمون است.

جدول ۱. شاخص‌های آماری و شاخص‌های اندازه‌گیری برای نمره‌ی خام خرده آزمون‌ها (شبیه‌سازی)

شاخص	خرده آزمون‌ها						
	فارسی	عربی	معارف	زبان	ریاضی	فیزیک	شیمی
میانگین	۱۲/۶۵ (۰/۵۰)	۹/۹۴ (۰/۳۹)	۱۱/۱۳ (۰/۴۴)	۹/۰۴ (۰/۳۶)	۲۲/۲۳ (۰/۴)	۱۶/۵۵ (۰/۳۶)	۱۳/۵۲ (۰/۳۸)
واریانس	۲۰/۹۵	۲۲/۳۳	۲۰/۵۳	۱۹/۰۳	۹۲/۹۸	۵۶/۳۲	۳۴/۱۵
چولگی	-۰/۰۳۱	۰/۲۹	۰/۱۴	۰/۳۲	۰/۲۹	۰/۳۸	۰/۲۹
کشیدگی	-۰/۵۳	-۰/۴۵	-۰/۴۹	-۰/۳۹	-۰/۴۵	-۰/۳۵	-۰/۴۲
SEM	۲/۱۴۹	۲/۱۸۹	۲/۱۷۳	۲/۱۶۰	۳/۲۵	۲/۹۲	۲/۵۵
KR20	۰/۷۲۸	۰/۷۲۴	۰/۷۷۱	۰/۷۱۴	۰/۸۶۴	۰/۸۴۰	۰/۸۷۷

جدول ۲ شاخص‌های آماری و برخی از ویژگی‌های اندازه‌گیری را برای ۱۰۰۰۰ داده واقعی که به صورت تصادفی از میان داوطلبان رشته ریاضی شرکت کننده در آزمون سراسری در سال ۱۳۹۵ انتخاب شده‌اند را نمایش می‌دهد.

1. KR20

جدول ۲. شاخص‌های آماری و شاخص‌های اندازه‌گیری
برای نمره‌ی خام خرده‌آزمون‌ها (داده‌های واقعی)

شاخص	خرده‌آزمون‌ها						
	فارسی	عربی	معارف	زبان	ریاضی	فیزیک	شیمی
میانگین	۷/۶۸	۵/۲۹	۹/۰۵	۵/۹۹	۴/۷۲	۵/۱۹	۳/۰۳
	(۰/۳۰)	(۰/۲۱)	(۰/۳۶)	(۰/۲۳)	(۰/۰۸)	(۰/۱۱)	(۰/۰۸)
واریانس	۱۸/۴۳	۲۰/۳۷	۳۱/۸۴	۴۰/۰۰	۴۱/۷۸	۴۶/۵۰	۱۵/۸۵
چولگی	۰/۳۷	۱/۲۰	۰/۴۱	۱/۰۳	۲/۴۳	۱/۹۷	۲/۰۱
کشیدگی	-۰/۱۲	۱/۶۳	-۰/۶۹	۰/۱۳	۷/۶۳	۴/۲۹	۵/۲۶
SEM	۲/۰۵۳	۱/۸۴۳	۲/۲۴۱	۲/۰۵۴	۲/۰۲۴	۲/۰۹۵	۱/۶۱۲
KR20	۰/۷۷۱	۰/۸۳۳	۰/۸۴۲	۰/۸۹۴	۰/۹۰۱	۰/۹۰۵	۰/۸۳۶

کمترین نمره خام برای دروس عمومی صفر و بیشترین نمره خام برای آن‌ها ۲۵ است. کمترین نمره خام برای خرده‌آزمون ریاضیات صفر و بیشترین آن ۵۵ است. همچنین کمترین نمره خام برای خرده‌آزمون‌های فیزیک و شیمی صفر و بیشترین نمره خام برای آن‌ها به ترتیب ۴۵ و ۳۵ است. از ضریب کودر ریچاردسون برای محاسبه پایایی نمره‌های خام استفاده گردید که طبق داده‌های جدول‌های ۱ و ۲ ضریب پایایی نمره‌های خام بین ۰/۷۱۴ تا ۰/۹۰۱ متغیر بوده است، تغییر مقادیر پایایی بین خرده‌آزمون‌ها ارتباط نزدیکی با تعداد سؤال‌های خرده‌آزمون داشته است، به طوری که چه در داده‌های شبیه‌سازی شده و چه در داده‌های واقعی آن دسته از خرده‌آزمون‌هایی که تعداد سؤال‌های بیشتری داشته‌اند، داری ضریب پایایی بزرگ‌تر بودند و خرده‌آزمون‌هایی که داری تعداد سؤال‌های کمتر بودند حائز ضریب پایایی کمتری شده‌اند. تأثیر دیگری که تعداد سؤال‌ها بر ویژگی نمره‌های خام داشته تغییر در واریانس نمره‌ها بود، به طوری که خرده‌آزمون‌هایی که تعداد سؤال بیشتری داشتند نسبت به خرده‌آزمون‌هایی که تعداد سؤال کمتری داشته‌اند، واریانس بزرگ‌تری را در بین نمره‌ها ایجاد کردند. البته این گفته تا حدودی بستگی به دشواری سؤال‌های هر خرده‌آزمون نیز دارد، به طوری که در داده‌های واقعی جدول ۲ خرده‌آزمون‌های دشوار مانند ریاضیات، فیزیک و شیمی علی‌رغم تعداد سؤال‌های بیشتر نتوانستند واریانس را خیلی نسبت به

خرده‌آزمون‌های آسان‌تر مثل ادبیات فارسی افزایش دهند. با دقت در جدول ۲ می‌بینید که خرده‌آزمون شیمی به دلیل دشواری زیاد با وجود تعداد سؤال بیشتر نسبت به ادبیات فارسی واریانس کمتری را در نمره‌های خام ایجاد کرده است.

از دو روش نرمال‌سازی و آرک‌سینوس برای تبدیل نمره‌های خام به نمره‌های مقیاس در این پژوهش بهره گرفته شده است. جدول ۳ و جدول ۴ برخی شاخص‌های آماری به همراه ضریب پایایی را برای نمره‌های شبیه‌سازی شده و واقعی ۱۰۰۰۰ نمونه تصادفی برای هفت خرده‌آزمون نمایش می‌دهند. برای محاسبه ضریب پایایی گزارش شده از رابطه پایایی کولن (۲۰۱۲) استفاده شد. این ضریب با فرض اینکه توزیع نمره‌های خام دو جمله‌ای است محاسبه شده است و برای محاسبه ضریب پایایی نمره‌های مقیاس به کار می‌رود (کولن، ۲۰۱۲). ضرایب پایایی مانند کرونباخ، کودر ریچاردسون و... تنها برای نمره‌های خام کاربرد داشته و بدون اطلاع از چگونگی پاسخ افراد به همه سؤال‌ها قادر به محاسبه پایایی نیستند، این در حالی است که ضریب پایایی کولن بدون دسترسی به پاسخ همه سؤال‌ها توسط افراد و تنها با داشتن توزیع نمره‌ها می‌تواند پایایی نمره‌های مقیاس را محاسبه کند.

جدول ۳: برخی شاخص‌های آماری به همراه ضریب پایایی نمره‌های نرمال (داده‌های شبیه‌سازی شده)

خرده‌آزمون	گشتاورهای اول تا چهارم			
	میانگین	واریانس	چولگی	کشیدگی
فارسی	۴۹۹۹/۸	۱۵۴۵۰۰۰	-/۰۳	-۰/۶۶
عربی	۵۰۰۱/۳	۱۵۴۱۶۰۰	۰/۸۵	-۰/۵۳
معارف	۵۰۰۰/۶	۱۵۴۴۰۰۰	۰/۴۵	-۰/۶۲
زبان خارجه	۵۰۰۱/۶	۱۵۳۸۸۰۰	۱/۰۶	-۰/۶۲
ریاضی	۵۰۰۰/۲	۱۵۵۷۴۰۰	۰/۶۱	-۰/۶۱
فیزیک	۵۰۰۰/۷	۱۵۵۳۱۰۰	۰/۷۷	-۰/۶۵
شیمی	۵۰۰۰/۹	۱۵۵۰۰۰۰	۰/۸۱	-۰/۶۶

جدول ۴. برخی شاخص‌های آماری به همراه ضریب پایایی نمره‌های نرمال (داده‌های واقعی)

خرده‌آزمون	گشتاورهای اول تا چهارم				ضریب پایایی
	میانگین	واریانس	چولگی	کشیدگی	
فارسی	۵۰۰۵/۳	۱۵۱۴۱۰۰	۱/۳۷	-۰/۵۳	۰/۹۵۲
عربی	۵۰۲۰/۳	۱۴۳۳۵۰۰	۱/۴۲	-۰/۷۱	۰/۹۴۰
معارف	۵۰۰۴/۶	۱۵۲۰۷۰۰	۱/۲۵	-۰/۰۷	۰/۹۵۷
زبان خارجه	۵۰۳۸/۶	۱۳۴۰۰۰۰	۱/۵۴	-۰/۲۰	۰/۹۴۰
ریاضی	۵۰۳۹/۷	۱۳۳۷۹۰۰	۱/۲۷	-۱/۳۶	۰/۹۲۲
فیزیک	۵۰۴۰/۸	۱۳۳۳۵۰۰	۱/۳۴	-۱/۱۳	۰/۹۲۴
شیمی	۵۰۵۳/۳	۱۲۶۳۴۰۰	۱/۲۹	-۱/۳۸	۰/۹۱۲

جدول ۵ برخی شاخص‌های آماری به همراه ضریب پایایی را برای تبدیل آرک‌سینوس نمایش می‌دهد. داده‌های این جدول حاصل از ۱۰۰۰۰ داده شبیه‌سازی شده است که به صورت تصادفی برای هفت خرده‌آزمون انتخاب شده‌اند.

جدول ۵. برخی شاخص‌های آماری و ضریب پایایی نمره‌های آرک‌سینوس (شبیه‌سازی)

خرده‌آزمون	گشتاورهای اول تا چهارم				ضریب پایایی
	میانگین	واریانس	چولگی	کشیدگی	
فارسی	۵۰۴۱/۸	۱۷۲۰۷۰۰	۰/۰۵	۰/۶۴	۰/۹۵۹
عربی	۴۲۵۴/۹	۱۹۳۴۳۰۰	۰/۸۸	-۰/۵۱	۰/۹۶۲
معارف	۴۶۰۶/۵	۱۷۰۱۱۰۰	۰/۵۰	-۰/۶۰	۰/۹۵۸
زبان خارجه	۳۹۹۲/۲	۱۶۸۳۱۰۰	۱/۱۰	-۰/۴۶	۰/۹۵۶
ریاضی	۴۲۹۷/۹	۱۶۸۰۱۰۰	۰/۸۲	-۰/۶۲	۰/۹۷۲
فیزیک	۴۰۳۴	۱۵۳۲۹۰۰	۱/۰۵	-۰/۵۴	۰/۹۶۵
شیمی	۴۱۷۷/۲	۱۴۹۶۹۰۰	۰/۹۴	-۰/۵۴	۰/۹۵۹

جدول ۶ برخی شاخص‌های آماری به همراه ضریب پایایی را برای نمره‌های واقعی ۱۰۰۰۰ نمونه تصادفی داوطلبان آزمون سراسری در سال ۱۳۹۵ را نمایش می‌دهد که در آن‌ها از تبدیل آرک‌سینوس برای تولید مقیاس استفاده شده است.

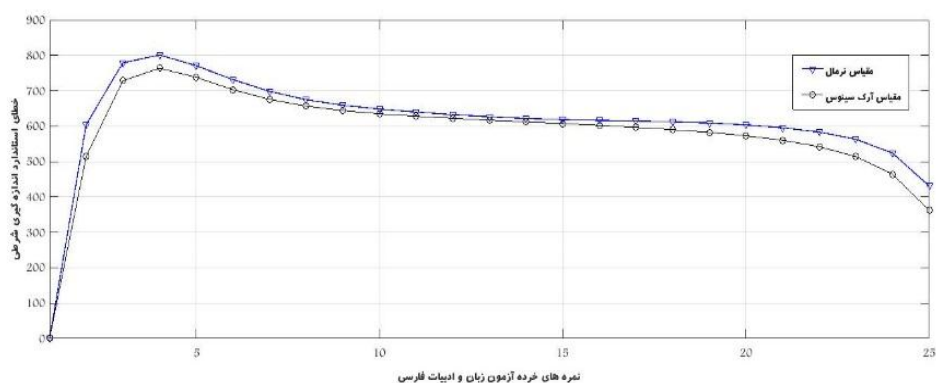
جدول ۶. برخی شاخص های آماری به همراه ضریب پایایی نمره های آرک سینوس (داده های واقعی)

خرده آزمون	گشتاورهای اول تا چهارم			
	میانگین	واریانس	چولگی	کشیدگی
فارسی	۳۵۵۰/۵۰	۱۹۱۶۴۰۰	۱/۳۵	-۰/۳۹
عربی	۲۷۰۰/۱	۲۵۴۹۲۰۰	۱/۵۰	-۰/۵۳
معارف	۳۹۳۹	۳۰۴۹۷۰۰	۱/۱۴	-۰/۴۳
زبان خارجه	۲۷۷۰/۳	۴۸۶۸۹۰۰	۱/۵۰	-۰/۵۱
ریاضی	۱۳۷۲/۸	۱۶۱۷۰۰۰	۱/۳۸	-۱/۰۳
فیزیک	۱۶۴۲/۹	۲۳۰۶۰۰۰	۱/۴۱	-۰/۹۴
شیمی	۱۳۹۵/۳	۱۶۰۱۱۰۰	۱/۴۰	-۱

تفاوت بین ویژگی های آماری و اندازه گیری در تبدیل نرمال و آرک سینوس را می توان در مقایسه جداول ۳ و ۴ با جداول ۵ و ۶ مشاهده کرد، ضریب پایایی نمره ها در هر دو روش مقدار بالایی را نشان داده است، دقت بالای نمره ها در تبدیل نمره های خام به نمره های مقیاس در همه خرده آزمون ها مشاهده می شود، هرچند که مقدار ضریب پایایی برای تبدیل آرک سینوس در نمره های واقعی و شبیه سازی شده مقداری بیشتر از ضریب پایایی در نمره های نرمال را نشان داده است ولی این تفاوت در بیشترین موارد به چند صدم نمی رسد. اگر چهار گشتاور اول (میانگین، واریانس، چولگی و کشیدگی) هر دودسته نمره را با هم مقایسه کنیم می بینیم که روش تبدیل آرک سینوس توانسته است، واریانس بیشتری را بین نمره ها ایجاد کند. در مقایسه واریانس خرده آزمون ها مشاهده می کنید که واریانس در روش تبدیل آرک سینوس از روش نرمال سازی در همه آن ها بیشتر بوده که نشان می دهد این روش قادر است تفکیک بهتری بین سطوح نمره ها ایجاد کند. میزان چولگی و کشیدگی نمره ها در هر دو روش تفاوت زیادی نداشته و نشان می دهد که دو روش بر روی چولگی و کشیدگی نمره ها تغییر چندانی ایجاد نمی کنند.

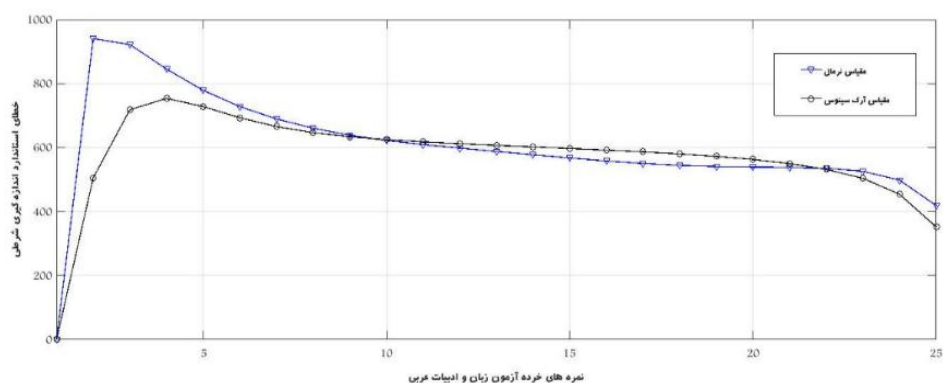
در این پژوهش از دو روش مقیاس سازی برای تبدیل نمره های خام به نمره های مقیاس استفاده شده است. روش اول که به روش نرمال سازی مشهور است و روش دوم که به تبدیل آرک سینوس مشهور است در این قسمت قصد بر آن است که با مقایسه نمودارهای خطای

استاندارد اندازه‌گیری شرطی و همچنین مقادیر میانگین این خطا برای هر روش، این دو نوع تبدیل با هم مقایسه شوند. نمودار شکل ۱ خطای استاندارد اندازه‌گیری شرطی برای نمره‌های خرده‌آزمون زبان و ادبیات فارسی را برای نمره‌های مقیاس نرمال و آرک‌سینوس با هم مقایسه کرده است.



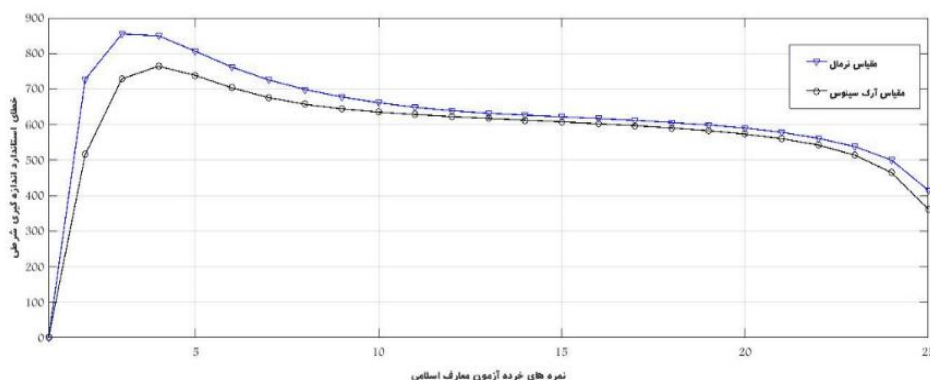
شکل ۱. مقایسه خطای استاندارد اندازه‌گیری شرطی نمره‌های مقیاس خرده‌آزمون فارسی

همان‌گونه که در شکل ۱ می‌بینید، خطای استاندارد اندازه‌گیری شرطی برای مقیاس آرک‌سینوس از مقیاس نرمال در بیشتر سطوح کمتر بوده هرچند در میانه طیف نمره‌ها این خطاها به هم نزدیک هستند. شکل ۲ خطای استاندارد اندازه‌گیری شرطی برای نمره‌های خرده‌آزمون زبان و ادبیات عربی را برای نمره‌های مقیاس نرمال و آرک‌سینوس با هم مقایسه کرده است.



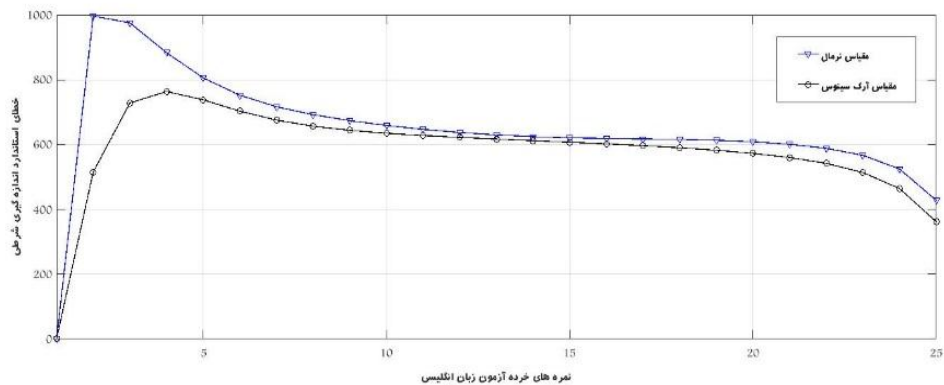
شکل ۲. مقایسه خطای استاندارد اندازه‌گیری شرطی نمره‌های مقیاس خرده‌آزمون عربی

در نمودار شکل ۲ می‌بینید، خطای استاندارد اندازه‌گیری شرطی برای مقیاس آرک‌سینوس از مقیاس نرمال در بیشتر سطوح کمتر بوده هرچند در میانه طیف نمره‌ها این خطاها به هم نزدیک هستند و در نقاطی هم این خطای مقیاس نرمال است که کمتر شده ولی نوسان خطاها در مقیاس آرک‌سینوس نسبت به مقیاس نرمال کمتر بوده به طوری که نمودار آن افقی‌تر از نمودار مقیاس نرمال به نظر می‌رسد که خود نشان‌دهنده‌ی هم‌ترازی خطاها برای بیشتر سطوح نمره‌ها است. شکل ۳ همین نمودار را برای خرده‌آزمون معارف اسلامی نشان داده است.



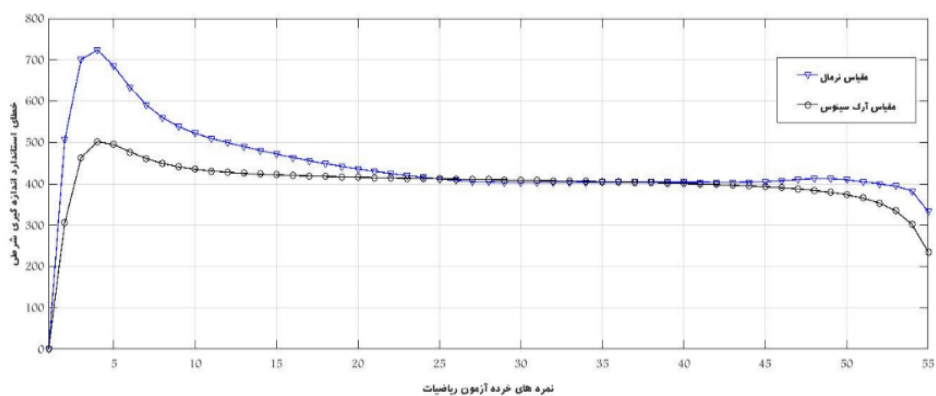
شکل ۳. مقایسه خطای استاندارد اندازه‌گیری شرطی نمره‌های مقیاس خرده‌آزمون معارف اسلامی

بازهم در شکل ۳ می‌توانید کاهش خطای استاندارد اندازه‌گیری شرطی را برای مقیاس آرک‌سینوس و مسطح شدن نمودار آن را نسبت به نمودار مقیاس نرمال مشاهده کنید. شکل ۴ نمودار را برای خرده‌آزمون زبان و ادبیات انگلیسی نشان داده است.



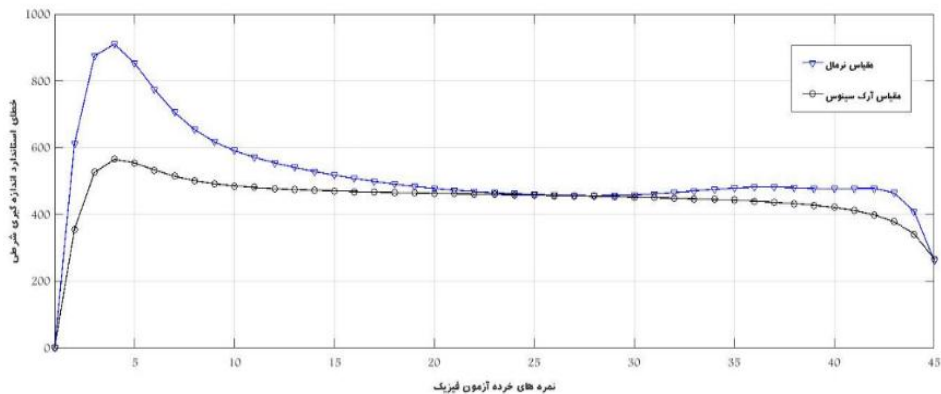
شکل ۴. مقایسه خطای استاندارد اندازه‌گیری شرطی
نمره‌های مقیاس خرده‌آزمون زبان و ادبیات انگلیسی

شکل ۵ نمودار را برای خرده‌آزمون ریاضی نمایش داده است.



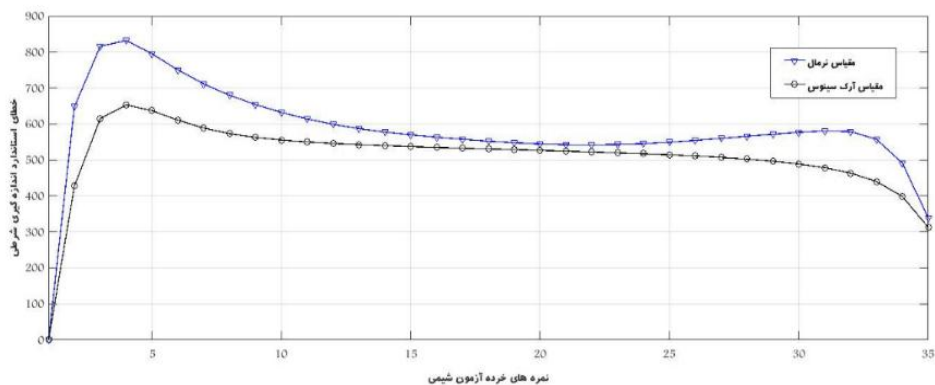
شکل ۵. مقایسه خطای استاندارد اندازه‌گیری شرطی نمره‌های مقیاس خرده‌آزمون ریاضی

در نمودار ۵ به کاهش چشم‌گیر خطا و همچنین مسطح شدن نمودار به خصوص در میانه طیف دقت کنید. نمودار ۶ مربوط به خرده‌آزمون فیزیک است.



شکل ۶. مقایسه خطای استاندارد اندازه‌گیری شرطی نمره‌های مقیاس خرده‌آزمون فیزیک

در نمودار ۶ نیز میانه طیف نمره‌ها، در هر دو مقیاس خطای مشابهی را به نمایش گذاشته‌اند. شکل ۷ نمودار را برای خرده‌آزمون شیمی نشان داده است.



شکل ۷. مقایسه خطای استاندارد اندازه‌گیری شرطی نمره‌های مقیاس خرده‌آزمون شیمی

نمودارهای ۱ تا ۷ مقایسه بین خطای استاندارد اندازه‌گیری شرطی بین نمره‌های خرده‌آزمون‌های مختلف را در دو نوع تبدیل نشان می‌دهند، دقت در این نمودارها دو نکته اساسی را نمایش می‌دهد، اول اینکه استفاده از روش آرک‌سینوس برای تبدیل نمره‌های خام به نمره‌های مقیاس باعث کاهش خطای استاندارد اندازه‌گیری شرطی برای کلیه سطوح نمره‌ها گردیده است به طوری که به عنوان مثال بیشینه خطای استاندارد اندازه‌گیری شده در درس ریاضی به روش نرمال‌سازی عدد ۵۹۶/۰۳۲ است، ولی بیشینه خطای استاندارد اندازه‌گیری

شرطی در درس ریاضی برای روش آرک‌سینوس $۵۶۱/۸۳$ است، دومین نکته اینکه روش آرک‌سینوس در تبدیل نمره‌های خام به نمره‌های مقیاس نوسان کمتری داشته و برای همه سطوح نمره‌ها خطا تقریباً یکسان است. نمودار در ابتدا اعداد بالایی را برای خطای استاندارد اندازه‌گیری شرطی در هر دو تبدیل نشان می‌دهند که وقتی به نقاط میانی توزیع نزدیک می‌شویم این مقدار کاهش می‌یابد، این کاهش خطا در روش آرک‌سینوس در میانه توزیع بیشتر بوده ضمن اینکه نمودار روش تبدیل آرک‌سینوس در نقاط میانی نوسان نداشته و به خط راست نزدیک‌تر است که این ویژگی نشان از تمایل این روش به همسان‌سازی خطاهای استاندارد اندازه‌گیری شرطی برای همه سطوح نمره‌ها دارد، یا به زبان ساده در این روش برای همه نمره‌ها تقریباً به یک اندازه مرتکب خطای استاندارد اندازه‌گیری می‌شویم. در روش نرمال‌سازی ضمن اینکه مقدار خطای استاندارد اندازه‌گیری شرطی بالایی را برای بیشتر سطوح نمره‌ها مرتکب می‌شویم، برای همه سطوح نمره‌ها به یک اندازه خطا نمی‌کنیم و این مقدار دچار نوسان است که در نمودارهای ۱ تا ۷ این نوسان را مشاهده می‌کنید.

جدول ۷ مقادیر مجذور میانگین خطای استاندارد اندازه‌گیری شرطی برای هفت خرده‌آزمون در دو مقیاس استفاده‌شده در این پژوهش را نشان می‌دهد.

جدول ۷. مجذور میانگین خطای استاندارد اندازه‌گیری شرطی نمره‌های مقیاس

مقیاس	نوع داده						
	فارسی	عربی	معارف	زبان	ریاضی	فیزیک	شیمی
خرده‌آزمون‌ها							
نرمال							
واقعی	۲۷۰/۴۷	۲۹۳/۰۷	۲۵۴/۲۰	۲۸۲/۹۸	۳۲۳/۶۰	۳۱۷/۳۰	۳۳۳/۲۳
شبه‌سازی	۲۵۸/۳۷	۲۵۹/۹۵	۲۶۲/۹۷	۲۷۰/۱۴	۲۱۷/۴۴	۲۳۸/۳۵	۲۵۲/۳۴
آرک‌سینوس							
واقعی	۲۵۳/۹۱	۲۴۷/۷۸	۲۵۲/۲	۲۳۵/۹۵	۱۹۳/۱۸	۲۰۳/۶۸	۲۱۴/۲۴
شبه‌سازی	۲۴۸/۵۷	۲۵۱/۲۳	۲۵۱/۲۱	۲۵۴/۶۴	۲۰۴/۱۳	۲۱۷/۵۶	۲۳۳/۹۱

با توجه به جدول ۷ مقادیر مجذور میانگین مربوط به روش آرک‌سینوس در همه خرده‌آزمون‌ها از روش نرمال‌سازی کمتر است، میانگین خطای استاندارد اندازه‌گیری شرطی

در روش نرمال‌سازی چه در داده‌های شبیه‌سازی شده و چه در داده‌های واقعی مقدار بیشتری را نشان می‌دهد که دلیل آن‌هم به خاطر مقدار بالای خطای استاندارد اندازه‌گیری شرطی در ابتدا و انتهای توزیع و هم به دلیل نوسان‌هایی است که این خطا در میانه‌های توزیع از خود نشان داده است.

بحث و نتیجه‌گیری

پژوهش حاضر به منظور بررسی ویژگی‌های اندازه‌گیری و آماری در روش‌های نرمال‌سازی و آرک‌سینوس برای تبدیل نمره‌های خام به نمره‌های مقیاس و همچنین مقایسه خطای استاندارد اندازه‌گیری شرطی نمره‌ها بین این دو روش طرح ریزی شده است. در این پژوهش نمره‌های خام حاصل از خرده‌آزمون‌های ادبیات فارسی، عربی، معارف اسلامی، زبان انگلیسی، ریاضی، فیزیک و شیمی که مربوط به ۱۰۰۰۰ نمونه شبیه‌سازی شده و ۱۰۰۰۰ نمونه واقعی بود را به دو روش نرمال‌سازی و تبدیل آرک‌سینوس به نمره‌های مقیاس تبدیل شدند. شاخص‌های مختلف آماری و شاخص‌های اندازه‌گیری مانند خطای استاندارد اندازه‌گیری و ضریب پایایی هم برای نمره‌های خام و هم برای نمره‌های مقیاس محاسبه و گزارش شدند. از نظریه نمره حقیقی قوی برای بررسی خطای استاندارد اندازه‌گیری شرطی نمره‌های مقیاس و همچنین محاسبه پایایی نمره‌ها استفاده شد و با استفاده از مجذور میانگین خطای استاندارد اندازه‌گیری شرطی و نمودارهای خطای استاندارد اندازه‌گیری شرطی اقدام به مقایسه دقت روش‌های ساخت نمره مقیاس (روش‌های نرمال‌سازی و تبدیل آرک‌سینوس) نمودیم. نتایج حاکی از برتری روش تبدیل آرک‌سینوس در کاهش خطای استاندارد اندازه‌گیری شرطی برای توزیع نمره‌ها داشت، به طوری که با وجود خطای استاندارد اندازه‌گیری شرطی بالایی که هر دو روش برای ابتدا و انتهای توزیع نمره‌ها از خود نشان دادند، روش تبدیل آرک‌سینوس توانسته ضمن کاهش میزان خطای استاندارد اندازه‌گیری شرطی برای میانه توزیع موجب همسان‌سازی محسوس در خطا گردیده به طوری که در نقاط میانی توزیع شاهد تبدیل این منحنی به یک خط تقریباً راست بودیم. میانگین خطای استاندارد اندازه‌گیری شرطی روش تبدیل آرک‌سینوس از روش نرمال‌سازی کمتر بوده است و به نظر می‌رسد

استفاده از تبدیل آرک‌سینوس برای تبدیل نمره‌های خام به نمره‌های مقیاس می‌تواند باعث کاهش خطای استاندارد اندازه‌گیری برای همه سطوح نمره‌ها شده و به همسانی خطا برای همه سطوح نمره‌ها کمک کند.

در این پژوهش با بررسی که در انواع روش‌های تبدیل مقیاس صورت گرفت دیده شده که هرچند نوسان خطا در مقیاس نرمال زیادتر از مقیاس آرک‌سینوس بوده است، ولی روش مقیاس نرمال نیز در برخی موارد مقدار خطای کمتری را نسبت آرک‌سینوس نشان داده است این خود نشان می‌دهد که نباید نتیجه‌گیری مطلقاً را گزارش داد و برترین مقیاس را به مقیاس خاصی اطلاق کرد، همچنین با تمرکز بر روی تحلیل نتایج مقیاس آرک‌سینوس مشخص شد که این مقیاس دارای دو برتری نسبت به مقیاس نرمال است، برتری اول کاهش مقدار خطا و برتری دوم کاهش نوسان خطاهاست، مقیاس آرک‌سینوس توانسته خطاها را برای سطوح مختلف نمره‌ها همانند کند، نتایج به‌دست آمده در این پژوهش با نتایج حاصل از پژوهش‌های کولن و همکاران (۱۹۸۹) همخوانی دارد. کولن (۲۰۱۴) بر همان نتایج در کتاب هم‌تراز سازی، پیوند زدن و مقیاس‌سازی اشاره مجدد کرده است، در پژوهش کولن و همکاران (۱۹۸۹) از آرک‌سینوس به‌عنوان مقیاسی نام برده شده که می‌تواند نوسان خطاها را کاهش داده و به همانندسازی خطاها برای همه سطوح نمره‌ها کمک کند، پژوهش حاضر نیز نتایج گزارش شده توسط آن پژوهش را تأیید کرده و نشان داد هم در داده‌های شبیه‌سازی و هم در داده‌های واقعی نوسان خطاها در مقیاس آرک‌سینوس کاهش یافته و می‌توان به‌طور تقریبی چنین بیان کرد که خطا تقریباً برای همه سطوح نمره‌ها یکسان است. همچنین در آن پژوهش به این موضوع اشاره شده که مقیاس خوب و برتری وجود ندارد و هرکدام از مقیاس‌ها داری ویژگی‌های خاص خودش است و مطلوب بودن هر مقیاس در معناداری و تفسیرپذیری بهتر آن است. همچنین نتایج این پژوهش نشان داد که هیچ‌کدام از روش‌های مقیاس‌سازی نتوانسته تفاوت آشکاری را در مقدار ضریب پایایی ایجاد کنند و برای هرکدام از روش‌ها مزایایی بیان شد، به‌عنوان‌مثال مقیاس نرمال در نمودار خطای استاندارد اندازه‌گیری دچار نوسان است و یا اینکه خطاها در مقیاس آرک‌سینوس به هم نزدیک هستند، چانگک در سال ۲۰۰۶ برای بررسی انواع روش‌های مقیاس‌سازی خطی و غیرخطی از

ضریب پایایی آن‌ها استفاده کرده است، در پژوهش چانگ (۲۰۰۶) از روش خطی تبدیل نمره‌ها، روش آرک‌سینوس و نرمال بهره برده شده، نتایج این پژوهش با پژوهش چانگ (۲۰۰۶) در برآورد مقدار ضریب پایایی به هم نزدیک بوده و هر دو پژوهش روی این موضوع تأکید دارند که روش‌های مقیاس‌سازی نتوانستند تأثیر چشم‌گیری بر مقدار پایایی داشته باشند. این گفته بیانگر این موضوع است که ضریب پایایی معیار مناسبی برای مقایسه نمره‌های مقیاس نیست. این پژوهش نشان داد که هرگاه معیار برای انتخاب تبدیل مقیاس را میانگین خطای استاندارد اندازه‌گیری شرطی قرار دهیم، می‌توانیم بین روش‌های متفاوت مقیاس‌سازی تفکیک و تفاوت قائل شویم. به گونه‌ای که جدول ۷ نمایش داد مجذور میانگین خطای استاندارد اندازه‌گیری شرطی در بین روش‌های متفاوت مقیاس‌سازی مقادیر متفاوتی خواهد داشت و نشان می‌دهد که می‌تواند معیار مناسبی برای مقایسه ایجاد کند.

به‌عنوان پیشنهاد کاربردی باید چنین گفت که استفاده از مقیاس آرک‌سینوس، باعث همانند شدن خطای استاندارد اندازه‌گیری برای همه سطوح نمره‌ها می‌شود، به همین منظور در صورتی که بخواهید نمره عادلانه‌تری به افراد اختصاص دهید از مقیاس آرک‌سینوس برای تبدیل نمره‌های خام به نمره‌های مقیاس استفاده کنید. ضمن اینکه برای تفسیرپذیری بهتر آزمون‌ها توصیه می‌شود که خطای استاندارد اندازه‌گیری شرطی هر نمره محاسبه شود و آزمون‌ساز در گزارش فنی مربوط به آزمون به این خطاها اشاره کند، این موضوع باعث می‌شود که دید بهتری نسبت به نمره‌ها داشته و در مواردی که قرار است از برش نمره‌ها استفاده شود با دقت مطلوب‌تری این کار صورت بگیرد.

در این پژوهش از روش‌های هموارسازی برای ساختن نمره‌های مقیاس استفاده نشد و همه نمره‌ها به صورت هموار نشده در تحلیل‌ها مورد استفاده قرار گرفته‌اند. برای هموار شدن توزیع نمره‌ها و کمک به کاهش خطاهای تصادفی می‌توان از یکی از روش‌های هموارسازی توزیع بتای چهار پارامتری^۱ (لرد ۱۹۶۵) و یا دو جمله‌ای کرنل (کولن ۱۹۹۱) نیز استفاده کرد. نکته آخر اینکه در این پژوهش از توزیع دو جمله‌ای نمره حقیقی قوی برای توزیع شرطی

1. Beta4

نمره‌های مشاهده شده و حقیقی بهره برده‌ایم که با بسط این توزیع می‌توان از توزیع دوجمله‌ای گسترش یافته (کولن ۱۹۹۲) نیز برای بررسی دقیق‌تر خطای استاندارد اندازه‌گیری شرطی استفاده کرد.

منابع

- سازمان سنجش آموزش کشور (۱۳۹۵). کارنامه آماری آزمون سراسری سال ۱۳۹۵. تهران: انتشارات سازمان سنجش آموزش کشور (دفتر طرح و آمار)
- نقی زاده، سیما (۱۳۹۴). نمره کل سازی آزمون سراسری در گروه آزمایشی علوم ریاضی و فنی سال ۱۳۹۱ بر اساس توزیع واقعی نمرات و مقایسه آن با روش فعلی. تهران: مرکز تحقیقات ارزشیابی، اعتبار سنجی و تضمین کیفیت آموزش عالی (سازمان سنجش آموزش کشور).
- Allen, M. J., & Wendy, Y. M. (1979). *Introduction to Measurement Theory*. California: Cole publishing company.
- Angoff, W.H. (1971). *Scales, norms, and equivalent scores*. In RL. Thorndike (Ed.).
- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education, (2014). *Standards for educational and psychological testing*. Washington, DC: APA
- Brennan, R. L., & Lee, W. C. (1999). Conditional scale-score standard errors of measurement under binomial and compound binomial assumptions. *Educational and Psychological Measurement*, 59(1), 5-24.
- Brooks, G. P., & Johnson, G. A. (2003). TAP: Test Analysis Program. *Applied Psychological Measurement*. 27(4), 303-304.
- Brooks, G. P., & Johnson, G. A. (2014). *TAP: Test Analysis Program* version (14.7.4) [computer software]. Retrieved from <http://www.ohio.edu/people/brooksg/software.htm>.
- Chang, S. W. (2006). Methods in Scaling the Basic Competence Test. *Educational and Psychological Measurement*, 66(6), 907-929.
- Dorans N. J., Pommerich, M. & Holland P. W. (2007). *A Framework and History for Score Linking*. In Holland P. W. (Eds.), *Linking and Aligning Scores and Scales* (pp 5-30). New York: Springer.
- Feldt, L. S., & Brennan, R. L. (1989). *Reliability*. In R. L. Linn (Ed.), *Educational measurement* (3rd ed., pp. 105-146). New York, NY: Macmillan.

- Feldt, L. S., & Quails, A. L. (1996). Estimation of measurement error variance at specific score levels. *Journal of Educational Measurement*, 33, 141-156. 156.
- Gulliksen, H. (1950). *Theory of mental test*. New York: John Wiley & sons.
- Haertel, H. E. (2006). *Reliability*. In R. L. Brennan (Ed.), *Educational measurement* (4th ed., pp. 65-86). CT: American Council on Education and Praeger.
- Iowa Assessment (2016). *Iowa Test of Basic Skills*. Iowa City: Author Retrieved: itp.education.uiowa.edu
- Kolen, M. J., Hanson, B. A., & Brennan, R. L. (1992). Conditional standard errors of measurement of scale scores. *Journal of Educational Measurement*, 29, 285-307.
- Kolen, M. J., & Hanson, B. A. (1989). *Scaling the ACT Assessment*. In R. L. Brennan (Ed.), *Methodology used in scaling the ACT Assessment and P-ACT+* (pp. 35-55). Iowa City, IA: American College Testing Program.
- Kolen, M. J. (1991). Smoothing methods for estimating test score distributions. *Journal of Educational Measurement*, 28, 257-282.
- Kolen, M. J., & Brennan, R. L. (2004). *Test Equating, Scaling and Linking* (2nd Ed.). New York: Springer.
- Kolen, M. J., Wang, T., Lee, W. Chon. (2012). Conditional Standard Errors of Measurement for Composite Scores Using IRT. *International Journal of Testing*, 12, 1-20.
- Kolen, M. J., & Brennan, R. L. (2014). *Test Equating, Scaling and Linking*, 3rd Ed. New York: Springer.
- Lee, W. C., Brennan, R. L. & Kolen, M. J. (2000), Estimators of Conditional Scale-Score Standard Errors of Measurement: A Simulation Study. *Journal of Educational Measurement*, 37, 1-20.
- Lord, F. M. (1965). A strong true-score theory with applications. *Psychometrika*, 30, 239-270.
- Lord, F. M., & Novick, M. R. (1968). *Statistical theory of mental test scores*. MA: Addison-wesley.
- Lord, F. M. (1969). Estimating true-score distributions in psychological testing (An empirical Bayes estimation problem). *Psychometrika*, 34, 259-299.
- Mood, M. A., Graybill, A. F. & Boes, C. D. (2008). *Introduction to the Theory of Statistics*. C.A: McGraw-Hill.
- Petersen, N. S., Kolen, M. J., & Hoover, H. D. (1989). *Scaling, norming, and equating*. In R. L. Linn (Ed.), *Educational measurement* (3rd ed., pp. 221-262). New York: American Council Education; and Macmillan.
- The SAT. (2016). *SAT technical manual*. New York: Author. Retrieved from collegereadiness.collegeboard.org.

- The ACT. (2014). *ACT assessment technical manual*. Iowa City: Author.
Retrieved from
<http://www.act.org/research/researchers/techmanuals.html>
- Woodruff, D., Traynor, A., Cui, Z., & Fang, Y. (2013). A Comparison of Three Methods for Computing Scale Score Conditional Standard Errors of Measurement. *ACT Research Report Series*, 2013 (7). ACT, Inc.

Archive of SID