

تعیین روش بهینه شناسایی کارکرد افتراقی در سنجش انطباقی کامپیوتوی

نگار شریفی یگانه^۱، محمدرضا فلسفی نژاد^۲، نورعلی فرخی^۳، احسان جمالی^۴

تاریخ دریافت: ۹۷/۰۲/۲۶

تاریخ پذیرش: ۹۷/۰۹/۱۵

چکیده

یکی از چالش‌های اساسی گذر از آزمون‌های مداد – کاغذی به انطباقی کامپیوتوی عادلانه بودن است که ارزیابی آن در چارچوب کارکرد افتراقی ضرورتی اجتناب‌ناپذیر می‌باشد. هدف مطالعه حاضر بررسی کارکرد افتراقی، ارزیابی عوامل مداخله‌گر در میزان آشکارسازی و معرفی روش بهینه مطالعه کارکرد افتراقی در سنجش انطباقی کامپیوتوی بود. با توجه به مسئله پژوهش از روش تجربی استفاده شد. گردآوری داده‌ها و دستکاری متغیرها با استفاده از روش شبیه‌سازی صورت گرفت. پاسخ‌های گروه نمونه ۱۰۰۰ نفری (گروه مرجع و کانونی با حجم یکسان ۵۰۰ نفری) به بانک ۵۵ سوالی دوارزشی براساس مدل لجستیک سه‌پارامتری در ۲۰ تکرار شبیه‌سازی شد. ۱۵ سوال بانک از نظر نوع و اندازه کارکرد افتراقی دستکاری شدند و اثر آزمون براساس تفاوت میانگین توانایی گروه‌های مقایسه تعیین گردید. آزمون انطباقی کامپیوتوی ۳۰ سوالی با نرم‌افزار Firestar اجرا شد. تحلیل کارکرد افتراقی با روش رگرسیون لجستیک و آزمون نسبت درستنایی صورت گرفت و روش‌ها براساس توان و خطای نوع اول مقایسه شدند. میزان خطای نوع اول روش آزمون نسبت درستنایی کمتر از رگرسیون لجستیک بود. توان هر دو روش متأثر از نوع، مقدار کارکرد افتراقی و اثر آزمون بود. روش آزمون نسبت درستنایی در شناسایی کارکرد افتراقی یکنواخت در هر دو موقعیت اثر و بدون اثر نسبت به روش رگرسیون لجستیک توان

۱. دانشجوی دوره دکتری سنجش و اندازه‌گیری دانشگاه علامه طباطبائی، تهران، ایران

۲. داشیار گروه سنجش و اندازه‌گیری دانشگاه علامه طباطبائی، تهران، ایران (نویسنده مسئول)

falsafinejad@yahoo.co.uk

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

۴. استادیار سازمان سنجش آموزش کشور، تهران، ایران

بیشتری داشته است و با افزایش شدت کارکرد افتراقی توان نیز افزایش یافته است. در ارزیابی کارکرد افتراقی غیریکنواخت تفاوتی بین روش‌ها مشاهده نشد و هر دو روش توان کمی داشتند. با توجه به توان و میزان خطا نوع اول، روش آزمون نسبت درستنمایی رویکرد مطلوب در بررسی کارکرد افتراقی یکنواخت است، در حالی که ارزیابی کارکرد افتراقی غیریکنواخت مستلزم مطالعات تكمیلی می‌باشد.

واژگان کلیدی: آزمون انطباقی کامپیوتری، کارکرد افتراقی سؤال ، روش رگرسیون لجستیک، روش آزمون نسبت درستنمایی سؤال - پاسخ

مقدمه

عدالت آموزشی ایجاب می‌کند آزمون‌ها برای تمامی افراد جامعه از هرجنس، نژاد، سن و موقعیت اجتماعی و اقتصادی عادلانه^۱ باشد و آزمودنی‌های متعلق به گروه‌های مختلف جامعه با توانایی‌های مشابه، امکان توفیق یکسانی در آزمون داشته باشند. هدف از اجرای آزمون‌ها تفکیک و تمایز افراد است، آزمون‌ها نبایستی منجر به تبعیض میان آزمودنی‌ها شود. دیدگاه‌های متفاوتی در خصوص عادلانه بودن آزمون وجود دارد. تاکنون تعریفی از عادلانه بودن آزمون که مورد توافق همگان باشد، ارائه نشده است. ویلنگهام و کول^۲ آزمون عادلانه را به عنوان آزمونی روا و معتر برای تمام آزمودنی‌ها تعریف کرده‌اند. به عقیده آن‌ها آزمون عادلانه می‌بایست فرصت‌های یکسانی را به منظور ارائه دانش و مهارت‌های مورد سنجش در اختیار تمامی آزمودنی‌ها قرار دهد (ناسیتروم^۳، ۲۰۰۳).

چالش اساسی عادلانه بودن آزمون و ضرورت بررسی آن به عنوان بخشی از فرایند آزمون‌سازی همزمان با آغاز عصر حقوق مدنی در راستای انتظارات مبنی بر یکسانی نتایج آزمون‌ها شکل گرفت (دانکن^۴، ۲۰۰۶). عادلانه بودن آزمون‌ها در چارچوب سوگیری سؤال‌ها^۵ مورد بررسی قرار می‌گیرد. به عقیده انگاف و فورد^۶ (۱۹۷۳) سؤال در صورتی دارای سوگیری است که شاخص دشواری آن برای یک گروه در مقایسه با گروه دیگر

1. fairness

2. Willingham & Cole

3. Näsström

4. Duncan

5. Item bias

6. Angoff & Ford

نسبتاً بیشتر یا کم تر باشد. همچنین شونمن^۱ سؤالی را سودار می‌داند که نسبت پاسخ درست افرادی که دارای توانایی یکسانی در آزمون هستند در هر زیر گروه از آزمودنی‌ها متفاوت باشد (کروکر و الجینا^۲، ۲۰۰۶، ترجمه فرزاد و زارع، ۱۳۸۸). به طور کلی سوگیری زمانی مطرح می‌شود که احتمال پاسخ‌گویی درست یک گروه از آزمودنی‌ها به دلیل مشخصات سؤال‌ها و شرایط اجرای آزمون از گروه دیگر آزمودنی‌ها کم تر یا بیشتر باشد (زومبو^۳، ۱۹۹۹).

تعیین کارکرد افتراقی سؤال روش آماری بررسی وجود سوگیری در سؤال‌های آزمون است. سؤال زمانی دارای کارکرد افتراقی است که آزمودنی‌ها با توانایی یکسان متعلق به گروه‌های مختلف، احتمال پاسخ‌گویی درست متفاوت داشته باشند (راجو^۴، ۱۹۹۰؛ هریرا و گومز^۵، ۲۰۰۸؛ همبلتون، سوامی‌ناتان و راجرز^۶، ۱۹۹۱، ترجمه فلسفی‌نژاد، ۱۳۸۹). در ادبیات کارکرد افتراقی این گروه‌ها مرجع^۷ (گروه اکثریت آزمودنی‌ها که عملکرد آن‌ها به عنوان نقطه مرجع محسوب می‌شود) و کانونی^۸ (گروه اقلیت آزمودنی‌ها که انتظار می‌رود سؤال‌های آزمون به نفع آنان نباشد) نامیده می‌شوند (سوامی‌ناتان و راجرز، ۱۹۹۰). وجود کارکرد افتراقی بیانگر آن است که عوامل مربوط به عضویت در گروه، احتمال پاسخ‌گویی درست را تحت تاثیر قرار می‌دهد. در حضور کارکرد افتراقی، وجود سوگیری سؤال باید به روش‌های مختلف مورد بررسی قرار گیرد (دانکن، ۲۰۰۶). به دلیل آن که سوگیری سؤال‌ها باعث ایجاد خطای سیستماتیک می‌شود، تحریف‌هایی در استباطه‌های مبتنی بر نمره آزمون به وجود می‌آورد و تهدیدی برای تفسیر روا نمره‌های آزمون محسوب می‌شود (کونولی^۹، ۲۰۰۳). لذا شناسایی و کنترل سوگیری سؤال‌های آزمون گام مهمی در سنجش علمی می‌باشد و باید به عنوان بخشی از فرایند تحلیل آزمون‌ها در نظر گرفته شود.

1. Scheuneman

2. Crocker & Algina

3. Zumbo

4. Raju

5. Herrera & Gómez

6. Hambleton, Swaminathan & Rogers

7. reference

8. focal

9. Conoley

مطالعه سوگیری و روش‌شناسی آن به طور قابل توجه‌ای تحت تاثیر سنجش و شرایط آن است. یکی از عوامل موثر بر سوگیری سؤال‌ها، فرمت اجرای آزمون است. با ظهور و پیشرفت فن‌آوری کامپیوتر و کاربرد آن در سنجش در چارچوب آزمون‌های کامپیوتری و انطباقی کامپیوتری^۱ (CAT) مسئله عادلانه بودن این نوع از آزمون‌ها به طور جدی مطرح شد. در آزمون انطباقی کامپیوتری، سؤال‌ها متناسب با توانایی آزمودنی‌ها در جریان آزمون انتخاب و اجرا می‌شوند (ریکیسی^۲، ۱۹۸۹، به نقل از دیویدسون^۳، ۲۰۰۳). آزمون‌های کامپیوتری خصوصاً انطباقی کامپیوتری با توجه به توانمندی بالقوه‌ای که در بهینه‌سازی فرایند سنجش دارند، جایگزین جدی آزمون‌های مداد-کاغذی قلمداد می‌شوند. بسیاری از سازمان‌ها و مؤسسات آزمون‌سازی بین‌المللی نظام سنجش و ارزشیابی خود را از فرمت مداد - کاغذی به کامپیوتری و انطباقی کامپیوتری تغییر داده‌اند یا در حال مطالعه و برنامه‌ریزی جهت انجام این تغییر هستند. آزمون تحصیلات تکمیلی^۴ (GRE)، آزمون پذیرش تحصیلات تکمیلی دوره مدیریت (GMAT)^۵، آزمون مجوز پرستاری^۶ (NCLEX)، مجموعه آزمون‌های استعداد شغلی ارتش^۷ و آزمون تافل^۸ از جمله آزمون‌های خطیری هستند که به صورت انطباقی کامپیوتری اجرا می‌شوند (کلندر^۹، ۲۰۱۱).

این در حالی است که گذر از آزمون‌های مداد-کاغذی به کامپیوتری و انطباقی کامپیوتری فرایندی پر مخاطره با چالش‌های نظری و عملی فراوان است. عادلانه بودن آزمون‌ها یکی از مهم‌ترین چالش‌های فرایند تغییر روش اجرای آزمون‌ها است که ضرورت ارزیابی آن در چارچوب کارکرد افتراقی سؤال‌ها مورد تأکید قرار گرفته است (انجمان تحقیقات آموزشی آمریکا^{۱۰}، انجمان روان‌شناسی آمریکا^{۱۱} و انجمان ملی سنجش

-
1. Computerized Adaptive Testing
 2. Reckase
 3. Davidson
 4. Graduate Record Examination
 5. Graduate Management Admission Test
 6. Nursing Licensing Examination (NCLEX)
 7. Armed Services Vocational Aptitude Battery (ASVAB)
 8. TOEFL
 9. Kalender
 10. American Educational Research Association (AERA)
 11. American Psychological Association (APA)

آموزشی^۱، ۱۹۹۹). اهمیت بررسی کارکرد افتراقی در آزمون‌های انطباقی کامپیوتری به دلیل وجود چندین منبع بالقوه کارکرد افتراقی و همچنین کم بودن تعداد سؤال‌ها دوچندان می‌شود. کم بودن تعداد سؤال‌ها در آزمون انطباقی کامپیوتری باعث می‌شود پاسخ به هر سؤال نقش مهم‌تری در نمره آزمودنی داشته باشد و هرگونه عیب و کاستی در سؤال‌ها عواقب جدی‌تری برای آزمودنی به همراه داشته باشد.

اگرچه روش‌های مطالعه کارکرد افتراقی سؤال ادبیات نسبتاً وسیعی را در متون سنجش به خود اختصاص داده است اما روش‌های ارائه شده بیشتر در آزمون‌های مداد – کاغذی که به لحاظ اجرایی دارای فرم بسته هستند کاربرد دارد. بررسی کارکرد افتراقی سؤال‌ها در آزمون‌های انطباقی کامپیوتری با چالش دشواری تعیین متغیر تطبیق^۲ به منظور جور کردن آزمودنی‌های گروه‌های مقایسه همراه است. در آزمون‌های انطباقی کامپیوتری تمامی آزمودنی‌ها به سؤال‌های یکسانی پاسخ نمی‌دهند، این مسئله منجر به نمره‌های کل متفاوتی می‌شود. بنابراین کاربرد نمره کل به عنوان متغیر تطبیق در سنجش انطباقی کامپیوتری بی معنا است (پیروم‌سومبات^۳، ۲۰۱۴؛ زویک^۴، ۲۰۱۰). به علاوه در سنجش انطباقی کامپیوتری ارزیابی کارکرد افتراقی مستلزم انتخاب روشی است که در نمونه‌های کوچک نتایج پایداری ارائه کند. اگرچه تعداد کل آزمودنی‌ها در آزمون انطباقی کامپیوتری زیاد است، تعداد پاسخ‌ها برای برخی از سؤال‌های آزمون ممکن است خیلی کم باشد (زویک، ۲۰۱۰). بر این اساس شناسایی کارکرد افتراقی سؤال در آزمون‌های انطباقی کامپیوتری مستلزم به کارگیری روش‌های جدید و مخصوص می‌باشد (پیروم‌سومبات، ۲۰۱۴؛ زویک، ۲۰۱۰).

رویکردهای فعلی تعیین کارکرد افتراقی در سنجش انطباقی کامپیوتری مبتنی بر روش‌های به کار رفته در آزمون‌های مداد – کاغذی است. به عنوان مثال زویک و همکارانش (زویک، تایر و وینگرسکی^۵، ۱۹۹۴ و زویک، تایر و لویس^۶، ۱۹۹۷) روش‌های زویک، تایر و وینگرسکی (ZTW) و زویک، تایر و لویس (ZTL) را براساس آماره

1. National Council on Measurement in Education (NCME)

2. matching variable

3. Piromsombat

4. Zwick

5. Thayer & Wingersky

6. Lewis

مانتل - هنزل^۱ (هالند^۲ و تایر، ۱۹۸۸) و روش استاندارسازی^۳ (دورانس و کولیک^۴، ۱۹۸۶) مطرح کردند. در این روش‌ها آزمودنی‌های گروه مرجع و کانونی براساس نمره واقعی مبتنی بر برآورد توانایی آزمون انطباقی کامپیوترا جور می‌شوند. تفاوت این دو روش مربوط به استفاده از برآورد بیزین تجربی^۵ در روش ZTL برای ارائه نتایج با ثبات تر در نمونه‌های کوچک می‌باشد (پیرومسمبات، ۲۰۱۴؛ زویک، ۲۰۱۰). هم‌چنین نانداکومار و رووسوس^۶ (۲۰۰۱) روش SIBTEST (شیلی و استوت^۷، ۱۹۹۳) را برای سنجش انطباقی کامپیوترا توسعه بخشیدند که منجر به ایجاد روش CATSIB شد. در روش آزمودنی‌ها براساس برآورد توانایی آزمون انطباقی کامپیوترا جور می‌شوند، سپس روش SIBTEST برای ارزیابی کارکرد افتراقی به کار می‌رود. روش‌های ZTL، ZTW و CATSIB صرفاً برای شناسایی کارکرد افتراقی یکنواخت^۸ (هماهنگ) طراحی شده‌اند. لی، چن و یو^۹ (۲۰۰۶) نسخه اصلاح شده رگرسیون لجستیک^{۱۰} و آزمون نسبت درستنمایی سؤال-پاسخ^{۱۱} را به منظور بررسی کارکرد افتراقی غیریکنواخت^{۱۲} (غیرهماهنگ) در آزمون انطباقی کامپیوترا در سؤال‌های پیش‌آزمون معرفی کردند. در کارکرد افتراقی یکنواخت، تفاوت احتمال پاسخ‌گویی درست گروه کانونی و مرجع در تمام سطوح توانایی یکسان است و تعامل بین سطح توانایی و عضویت در گروه وجود ندارد. در حالی که در کارکرد افتراقی غیریکنواخت احتمال پاسخ‌گویی درست گروه کانونی و مرجع در تمام سطوح توانایی یکسان نمی‌باشد و بین سطح توانایی و عضویت در گروه تعامل وجود دارد (دریانا^{۱۳}، ۲۰۰۷).

-
1. Mantel - Haenszel
 2. Holland & Thayer
 3. Standardization
 4. Dorans & Kulick
 5. Empirical bayesian estimation
 6. Nandakumar & Roussos
 7. Shealy & Stout
 8. uniform
 9. Lei, Chen & Yu
 10. Logistic Regression (LR)
 11. Item Response Theory-Likelihood Ratio Tests (IRT-LRT)
 12. non uniform
 13. Driana

از آنجایی که بیشتر الگوریتم‌های آزمون انطباقی کامپیوتروی مستلزم برآورد مبتنی بر نظریه سؤال – پاسخ است به نظر می‌رسد روش نسبت درستنمایی انتخاب منطقی برای تحلیل کارکرد افتراقی در چارچوب آزمون‌های انطباقی کامپیوتروی باشد (Miller^۱، ۱۹۹۲). کاربرد آزمون نسبت درستنمایی سؤال – پاسخ در ارزیابی کارکرد افتراقی توسط تیسن، استینبرگ و واینر^۲ (۱۹۹۳) مطرح شد. در این روش تفاوت پارامتر سؤال‌ها بین گروه‌ها مشروط بر ناتغیر بودن سایر سؤال‌های آزمون مدل‌سازی می‌شود. برآش مدل افزایشی^۳ و مدل فشرده^۴ براساس تفاوت لگاریتم درستنمایی دو مدل مورد مقایسه قرار می‌گیرد. در مدل افزایشی پارامترهای سؤال‌های مورد بررسی امکان تغییر در گروه‌های مورد مقایسه را دارند، در حالی که در مدل فشرده پارامترهای سؤال محدود می‌شوند (کوهن، کیم و ولак^۵، ۱۹۹۶؛ لی، چن و یو، ۲۰۰۶). اگر تفاوت لگاریتم درستنمایی دو مدل که دارای توزیع مجدور کای است معنادار باشد، آزمون‌های بعدی به منظور بررسی تفاوت‌های پارامتر سؤال‌ها انجام می‌شود (پیرومومبات، ۲۰۱۴).

یک مزیت روش نسبت درستنمایی امکان کاربرد مستقیم آن در آزمون انطباقی کامپیوتروی است. روش نسبت درستنمایی برخلاف سایر روش‌ها مبتنی بر نمره کل آزمون نیست (کیم، ۲۰۰۶). با وجود مطلوبیت‌های نظری این روش کاربرد آن به دلیل پیچیدگی‌های ناشی از داده‌های ناقص حاصل از اجرای انطباقی آزمون‌ها، با مشکلاتی مواجه شد. هم‌چنین برآش مدل‌های متعدد و محاسبات سنگین مورد نیاز روش نسبت درستنمایی استفاده از آن را خصوصاً در مطالعات شبیه‌سازی غیرعملی ساخت. تیسن (۲۰۰۱) نرم‌افزار IRTLRDIF را به منظور افزایش کاربرد روش نسبت درستنمایی در تحلیل کارکرد افتراقی آزمون‌های انطباقی کامپیوتروی ارائه کرد. این نرم‌افزار با تسريع انجام محاسبات امکان استفاده از روش نسبت درستنمایی را در مطالعات شبیه‌سازی ممکن ساخته است (لی، چن و یو، ۲۰۰۶).

1. Miller

2. Thissen, Steinberg & Wainer

3. augmented

4. compact

5. Cohen, Kim & Wollack

روش رگرسیون لجستیک نیز با توجه به قابلیت شناسایی کارکرد افتراقی غیریکنواخت می‌تواند در چارچوب آزمون انطباقی کامپیوترا مورد استفاده قرار گیرد. رگرسیون لجستیک مبتنی بر مدل‌سازی آماری احتمال پاسخ صحیح به سؤال براساس عضویت در گروه و ملاک است. در این روش پاسخ سؤال‌ها به عنوان متغیر وابسته در نظر گرفته می‌شود، متغیر گروهی به عنوان متغیری که قصد بررسی پاسخ‌ها را در آن دارند به عنوان گروه‌های مرجع و کانونی در نظر گرفته می‌شود. نمره کل آزمون نیز به عنوان ملاک به کار می‌رود. در روش رگرسیون لجستیک وجود کارکرد افتراقی از طریق بررسی بهبود ایجاد شده در برآشش مدل رگرسیون پس از اضافه کردن عضویت در گروه و تعامل بین نمره آزمون و عضویت در گروه تعیین می‌شود. تعامل به منظور تعیین کارکرد افتراقی غیریکنواخت استفاده می‌شود (ملیسپ^۱، ۲۰۱۱ به نقل از گرامی‌پور، ۱۹۹۰). در نسخه اصلاح شده رگرسیون لجستیک به منظور بررسی کارکرد افتراقی در آزمون‌های انطباقی کامپیوترا، برآورد توانایی آزمون انطباقی کامپیوترا در معادله رگرسیون لجستیک جایگزین می‌شود (پیرو مسومبات، ۲۰۱۴).

به طور کلی عملکرد روش‌های شناسایی کارکرد افتراقی سؤال‌ها مشابه نمی‌باشد و تحت تاثیر عوامل مختلفی از جمله شدت کارکرد افتراقی، توزیع توانایی گروه‌های مقایسه، طول آزمون، حجم نمونه و نسبت سؤال‌های دارای کارکرد افتراقی می‌باشد (کیم، ۲۰۱۰؛ کاباساکال، گوک، آرسان و کلیسی‌گل، ۲۰۱۴). کارایی روش‌های شناسایی کارکرد افتراقی در سنجش انطباقی کامپیوترا و عوامل موثر بر آن در مطالعات محدودی مورد بررسی قرار گرفته است. به عنوان مثال زویک، تایر و وینگرسکی (۱۹۹۴) ویژگی‌های روش ZTW را با شبیه‌سازی پاسخ‌ها به سه بانک ۷۵ سؤالی (بانک اول سؤال‌ها فاقد کارکرد افتراقی، بانک دوم سؤال‌ها دارای کارکرد افتراقی ناهمبسته با دشواری سؤال و بانک سوم سؤال‌ها دارای کارکرد افتراقی همبسته با دشواری) ارزشیابی کردند. به هر آزمودنی ۲۵ سؤال از هر یک از بانک‌ها ارائه شد. آزمودنی‌های شبیه‌سازی شده گروه‌های مرجع و کانونی براساس نمره واقعی مورد انتظار برای کل بانک جور شدند (زویک، ۲۰۱۰). حجم گروه‌های مرجع و کانونی، توزیع توانایی گروه‌ها و الگوی کارکرد افتراقی متغیرهای بودند که دستکاری شدند. عملکرد روش ZTW با آماره مانتل هنزل مرسوم

1. Millsap

مورد مقایسه قرار گرفت. نتایج بیانگر آن بود که برآوردهای کارکرد افتراقی روش ZTW همبستگی بالایی با مقدار کارکرد افتراقی واقعی دارد. به علاوه میزان خطای نوع اول روش ZTW در اکثر موقعیت‌های شبیه‌سازی شده قابل قبول بود (پیرو مسومبات، ۲۰۱۴؛ زویک، ۲۰۱۰، زویک، تایر و وینگرسکی، ۱۹۹۴). زویک و تایر (۲۰۰۲) عملکرد روش ZTL را با شبیه‌سازی داده‌ها با استفاده از بانک سؤال، حجم نمونه، طول آزمون، اندازه کارکرد افتراقی یکنواخت و اندازه اثر آزمون (تفاوت میانگین توانایی گروه‌ها) مورد بررسی قرار دادند. نتایج بیانگر آن بود که روش ZTL در شناسایی کارکرد افتراقی یکنواخت در سؤال‌های پیش‌آزمون خصوصاً در نمونه‌های کوچک موثر است.

نانداکومار و روسوس (۲۰۰۴؛ ۲۰۰۱) عملکرد روش CATSIB را در شناسایی کارکرد افتراقی در سؤال‌های پیش‌آزمون با استفاده از روش شبیه‌سازی بررسی کردند. نتایج بیانگر رضایت‌بخش بودن میزان خطای نوع اول CATSIB با تصحیح رگرسیون بود. توان این روش دامنه گسترده‌ای (۱۷/۰ تا ۱) داشت و با افزایش حجم گروه نمونه افزایش یافته بود. حساسیت روش CATSIB در شرایط کوچک بودن حجم هر یک از گروه‌ها، نابرابر بودن حجم گروه‌ها و زیاد بودن اندازه اثر آزمون (تفاوت میانگین توزیع توانایی گروه مرجع و کانونی) کم بود

لی، چن و یو (۲۰۰۶) نیز عملکرد CATSIB، رگرسیون لجستیک و آزمون نسبت درستنمایی سؤال – پاسخ را به طور همزمان با استفاده از مطالعه شبیه‌سازی در سؤال‌های پیش‌آزمون مورد بررسی قرار دادند. داده‌ها تحت شرایط نسبت حجم نمونه و بزرگی اثر آزمون شبیه‌سازی شدند. نتایج مطالعه نشان داد که روش CATSIB، رگرسیون لجستیک و آزمون نسبت درستنمایی از نظر توان شناسایی کارکرد افتراقی سوال یکنواخت قابل مقایسه هستند در حالی که در شناسایی کارکرد افتراقی غیریکنواخت روش‌های رگرسیون لجستیک و آزمون نسبت درستنمایی نسبت به CATSIB توانمندتر هستند. آزمون نسبت درستنمایی کنترل خوبی بر میزان خطای نوع اول در شرایط حجم نمونه نابرابر و اثر آزمون داشته است.

با توجه به مطالب مطرح شده می‌توان نتیجه گرفت که با وجود اهمیت عادلانه بودن آزمون‌های انطباقی کامپیوتری و ضرورت ارزیابی آن در چارچوب کارکرد افتراقی، مطالعات کمی در این حوزه انجام شده است. به علاوه در تمامی مطالعات، کارکرد افتراقی

در سؤال‌های پیش‌آزمون با جور کردن آزمودنی‌ها براساس توانایی مبتنی بر سنجش انطباقی کامپیوتري مورد بررسی قرار گرفته است. سؤال‌های پیش‌آزمون به صورت غیرانطباقی اجرا می‌شوند و بنابراین نمی‌تواند بیانگر شرایط و موقعیت کارکرد افتراقی در بافت آزمون‌های انطباقی کامپیوتري باشند. کم‌تر مطالعه‌ای به طور خاص کارکرد افتراقی را در سؤال‌های عملیاتی در جریان آزمون انطباقی کامپیوتري مورد بررسی قرار داده است. از این رو اهمیت مطالعه در این حوزه کاملاً مشهود است.

در جامعه ما نیز ضرورت بهینه‌سازی آزمون‌ها و کاربرد فن‌آوری کامپیوتري در سنجش با توجه به حجم بالاي کاربرد آزمون‌ها در موقعیت‌های مختلف به منظور اخذ تصمیمات سرنوشت‌ساز کاملاً احساس می‌شود. آزمون‌های انطباقی کامپیوتري از جمله حوزه‌هایي آزمون‌سازی محسوب می‌شوند که با خلاء مطالعات نظامند روبه رو هستند. از اين رو انجام مطالعه در خصوص چالش‌های تغیير روش اجرای آزمون‌ها از فرمت مداد- کاغذی به کامپیوتري و خصوصاً انطباقی کامپیوتري و بررسی راه حل‌های ممکن پيش از هر گونه تصمیم‌گيری برای تغیير فرمت آزمون‌ها بسيار حائز اهمیت است. در راستاي اين هدف، پژوهش حاضر به عنوان يكى از اولين مطالعات به بررسی کارکرد افتراقی در سؤال‌های عملیاتی در اجرای انطباقی کامپیوتري براساس روش آزمون نسبت درستنمايی سؤال - پاسخ و رگرسيون لجستيك پرداخته است. اين دو روش به دليل توانمندي در شناسايي کارکرد افتراقی يکنواخت و غيريکنواخت و همچنين ويژگي‌های نظری و عملی که پيشتر به آن اشاره شد انتخاب شدند.

روش

به منظور ارزیابی کارکرد افتراقی سؤال‌های عملیاتی در اجرای انطباقی کامپیوتري و مقایسه دو روش رگرسيون لجستيك و آزمون نسبت درستنمايی از روش‌های شبیه‌سازی مونت‌کارلو^۱ و پس‌تجربی^۲ (واقعي) برای ایجاد داده‌هایي مشابه با داده‌های واقعي، دستکاری متغيرهای مورد بررسی و اجرای آزمون انطباقی کامپیوتري استفاده شد. براین اساس پژوهش حاضر از نظر روش، تحقیق تجربی است. کاربرد روش‌های شبیه‌سازی در

1. Monte Carlo Simulations
2. Post-hoc Simulations

تحقیقات سنجش آموزشی بسیار متداول شده است و نتایج آن‌ها معمولاً به عنوان شواهد روایی برای تصمیم‌گیری درباره برنامه‌های سنجش در نظر گرفته می‌شوند. سنجش انطباقی کامپیوتروی یکی از حوزه‌هایی است که به طور فراینده‌ای با استفاده از روش‌های شبیه‌سازی مورد بررسی قرار می‌گیرد (هارول، استون، هسو و کرسکی^۱، ۱۹۹۶).

داده‌ها با استفاده از روش شبیه‌سازی مونت کارلو تولید شدند. روش کار به این صورت است که پارامتر توانایی آزمودنی‌ها و پارامترهای سؤال‌ها براساس توزیع‌های تعیین شده و با توجه به مدل سؤال – پاسخ انتخابی ایجاد شدند. در ادامه ماتریس پاسخ سؤال‌ها برای هر ترکیب سؤال و آزمودنی با مقایسه احتمال پاسخ‌گویی درست به سؤال با اعداد تصادفی استخراج شده از توزیع یکنواخت تشکیل شد. در صورتی که احتمال پاسخ‌گویی محاسبه شده بزرگ‌تر از عدد تصادفی بود نمره ۱ و در صورتی که کوچک‌تر از عدد تصادفی بود نمره ۰ در نظر گرفته شد. (نیدیک و ویس، ۲۰۰۹).

بر این اساس در مطالعه حاضر پاسخ به سؤال‌های بانک ۵۵ سؤالی دوارزشی براساس مدل لجستیک سه پارامتری با استفاده از برنامه WinGen (هان^۲، ۲۰۰۷) در ۱۸ شرایط آزمایشی مختلف با توجه به عوامل مداخله‌گر شبیه‌سازی شد. ویژگی‌های ۱۸ مجموعه داده حاصل در جدول شماره ۱ ارائه شده است. منطبق بر مطالعات صورت گرفته، پارامتر توانایی آزمودنی‌ها از توزیع نرمال با میانگین ۰ و انحراف استاندارد ۱ ($N(0, 1)$)، پارامتر دشواری سؤال‌ها از توزیع یکنواخت (۳، -۳)، پارامتر تشخیص از توزیع لگ نرمال با میانگین $1/0$ و انحراف استاندارد $1/0$ ($LnN(-0/1, 0/1)$) و پارامتر حدس از توزیع یکنواخت (۰، ۰/۲۵) انتخاب شدند. توزیع پارامتر دشواری سؤال‌ها به منظور پوشش دامنه وسیعی از سطوح توانایی واقعی انتخاب شده است و منجر به تابع آگاهی قابل مقایسه‌ای در سطوح توانایی می‌شود. این توزیع‌ها دامنه وسیعی از پارامترهای سؤال‌ها را که در عمل و در شبیه‌سازی‌های قبلی آزمون انطباقی کامپیوتروی به کار رفته‌اند را تولید می‌کنند.

1. Harwell, Stone, Hsu & Kirisci
2. Han

(پیرومومبات، ۲۰۱۴؛ چانگ و ینگ^۱، ۱۹۹۶؛ زویک، ۱۹۹۴؛ گیور^۲، ۲۰۰۸؛ وانگ و ویسپول^۳، ۱۹۹۸).

۱۵ سؤال بانک از نظر نوع کارکرد افتراقی در سه سطح بدون کارکرد افتراقی، دارای کارکرد افتراقی یکنواخت و کارکرد افتراقی غیریکنواخت دستکاری شدند. در کارکرد افتراقی یکنواخت سؤال‌ها از لحاظ سطح دشواری در دو گروه متفاوت بودند، اما از نظر قدرت تشخیص تفاوتی نداشتند. در کارکرد افتراقی غیریکنواخت سؤال‌ها صرفاً از نظر قدرت تشخیص در دو گروه متفاوت بودند. علاوه بر نوع کارکرد افتراقی، اندازه کارکرد افتراقی و اثر آزمون نیز دستکاری شد. اندازه کارکرد افتراقی در چهار سطح $0/4$ ، $0/6$ ، $0/8$ و 1 دستکاری شد. اثر آزمون نیز به صورت تفاوت میانگین توانایی گروه‌ها تعریف می‌شود. در شرایط داده‌های واقعی یافتن گروه‌های مرجع و کانونی با توزیع توانایی یکسان دشوار می‌باشد. تفاوت توزیع توانایی گروه‌ها در شناسایی کارکرد افتراقی سؤال‌ها اثر گذار است (کاباساکال و همکاران، ۲۰۱۴). وقتی که اثر آزمون وجود ندارد هر دو گروه مرجع و کانونی از جامعه‌ای با توزیع نرمال استاندارد با میانگین 0 و انحراف استاندارد 1 ($N=100$) نمونهبرداری شدند. هنگامی که اثر آزمون وجود دارد فرض می‌شود میانگین جامعه گروه کانونی یک انحراف استاندارد زیر میانگین جامعه گروه مرجع ($N=11$) است.

به منظور کاهش سوگیری بالقوه چندین مجموعه داده شبیه‌سازی شد. تعداد تکرارها با توجه به هدف مطالعه، میزان کاهش مورد نظر در واریانس نمونه‌گیری پارامتر و توان مورد نیاز آزمون‌های آماری انتخاب می‌شود. در مطالعات مقایسه روش‌شناسی رویکرد سؤال – پاسخ به عنوان مثال بررسی کارکرد افتراقی سؤال‌ها تعداد تکرارهای کمتری لازم است. هارول و همکاران (۱۹۹۶) حداقل ۱۰ تکرار را در این موارد توصیه کرده‌اند. براین اساس در مطالعه حاضر برای کسب اطمینان از حداقل بودن سوگیری نمونه و به منظور دستیابی به توان آماری مطلوب ۲۰ تکرار برای یک از شرایط آزمایشی صورت گرفت.

1.Chang & Ying

2. Guyer

3. Wang & Vispoel

جدول ۱. ویژگی‌های شرایط آزمایشی

شرط	نوع کارکرد افتراقی	توزیع گروه کانونی	مقدار کارکرد افتراقی
۱	بدون کارکرد افتراقی	$N(0,1)$	-
۲	بدون کارکرد افتراقی	$N(-1,1)$	-
۳	کارکرد افتراقی یکنواخت	$N(0,1)$	۰/۴
۴	کارکرد افتراقی یکنواخت	$N(0,1)$	۰/۶
۵	کارکرد افتراقی یکنواخت	$N(0,1)$	۰/۸
۶	کارکرد افتراقی یکنواخت	$N(0,1)$	۱
۷	کارکرد افتراقی یکنواخت	$N(-1,1)$	۰/۴
۸	کارکرد افتراقی یکنواخت	$N(-1,1)$	۰/۶
۹	کارکرد افتراقی یکنواخت	$N(-1,1)$	۰/۸
۱۰	کارکرد افتراقی یکنواخت	$N(-1,1)$	۱
۱۱	کارکرد افتراقی غیر یکنواخت	$N(0,1)$	۰/۴
۱۲	کارکرد افتراقی غیر یکنواخت	$N(0,1)$	۰/۶
۱۳	کارکرد افتراقی غیر یکنواخت	$N(0,1)$	۰/۸
۱۴	کارکرد افتراقی غیر یکنواخت	$N(0,1)$	۱
۱۵	کارکرد افتراقی غیر یکنواخت	$N(-1,1)$	۰/۴
۱۶	کارکرد افتراقی غیر یکنواخت	$N(-1,1)$	۰/۶
۱۷	کارکرد افتراقی غیر یکنواخت	$N(-1,1)$	۰/۸
۱۸	کارکرد افتراقی غیر یکنواخت	$N(-1,1)$	۱

شبیه‌سازی پس تجربی آزمون‌های انطباقی کامپیوتروی ۳۰ سؤالی (با طول ثابت) براساس بانک سؤال تولید شده و پارامتر توانایی آزمودنی‌ها با استفاده از نرم‌افزار Firestar چویی و پودرابسکی و مکینی، ۲۰۱۱) انجام شد. طول آزمون مشابه با آزمون‌های به کار رفته در مطالعات قبلی که در دامنه ۱۵ تا ۳۰ سؤال بودند، انتخاب شده است (زویک، تایر و وینگرسکی، ۱۹۹۴؛ لی، چن و یو، ۲۰۰۶؛ نامداکومار و روسوس، ۲۰۰۴). فرایند شبیه‌سازی به این ترتیب است که براساس قاعده شروع و ملاک گزینش تعیین شده، سؤالی از بانک انتخاب و به آزمودنی ارائه می‌شود. سپس پاسخ آزمودنی به سؤال که بر اساس پارامتر توانایی و سؤال‌ها تولید شده است مورد بررسی قرار می‌گیرد و با توجه به پاسخ و روش برآورد توانایی تعیین شده توانایی آزمودنی برآورد می‌شود. در ادامه براساس پاسخ آزمودنی و ملاک گزینش سؤال بعدی از بانک انتخاب می‌شود و مجدداً پاسخ جدید مورد بررسی قرار می‌گیرد و برآورد توانایی آزمودنی اصلاح می‌شود. این فرایند تا زمان

تحقیق ملاک‌های خاتمه آزمون ادامه می‌باید. در این مرحله آزمون خاتمه می‌یابد و برآوردهای توانایی آزمودنی ارائه می‌شود. ۷۲۰ شبیه‌سازی آزمون انطباقی کامپیوتری انجام شد. برای تمام آزمودنی‌ها نقطه شروع آزمون توانایی صفر در نظر گرفته شد و گزینش سؤال‌ها براساس بیشینه آگاهی انجام گرفت. در این مرحله توانایی آزمودنیها با استفاده از روش بیزین پسین مورد انتظار (EAP) برآورد شد.

از آنجایی که در آزمون‌های انطباقی کامپیوتری هر آزمودنی به مجموعه‌ای متفاوت از سؤال‌ها به تناسب توانایی خود پاسخ می‌دهد، پاسخ‌های جمع‌آوری شده هر آزمودنی بخش کوچکی از داده‌هایی است که در آزمون غیرانطباقی متناظر به دست می‌آید. به عبارت دیگر ماتریس داده‌ها ناکامل می‌باشد. پیش از تحلیل کارکرد افتراقی سؤال‌ها می‌بایست ماتریس داده‌ها کامل شود. تکمیل ماتریس داده‌ها به روش جایگذاری براساس احتمال پاسخ‌گویی درست به سؤال‌ها مبتنی بر برآورد توانایی آزمون انطباقی کامپیوتری برای هر آزمودنی صورت گرفت. ابتدا با توجه به این که هر یک از آزمودنی‌ها به مجموعه سؤال‌های متفاوتی در شبیه‌سازی آزمون انطباقی کامپیوتری پاسخ داده‌اند، ۱۴۴۰ فایل "سؤال‌های ارائه شده" و "پاسخ سؤال‌ها" (۱۸ شرایط مطالعه ۲ گروه مرجع و کانونی ۲۰ تکرار ۲ مجموعه فایل) با استفاده از سامانه مدیریت پایگاه داده‌ها متن باز MySQL و برنامه‌نویسی PHP مرتب شد و به این ترتیب ماتریس اولیه داده‌ها تشکیل شد. در ادامه احتمال پاسخ‌گویی صحیح به سؤال‌ها براساس مدل لجستیک سه پارامتری با استفاده از پارامتر سؤال‌ها و برآوردهای توانایی آزمون انطباقی کامپیوتری برای هر مجموعه آزمودنی‌های شبیه‌سازی شده محاسبه شد. سپس مقادیر احتمال پاسخ‌گویی به دست آمده با اعداد تصادفی از توزیع یکنواخت (۰ و ۱) (مورد مقایسه قرار گرفت. در صورتی که مقادیر احتمال محاسبه شده مساوی و یا بزرگ‌تر از عدد تصادفی بود، پاسخ ۱ و در صورتی که کوچک‌تر بود پاسخ ۰ لحاظ شد و به این ترتیب ماتریس داده‌ها جایگذاری و کامل شد.

پس از تشکیل ماتریس کامل داده‌ها کارکرد افتراقی سؤال‌های دستکاری شده با استفاده از روش رگرسیون لجستیک و آزمون نسبت درستنمایی سؤال – پاسخ تحلیل شدند. در روش رگرسیون لجستیک نمره واقعی مبتنی بر برآورد توانایی آزمون انطباقی کامپیوتری به عنوان متغیر تطبیق مورد استفاده قرار گرفت. در روش آزمون نسبت

درستنمایی نیز سایر سؤال‌های بانک به عنوان سؤال‌های لنگر در نظر گرفته شدند. در مطالعه حاضر مقایسه دو رویکرد رگرسیون لجستیک و آزمون نسبت درستنمایی در شناسایی کارکرد افتراقی براساس توان (نسبت مثبت واقعی) و خطای نوع اول (نسبت مثبت کاذب) در ۲۰ تکرار برای هر موقعیت صورت گرفت. در صورتی که سؤال دارای کارکرد افتراقی باشد و روش‌های شناسایی کارکرد افتراقی به درستی کارکرد افتراقی را مشخص سازند، در این حالت کارکرد افتراقی مثبت واقعی تشخیص داده شده است. حال در صورتی که سؤال فاقد کارکرد افتراقی باشد ولی روش‌های شناسایی کارکرد افتراقی به اشتباه کارکرد افتراقی را در سؤال مشخص سازند در این حالت کارکرد افتراقی منفی کاذب تشخیص داده شده است.

یافته‌ها

توان آماری کارکرد افتراقی یکنواخت (نسبت مثبت واقعی): توان آماری شناسایی کارکرد افتراقی یکنواخت (نسبت مثبت واقعی) دو رویکرد رگرسیون لجستیک و آزمون نسبت درستنمایی در ۲۰ تکرار برای سطوح مختلف کارکرد افتراقی و تحت شرایط اثر و بدون اثر آزمون در جداول شماره ۲ و ۳ ارائه شده است. هر دو روش رگرسیون لجستیک و آزمون نسبت درستنمایی هنگامی که شدت کارکرد افتراقی بیشتر است و توزیع توانایی گروه مرجع و کانونی یکسان است (بدون اثر آزمون) توان بیشتری در شناسایی سؤال‌های دارای کارکرد افتراقی یکنواخت داردند. با افزایش شدت کارکرد افتراقی توان روش رگرسیون لجستیک و آزمون نسبت درستنمایی در هر دو موقعیت اثر و بدون اثر آزمون افزایش یافته است.

مقادیر توان مساوی و بزرگ‌تر از ۰/۸ بیانگر توانمندی روش در شناسایی کارکرد افتراقی است (کاباساکال و همکاران، ۱۴؛ ۲۰۱۴؛ لی، چن و یو، ۲۰۰۶). در روش رگرسیون لجستیک در شرایط بدون اثر هنگامی که اندازه کارکرد افتراقی ۰/۴ (کم) است نسبت مثبت واقعی (توان) برای هیچ یک از سؤال‌های آزمون مساوی و بزرگ‌تر از ۰/۸ نیست. هنگامی که اندازه کارکرد افتراقی برابر یک است نسبت مثبت واقعی ۷ سؤال مساوی و بزرگ‌تر از ۰/۸ است. این نتایج بیانگر عدم کفایت روش رگرسیون لجستیک در شرایط کارکرد افتراقی کم می‌باشد. در روش نسبت درستنمایی در شرایط بدون اثر هنگامی که

اندازه کار کرد افتراقی ۰/۴ است نسبت مثبت واقعی یک سؤال مساوی و بزرگتر از ۰/۸ است. حال آن که با افزایش شدت کار کرد افتراقی تعداد سؤال‌هایی که نسبت مثبت واقعی مساوی و بزرگتر از ۰/۸ دارند افزایش می‌یابد.

جدول ۲. توان کار کرد افتراقی یکنواخت شرایط بدون اثر

آزمون نسبت درستنما					رگرسیون لجستیک			سوال
۱	۰/۸	۰/۶	۰/۴	۱	۰/۸	۰/۶	۰/۴	
۰/۹۵	۱	۰/۹	۰/۴۵	۱	۰/۸۵	۰/۷۵	۰/۴	۴۱
۱	۰/۹۵	۰/۸	۰/۶	۰/۹۵	۰/۹	۰/۷	۰/۴۵	۴۲
۰/۸۵	۰/۸۵	۰/۷۵	۰/۳۵	۰/۸۵	۰/۷۵	۰/۴	۰/۳۵	۴۳
۰/۲	۰/۱	۰/۱	۰/۰۵	۰/۱۵	۰/۱	۰/۱	۰/۰۵	۴۴
۰/۸	۰/۶۵	۰/۵	۰/۲۵	۰/۵	۰/۴	۰/۲۵	۰/۲	۴۵
۰/۹	۱	۰/۹	۰/۷	۱	۰/۹۵	۰/۹	۰/۴۵	۴۶
۰/۷	۰/۶۵	۰/۳	۰/۳	۰/۳۵	۰/۳۵	۰/۳	۰/۲	۴۷
۰/۷	۰/۳۵	۰/۳۵	۰/۱۵	۰/۳۵	۰/۲	۰/۱۵	۰/۱	۴۸
۰/۸	۰/۸	۰/۸	۰/۵۵	۰/۹۵	۰/۸	۰/۵	۰/۳۵	۴۹
۰/۷	۰/۴۵	۰/۵	۰/۵	۰/۳۵	۰/۲۵	۰/۲۵	۰/۲	۵۰
۰/۷۵	۰/۵	۰/۴۵	۰/۱	۰/۳۵	۰/۳	۰/۲۵	۰/۱	۵۱
۱	۱	۱	۰/۸	۱	۱	۰/۷۵	۰/۶۵	۵۲
۰/۷۵	۰/۷۵	۰/۴	۰/۲	۰/۶۵	۰/۳	۰/۱	۰/۱	۵۳
۰/۶۵	۰/۷۵	۰/۴	۰/۲	۰/۷	۰/۵	۰/۲۵	۰/۱	۵۴
۱	۰/۹۵	۰/۹	۰/۶	۱	۰/۹۵	۰/۸۵	۰/۴۵	۵۵
۰/۷۸	۰/۷۲	۰/۶۱	۰/۳۸	۰/۶۸	۰/۵۷	۰/۴۳	۰/۲۸	میانگین

در تمامی شرایط (برای تمامی مقادیر کار کرد افتراقی و شرایط اثر آزمون و بدون اثر آزمون) توان روش نسبت درستنما برای کار کرد افتراقی یکنواخت بیشتر از روش رگرسیون لجستیک می‌باشد. در موقعیت بدون اثر روش نسبت درستنما با نرخ متوسط آشکارسازی ۶۲٪ نسبت به روش رگرسیون لجستیک با نرخ متوسط آشکارسازی ۴۹٪ قدرت بیشتری در رد فرض صفر مبنی بر عدم کار کرد افتراقی یکنواخت سؤال و شناسایی کار کرد افتراقی یکنواخت دارد. تحت شرایط اثر آزمون نیز روش نسبت درستنما برای کار کرد افتراقی یکنواخت با نرخ

متوسط آشکارسازی ۵۲٪ نسبت به روش رگرسیون لجستیک با نرخ متوسط آشکارسازی ۴۲٪ توان بیشتری در شناسایی کارکرد افتراقی دارد. میزان شناسایی کارکرد افتراقی یکنواخت در هر دو روش در شرایط اثر آزمون کمتر است. با بررسی پارامتر سوال‌ها مشخص می‌شود که روش نسبت درستنمایی معمولاً در شناسایی کارکرد افتراقی یکنواخت سوال‌هایی که پارامتر تشخیص بالا و دشواری متوسط دارند، توانمندتر است.

جدول ۳. توان کارکرد افتراقی یکنواخت شرایط اثر

کارکرد افتراقی یکنواخت					سؤال		
آزمون نسبت درستنمایی				رگرسیون لجستیک			
۱	۰/۸	۰/۶	۰/۴	۱	۰/۸	۰/۶	۰/۴
۰/۹۵	۱	۰/۸۵	۰/۶۵	۰/۹۵	۰/۸	۰/۷۵	۰/۶
۱	۰/۹۵	۰/۷۵	۰/۶۵	۱	۰/۹۵	۰/۵۵	۰/۴
۰/۳۵	۰/۳	۰/۱۵	۰/۱	۰/۵	۰/۲۵	۰/۱۵	۰/۱
۰/۰۵	۰/۰۵	۰/۱	۰/۰۵	۰	۰/۱۵	۰/۱	۰/۱۵
۰/۳	۰/۱	۰/۲۵	۰/۱۵	۰/۱۵	۰/۰۵	۰/۲	۰/۱۵
۰/۹۵	۰/۹۵	۰/۷	۰/۵	۰/۸۵	۰/۸۵	۰/۶	۰/۴
۰/۱	۰/۱	۰/۰۵	۰/۱	۰/۱	۰/۱	۰/۰۵	۰/۱
۰/۸	۰/۶	۰/۴	۰/۳۵	۰/۵	۰/۲	۰/۲	۰/۱
۰/۷۵	۰/۶	۰/۴	۰/۲۵	۰/۷	۰/۴۵	۰/۳	۰/۲۵
۰/۸۵	۰/۷	۰/۶	۰/۵	۰/۶۵	۰/۵۵	۰/۴۵	۰/۳۵
۰/۷۵	۰/۶	۰/۶	۰/۲۵	۰/۴۵	۰/۳	۰/۲	۰/۱۵
۰/۹۵	۰/۹۵	۰/۹۵	۰/۶۵	۱	۰/۹۵	۰/۸	۰/۳۵
۰/۸۵	۰/۸۸	۰/۹۵	۰/۲۵	۰/۸	۰/۵	۰/۴۵	۰/۱
۰/۲	۰/۲۵	۰/۱۵	۰/۱۵	۰/۰۵	۰/۲۵	۰/۱	۰/۲
۰/۹۵	۰/۹۵	۰/۶۵	۰/۶	۰/۸۵	۰/۲۵	۰/۴۵	۰/۵
۰/۶۴	۰/۶	۰/۵	۰/۳۵	۰/۵۷	۰/۴۸	۰/۳۶	۰/۲۶
میانگین							

توان آماری کارکرد افتراقی غیریکنواخت: توان آماری شناسایی کارکرد افتراقی غیریکنواخت (نسبت مثبت واقعی) دو رویکرد رگرسیون لجستیک و آزمون نسبت درستنمایی در ۲۰ تکرار برای سطوح مختلف کارکرد افتراقی و تحت شرایط اثر آزمون و بدون اثر در جداول شماره ۴ و ۵ ارائه شده است توان آماری روش رگرسیون لجستیک و

آزمون نسبت درستنمايی در شناسايي کارکرد افتراقی غيريکنواخت به شدت کارکرد افتراقی بستگی دارد. با افزایش شدت کارکرد افتراقی ميانگين نسبت مثبت واقعی هر دو روش رگرسيون لجستيک و آزمون نسبت درستنمايی تحت شرایط يکسانی و عدم يکسانی توزيع توانايی گروه مرجع و کانوني افزایش یافته است. در شرایط بدون اثر زمانی که شدت کارکرد افتراقی متوسط بود (۰/۶) رویکرد نسبت درستنمايی در مقایسه با روش رگرسيون لجستيک در شناسايي کارکرد افتراقی غيريکنواخت توانمندتر بود. هنگامی که شدت کارکرد افتراقی ۰/۴، ۰/۸ و ۱ بود، توان دو روش در شناسايي کارکرد افتراقی سؤال‌ها تقریباً يکسان است. نرخ متوسط آشکارسازی کارکرد افتراقی غيريکنواخت روش رگرسيون لجستيک و نسبت درستنمايی به ترتیب برابر ۲۹ و ۳۱ درصد است که مقادير کمی می‌باشد.

جدول ۴. توان کارکرد افتراقی غيريکنواخت شرایط بدون اثر

کارکرد افتراقی غيريکنواخت										سؤال	
آزمون نسبت درستنمايی					رگرسيون لجستيک						
۱	۰/۸	۰/۶	۰/۴	۱	۰/۸	۰/۶	۰/۴	۰/۴	۰/۴		
۰/۷۵	۰/۸۵	۰/۸۵	۰/۳۵	۰/۸۵	۰/۸	۰/۶۵	۰/۳	۰/۳	۴۱		
۰/۴۵	۰/۳	۰/۲	۰/۰۵	۰/۳	۰/۳۵	۰/۱	۰/۱	۰/۱	۴۲		
۰/۲۵	۰/۲	۰/۲۵	۰/۰۵	۰/۳۵	۰/۲۵	۰/۱	۰/۰۵	۰/۰۵	۴۳		
۰/۱۵	۰/۱	۰/۲۵	۰/۰۵	۰/۲۵	۰/۲۵	۰/۲۵	۰/۰۵	۰/۰۵	۴۴		
۰/۳	۰/۱	۰/۱	۰/۰۵	۰/۱۵	۰/۱	۰/۱	۰/۰۵	۰/۰۵	۴۵		
۰/۹۵	۰/۸۵	۰/۷	۰/۵	۱	۰/۹	۰/۸	۰/۵	۰/۵	۴۶		
۰/۱۵	۰/۱۵	۰/۱	۰/۰۵	۰/۰۵	۰/۱	۰/۱	۰/۰۵	۰/۰۵	۴۷		
۰/۱۵	۰/۲	۰/۲	۰/۲۵	۰/۱	۰/۰۵	۰/۱	۰/۱۵	۰/۱۵	۴۸		
۰/۳۵	۰/۵	۰/۲۵	۰/۱۵	۰/۴	۰/۵	۰/۲	۰/۱	۰/۱	۴۹		
۰/۱۵	۰/۱	۰/۱۵	۰/۱۵	۰/۱	۰/۰۵	۰/۱	۰/۰۵	۰/۰۵	۵۰		
۰/۲	۰/۱۵	۰/۱	۰/۰۵	۰/۱۵	۰/۱	۰/۰۵	۰/۲	۰/۲	۵۱		
۱	۰/۹	۰/۸	۰/۳۵	۱	۰/۹	۰/۷	۰/۲۵	۰/۲۵	۵۲		
۰/۱	۰/۰۵	۰/۱۵	۰/۰۵	۰/۰۵	۰/۰۵	۰/۱	۰/۱	۰/۱	۵۳		
۰/۰۵	۰/۱۵	۰/۱	۰/۱۵	۰/۲۵	۰/۱۵	۰/۱	۰/۱	۰/۱	۵۴		
۰/۸	۰/۸	۰/۸	۰/۳۵	۰/۸۵	۰/۷۵	۰/۶۵	۰/۲۵	۰/۲۵	۵۵		
۰/۳۹	۰/۳۶	۰/۳۳	۰/۱۷	۰/۳۹	۰/۳۵	۰/۲۷	۰/۱۵	۰/۱۵	ميانگين		

در شرایطی که توزیع توانایی دو گروه مرجع و کانونی متفاوت بود (اثر آزمون) به استثناء زمانی که شدت کارکرد افتراقی متوسط (۰/۶) بود در سایر شرایط توان دو روش در شناسایی کارکرد افتراقی سؤال‌ها تقریباً مشابه بود. نرخ متوسط آشکارسازی کارکرد افتراقی غیریکنواخت روش رگرسیون لجستیک و نسبت درستنمایی در حالت اثر به ترتیب برابر ۳۴ و ۳۱ درصد است. میزان شناسایی کارکرد افتراقی غیریکنواخت روش رگرسیون لجستیک در موقعیت اثر آزمون بیشتر از بدون اثر است. در حالی که میزان شناسایی کارکرد افتراقی غیریکنواخت روش نسبت درستنمایی هنگامی که شدت کارکرد افتراقی (کم) است در موقعیت اثر آزمون بیشتر از بدون اثر است و برای سایر مقادیر کارکرد افتراقی تفاوتی دو موقعیت مشاهده نمی‌شود. با بررسی پارامتر سؤال‌هایی که نسبت مثبت واقعی بالایی دارند مشخص می‌شود این سؤال‌ها پارامتر تشخیص بالا و دشواری متوسطی دارند.

جدول ۵. توان کارکرد افتراقی غیریکنواخت شرایط اثر

کارکرد افتراقی غیریکنواخت								
آزمون نسبت درستنمایی				رگرسیون لجستیک				سوال
۱	۰/۸	۰/۶	۰/۴	۱	۰/۸	۰/۶	۰/۴	
۱	۰/۹	۰/۸۵	۰/۷	۰/۸	۰/۷	۰/۶۵	۰/۳۵	۴۱
۰/۸	۰/۸	۰/۴	۰/۰۵	۰/۵	۰/۶۵	۰/۴۵	۰/۱۵	۴۲
۰/۱	۰/۰۵	۰/۰۵	۰/۱۵	۰/۵	۰/۲۵	۰/۲۵	۰/۱۵	۴۳
۰/۲۵	۰/۳	۰/۳	۰/۳	۰/۴۵	۰/۴۵	۰/۴۵	۰/۵۵	۴۴
۰/۲	۰/۱	۰/۱	۰/۱۵	۰/۴۵	۰/۳	۰/۳	۰/۱۵	۴۵
۰/۶	۰/۴۵	۰/۴۵	۰/۳۵	۰/۱۵	۰/۱۵	۰/۵	۰/۱۵	۴۶
۰/۳۵	۰/۳	۰/۴	۰/۳۵	۰/۸۵	۰/۹	۰/۴۵	۰/۶	۴۷
۰/۱	۰/۱۵	۰/۱۵	۰/۰۵	۰/۲	۰/۱	۰/۲	۰/۱	۴۸
۰/۰۵	۰/۱	۰/۱	۰/۱	۰/۲۵	۰/۲۵	۰/۱۵	۰/۰۵	۴۹
۰/۲۵	۰/۱	۰/۱۵	۰/۱۵	۰/۱۵	۰/۰۵	۰/۰۵	۰/۱	۵۰
۰/۲۵	۰/۱۵	۰/۰۵	۰/۱	۰/۱	۰/۰۵	۰/۰۵	۰/۰۵	۵۱
۰/۹۵	۰/۹۵	۰/۹	۰/۶	۰/۹۵	۰/۹۵	۰/۹۵	۰/۲	۵۲
۰/۱۵	۰/۱۵	۰/۱۵	۰/۰۵	۰/۰۵	۰/۱	۰/۰۵	۰/۰۵	۵۳
۰/۱۵	۰/۱۵	۰/۱	۰/۱۵	۰/۵	۰/۶	۰/۵۵	۰/۵	۵۴
۰/۴	۰/۴۵	۰/۲۵	۰/۱۵	۰/۲	۰/۲	۰/۲	۰/۰۵	۵۵

میانگین

۰/۳۷

۰/۳۴

۰/۲۹

۰/۲۳

۰/۴۱

۰/۳۸

۰/۳۵

۰/۲۰

با مقایسه نتایج جداول ۲ الی ۵ مشخص می‌شود که به طور کلی توانمندی هر دو روش نسبت درستنمایی سؤال – پاسخ و رگرسیون لجستیک در شناسایی کارکرد افتراقی یکنواخت بیشتر از کارکرد افتراقی غیریکنواخت می‌باشد.

خطای نوع اول (نسبت مثبت کاذب): خطای نوع اول (نسبت مثبت کاذب) دو رویکرد رگرسیون لجستیک و آزمون نسبت درستنمایی در ۲۰ تکرار برای موقعیت اثر آزمون و بدون اثر آزمون در جدول شماره ۶ ارائه شده است. همان‌گونه که در جدول مشخص شده است کنترل خطای نوع اول در رگرسیون لجستیک تحت تاثیر اثر آزمون (توزیع متفاوت گروه‌های کانونی و مرجع) است. در شرایط بدون اثر روش رگرسیون لجستیک کنترل مناسبی بر خطای نوع اول داشته است، میانگین نسبت مثبت کاذب روش رگرسیون لجستیک در شرایط بدون اثر برابر با ۰/۰۵ با بیشینه ۰/۲ برای سؤال شماره ۴۴ است. در شرایط بدون اثر میزان خطای نوع اول روش رگرسیون لجستیک برای ۹ سؤال از ۱۵ سؤال مساوی یا کمتر از ۰/۰۵ است. در شرایط اثر آزمون میانگین نسبت مثبت کاذب روش رگرسیون لجستیک برابر ۰/۱۶ با بیشینه ۰/۵ برای سؤال شماره ۴۴ است. سؤال ۴۴ با داشتن پارامتر دشواری ۲/۳۱۷ یکی از دشوارترین سؤال‌های بانک می‌باشد به این ترتیب مشخص می‌شود که تورم خطای نوع اول در روش رگرسیون لجستیک در شرایط اثر آزمون به پارامتر سؤال‌ها بستگی داشته است. در شرایط اثر میزان خطای نوع اول روش رگرسیون لجستیک ۵ سؤال مساوی یا کمتر از ۰/۰۵ است.

جدول ۶. نرخ خطای نوع اول (مثبت کاذب) در بانک سؤال اول

سؤال	نرخ خطای نوع اول			
	بدون اثر	با اثر	رگرسیون لجستیک	آزمون نسبت درستنمایی
۴۱	۰/۱	۰/۱	۰/۰۵	۰/۱
۴۲	۰/۰۵	۰/۱	۰	۰/۱
۴۳	۰/۱	۰/۱۵	۰/۱۵	۰/۱
۴۴	۰/۱۵	۰/۵	۰/۱۵	۰/۲
۴۵	۰/۱۵	۰/۱	۰	۰

۰/۰۵	۰/۲	۰	۰/۱	۴۶
۰	۰/۴	۰/۰۵	۰	۴۷
۰/۱۵	۰/۰۵	۰/۰۵	۰	۴۸
۰/۱۵	۰/۳	۰	۰	۴۹
۰	۰/۰۵	۰/۰۵	۰/۰۵	۵۰
۰	۰	۰	۰/۰۵	۵۱
۰/۰۵	۰/۰۵	۰	۰	۵۲
۰	۰/۰۵	۰/۰۵	۰	۵۳
۰/۱	۰/۲	۰/۰۵	۰/۱	۵۴
۰/۱۵	۰/۱۵	۰/۰۵	۰	۵۵
۰/۰۸	۰/۱۶	۰/۰۴	۰/۰۵	میانگین

در روش آزمون نسبت درستنماهی نیز کنترل خطای نوع اول تحت تاثیر اثر آزمون است. میزان خطای نوع اول در شرایط بدون اثر کمتر از موقعیت اثر آزمون است. آزمون نسبت درستنماهی نیز در شرایط بدون اثر کنترل خوبی بر میزان خطای نوع اول داشته است. میانگین نسبت مثبت کاذب روش نسبت درستنماهی در شرایط بدون اثر برابر با $0/04$ است. در شرایط اثر آزمون میانگین نسبت مثبت کاذب برابر $0/08$ است. همان‌گونه که مشاهده می‌شود روش نسبت درستنماهی خطای نوع اول را بهتر کنترل می‌کند در این روش تحت شرایط بدون اثر میزان خطای نوع اول 13 مساوی یا کمتر از $0/05$ است. حال آن که در شرایط اثر آزمون میزان خطای نوع اول 7 سؤال در دامنه مطلوب قرار گرفته است. در مقایسه دو روش مشخص می‌شود که در موقعیت بدون اثر دو روش رگرسیون و آزمون نسبت درستنماهی از نظر کنترل خطای نوع اول تقریباً مشابه هستند و در شرایط اثر نیز روش آزمون درستنماهی کنترل نسبتاً خوبی بر خطای نوع اول دارد.

بحث و نتیجه‌گیری

نظر به تمایل فراینده کاربرد کامپیوتروها در سنجش در قالب آزمون‌های کامپیوتروی و انطباقی کامپیوتروی، بررسی عادلانه بودن این نوع آزمون‌ها در چارچوب کارکرد افتراقی اهمیت بسزایی یافته است. تدوین و به کارگیری روش‌های شناسایی کارکرد افتراقی سؤال، پاسخی به ضرورت سنجش بدون سوگیری آزمودنی‌ها می‌باشد. ضرورت اجرا سؤال‌های فاقد کارکرد افتراقی در آزمون‌های انطباقی کامپیوتروی از ماهیت انطباقی آزمون نشات

می‌گیرد (کیم، ۲۰۱۰). در آزمون‌های انطباقی کامپیوتری سؤال‌های کمتری اجرا می‌شود در نتیجه هر سؤال نقش برجسته‌تری در برآورد توانایی آزمودنی ایفا می‌کند و از این رو عدم سوگیری سؤال‌ها به منظور برآورد بهینه توانایی آزمودنی‌ها ضرورتی اجتناب‌ناپذیر می‌باشد. علی‌الرغم اهمیت بررسی عملکرد متفاوت سؤال‌ها در آزمون‌های انطباقی کامپیوتری، انجام آن در مقایسه با آزمون‌های مداد – کاغذی دشوارتر می‌باشد و با چالش‌های نظری و عملی همراه است که می‌توان به دشواری تعیین متغیر تطبیق به منظور جور کردن آزمودنی‌های گروه‌های مورد مقایسه اشاره کرد (زویک، ۲۰۱۰).

در ادبیات تحقیق شاخص‌هایی مبتنی بر روش‌های به کار رفته در آزمون‌های مداد – کاغذی به منظور بررسی کارکرد افتراقی در چارچوب آزمون‌های انطباقی کامپیوتری توسعه یافته‌اند. این روش‌ها کارکرد افتراقی را در سؤال‌های پیش‌آزمون که معمولاً به صورت غیرانطباقی اجرا می‌شوند، مورد بررسی قرار می‌دهند و تلاشی برای ارزیابی کارکرد افتراقی در سؤال‌های عملیاتی در جریان آزمون انطباقی کامپیوتری صورت نگرفته است. کارکرد افتراقی و روش‌های شناسایی آن در بافت آزمون‌های انطباقی کامپیوتری از جمله حوزه‌های آزمون‌سازی محسوب می‌شوند که علی‌الرغم اهمیت با کمبود مطالعات نظامند موadge هستند.

در مطالعه حاضر کارکرد افتراقی سؤال‌های عملیاتی در اجرای انطباقی کامپیوتری براساس روش آزمون نسبت درستنمایی سؤال – پاسخ و رگرسیون مورد بررسی قرار گرفت. به طور کلی این دو روش از جمله روش‌های انتخابی مطالعه کارکرد افتراقی در آزمون‌های انطباقی کامپیوتری می‌باشند. اثربخشی هر یک از روش‌ها تابعی از شرایط آزمون و نوع سوگیری مندرج در سؤال‌ها است. از روش‌های شبیه‌سازی به منظور تولید داده‌ها و اجرای آزمون انطباقی کامپیوتری استفاده شد. نوع کارکرد افتراقی، مقدار کارکرد افتراقی و توزیع توانایی گروه مرجع و کانونی عواملی بودند که در این مطالعه دستکاری شدند. تکمیل ماتریس داده‌های ناقص حاصل از اجرای مجموعه متفاوتی از سؤال‌ها در اجرای انطباقی کامپیوتری به روش جایگذاری براساس احتمال پاسخ‌گویی درست به سؤال‌ها مبتنی بر برآورد توانایی آزمون انطباقی کامپیوتری برای هر آزمودنی صورت گرفت. مقایسه دو رویکرد رگرسیون لجستیک و آزمون نسبت درستنمایی براساس ملاک نسبت

مثبت کاذب (خطای نوع اول) و مثبت واقعی (توان) در سطح آلفا ۰/۰۵ در ۲۰ تکرار برای هر موقعیت انجام شد.

بررسی ادبیات پژوهشی بیانگر آن است که میزان توان و خطای نوع اول روش‌های شناسایی کارکرد افتراقی تحت تاثیر عوامل مختلفی از جمله طول آزمون، حجم نمونه، شدت کارکرد افتراقی، نسبت سوال‌های دارای کارکرد افتراقی و تفاوت میانگین گروه‌های مقایسه (اثر آزمون) می‌باشد (کاباساکال و همکاران، ۲۰۱۴؛ کیم، ۲۰۱۰). نتایج مطالعه حاضر نشان داد که میزان کنترل خطای نوع اول روش رگرسیون لجستیک و آزمون نسبت درستنمایی سؤال – پاسخ متاثر از اثر آزمون است. در شرایطی که توزیع توانایی دو گروه مرجع و کانونی یکسان است (بدون اثر) میزان خطای نوع اول هر دو روش مطلوب است. همان‌گونه که لی، چن و یو (۲۰۰۶) گزارش کردند در این مطالعه نیز میزان خطای نوع اول روش رگرسیون لجستیک تحت شرایط بدون اثر به گونه کافی کنترل شده است اما تحت شرایط اثر آزمون خطای نوع اول متورم شده است. کفایت کنترل خطای نوع اول روش رگرسیون لجستیک در شرایط یکسانی میانگین گروه‌های مقایسه مرجع و کانونی و تورم خطای نوع اول در شرایط تفاوت میانگین گروه‌ها در چارچوب آزمون‌های مداد – کاغذی نیز توسط نارایانان^۱ و سومانی ناتان (۱۹۹۶)، لی و استوت^۲ (۱۹۹۶) نیز گزارش شده است.

با افزایش شدت کارکرد افتراقی، توان هر دو روش رگرسیون لجستیک و آزمون نسبت درستنمایی در شناسایی کارکرد افتراقی یکنواخت و غیریکنواخت افزایش یافته است. این روند در هر دو موقعیت یکسانی و عدم یکسانی توزیع توانایی گروه‌های مقایسه مشاهده شده است. این نتایج همسو با یافته‌های مستند شده برای آزمون‌های انطباقی کامپیوتری و مداد – کاغذی در پیشینه تحقیق است (لی، چن و یو، ۲۰۰۶؛ نانداکومار و رووس، ۲۰۰۴).

توان هر دو روش نسبت درستنمایی سؤال – پاسخ و رگرسیون لجستیک در شناسایی کارکرد افتراقی یکنواخت در شرایط یکسانی میانگین گروه‌های مرجع و کانونی بیشتر است. هم‌چنین توان روش آزمون نسبت درستنمایی سؤال – پاسخ در تشخیص سوال‌های

1. Narayanan

2. Li & Stout

دارای کارکرد افتراقی یکنواخت در هر دو موقعیت یکسانی و عدم یکسانی توزیع توانایی گروه‌ها از روش رگرسیون لجستیک بیشتر می‌باشد. در تمامی شرایط نرخ متوسط آشکارسازی کارکرد افتراقی یکنواخت روش نسبت درستنایی بیشتر از روش رگرسیون لجستیک است. این یافته با نتایج مطالعه‌ای، چن و یو (۲۰۰۶) که یانگر مشابهت دو روش نسبت درستنایی سؤال – پاسخ و رگرسیون لجستیک در شناسایی کارکرد افتراقی یکنواخت است، هم‌خوانی ندارد. یافته‌های مطالعات کارکرد افتراقی در چارچوب آزمون‌های مداد – کاغذی نیز موید حساسیت بیشتر روش نسبت درستنایی در مقایسه با روش رگرسیون لجستیک می‌باشد (کیم، ۲۰۱۰).

هر دو روش رگرسیون لجستیک و آزمون نسبت درستنایی در شناسایی کارکرد افتراقی غیریکنواخت عملکرد ضعیفی داشتند. توان هر دو روش در شناسایی کارکرد افتراقی غیریکنواخت به شدت کارکرد افتراقی بستگی داشته است و با افزایش شدت کارکرد افتراقی افزایش یافته است. در شرایط بدون اثر زمانی که شدت کارکرد افتراقی ۰/۸، ۰/۴ و ۱ می‌باشد توان دو روش در شناسایی کارکرد افتراقی سؤال‌ها تقریباً یکسان بوده است، در حالی که تحت شرایط کارکرد افتراقی متوسط (۰/۶) رویکرد نسبت درستنایی توانمندتر است. هنگامی که توزیع توانایی دو گروه مرجع و کانونی متفاوت است، به استثناء زمانی که شدت کارکرد افتراقی متوسط (۰/۶) بود در سایر شرایط توان دو روش در شناسایی کارکرد افتراقی سؤال‌ها تقریباً مشابه است. این یافته با نتایج مطالعه‌ای، چن و یو (۲۰۰۶) هم‌خوانی ندارد.

با توجه به توان و میزان خطای نوع اول دوریکرد رگرسیون لجستیک و آزمون نسبت درستنایی در شناسایی کارکرد افتراقی یکنواخت در چارچوب آزمون انطباقی کامپیوتري می‌توان نتیجه گرفت که در شناسایی کارکرد افتراقی یکنواخت، روش آزمون نسبت درستنایی از کارایی بیشتری برخوردار است در حالی که در شناسایی کارکرد افتراقی غیریکنواخت تفاوتی میان دو روش مشاهده نشد و هر دو روش به طور مشابهی ضعیف عمل کرده‌اند.

به نظر می‌رسد نتایج مطالعه حاضر و عدم هم‌خوانی آن با یافته‌ها پژوهش لی، چن و یو (۲۰۰۶) به دلیل بررسی ارزیابی کارکرد افتراقی در سؤال‌های عملیاتی، فرایند جایگذاری داده‌های ناقص و هم‌چنین تغییرپذیری پارامتر تشخیص سؤال‌های مورد بررسی باشد. در

مطالعات قبلی کارکرد افتراقی در سؤال‌های پیش آزمون مورد بررسی قرار گرفته‌اند. سؤال‌های عملیاتی برخلاف سؤال‌های پیش آزمون که برای تمامی آزمودنی‌ها اجرا می‌شود، به طور انطباقی برای آزمودنی‌ها اجرا می‌شود. فرایند جایگذاری پاسخ‌ها در سؤال‌های عملیاتی براساس روش احتمال پاسخ‌گویی درست به سؤال‌ها و ماهیت انطباقی اجرای سؤال‌ها ممکن است توان و خطای روش‌های ارزیابی رگرسیون لجستیک و آزمون نسبت درستنمایی را تحت تاثیر قرار دهد و منجر به کاهش توان دو روش مذکور شده باشد. ترکیب مختلف پارامتر تشخیص سؤال‌ها مورد بررسی نیز ممکن است بر نتایج به دست آمده تاثیرگذاشته باشد.

به علاوه به نظر می‌رسد نتایج مطالعه حاضر به دلیل اثر تفاوت مرحله اجرای سؤال‌های دارای کارکرد افتراقی در اجرای انطباقی کامپیوتروی آزمون هم باشد. از آنجایی که در آزمون انطباقی کامپیوتروی سؤال‌ها متناسب با توانایی آزمودنی‌ها اجرا می‌شود، ترتیب اجرای سؤال‌ها برای آزمودنی‌های مختلف، متفاوت است در نتیجه سؤال دارای کارکرد افتراقی در جریان آزمون انطباقی کامپیوتروی در مراحل مختلفی برای آزمودنی‌ها اجرا می‌شود. به عنوان مثال سؤال ممکن است برای یک آزمودنی در ابتدا آزمون برای آزمودنی دیگر در اواسط و یا در مراحل پایانی اجرا شود. تفاوت در زمان اجرا سؤال ممکن است عملکرد افتراقی سؤال را تحت تاثیر قرار دهد. به نظر می‌رسد وقتی سؤال دارای کارکرد افتراقی در ابتدا آزمون اجرا می‌شود، عملکرد الگوریتم آزمون انطباقی کامپیوتروی قادر به اصلاح اثرات کارکرد افتراقی سؤال باشد، حال اگر این سؤال در انتهای آزمون اجرا شود اثرات سوگیرانه آن نمود بیشتری پیدا کند. تعداد سؤال‌های دارای کارکرد افتراقی و شدت کارکرد افتراقی نیز ممکن است عملکرد اصلاحی الگوریتم آزمون انطباقی کامپیوتروی را تحت تاثیر قرار دهد. بنابراین تاثیر مرحله اجرای سؤال‌های دارای کارکرد افتراقی و عوامل موثر بر آن می‌باشد در مطالعات بعدی مورد بررسی قرار گیرد.

به منظور بررسی دقیق‌تر پیشنهاد می‌شود عملکرد سایر روش‌های ارزیابی کارکرد افتراقی از جمله روش مانتل هنزل و CATSIB نیز در شناسایی کارکرد افتراقی سؤال‌های عملیاتی آزمون‌های انطباقی کامپیوتروی به شیوه جایگذاری پاسخ‌ها در سؤال‌های عملیاتی براساس روش احتمال پاسخ‌گویی درست به سؤال‌ها مورد بررسی قرار گیرد. به علاوه برای کسب تصویری روشن‌تر از روش بهینه ارزیابی کارکرد افتراقی در آزمون‌های

انطباقی کامپیوتری تاثیر سایر عوامل از جمله درصد تعداد سؤال‌های دارای کارکرد افتراقی و نسبت حجم نمونه در شناسایی کارکرد افتراقی در مطالعات دیگر مورد بررسی قرار گیرد. هم‌چنین روش نسبت درستنایی سؤال – پاسخ همانند هر روش مبتنی بر مدل ممکن است در صورت انتخاب مدل نامناسب عملکرد خوبی نداشته باشد. بنابراین لازم است در آینده میزان مقاوم بودن این روش در صورت کاربرد مدل نادرست نیز مورد ارزیابی قرار گیرد.

منابع

- کروکر، لیندا و آلجینا، جیمز. (۲۰۰۶). مباحث نو در روان‌سنجی. ترجمه و گردآوری ولی الله فرزاد، حسین زارع (۱۳۸۸). تهران: آیث
- گرامی پور، مسعود. (۱۳۹۰). مقایسه قدرت آزمون نسبت درستنایی مبتنی بر مدل سؤال – پاسخ با روش‌های تحلیل عاملی تاییدی و رگرسیون لجستیک در شناسایی کنش افتراقی سؤال به منظور اطمینان از عادلانه بودن سنجش آزمون‌های سرنوشت ساز.
- پایان نامه دکتری سنجش و اندازه‌گیری، دانشگاه علامه طباطبائی همبلتون، رونالد ک.، سوامیناتان، اچ. و راجرز، جین (۱۹۹۱). مبانی نظریه پرسش – پاسخ. ترجمه محمدرضا فلسفی نژاد، (۱۳۸۹). انتشارات دانشگاه علامه طباطبائی

- American Educational Research Association, American Psychological Association, National Council on Measurement in Education, Joint Committee on Standards for Educational, & Psychological Testing (US). (1999). Standards for educational and psychological testing. Amer Educational Research Assn.
- Chang, H. H., & Ying, Z. (1996). A global information approach to computerized adaptive testing. *Applied Psychological Measurement*, 20(3), 213-229.
- Choi, S. W., Podrabsky, T., & McKinney, N. (2011). Firestar-D: Computerized Adaptive Testing Simulation Program for Dichotomous Item Response Theory Models. *Applied Psychological Measurement*, 0146621611406107.
- Cohen, A. S., Kim, S. H., & Wollack, J. A. (1996). An investigation of the likelihood ratio test for detection of differential item functioning. *Applied Psychological Measurement*, 20(1), 15-26.
- Conoley, C. A. (2003). *Differential item functioning in the Peabody Picture Vocabulary Test–Third Edition: Partial correlation versus expert judgment* (Doctoral dissertation, Texas A&M University).

- Davidson, P. (2003). Why technology had had only a minimal impact on testing in education? Proceedings from the 2nd Education Technology Conference and Exhibition. Oman: Sultan Qaboos University
- Dorans, N. & Kulick, E. (1986). Demonstrating the utility of the standardization approach to assessing unexpected differential item performance on the Scholastic Aptitude Test. *Journal of Educational Measurement*, 23, 355-368
- Driana, Elin. (2007). *Gender item functioning on a ninth-grade mathematics proficiency test in Appalachian Ohio*. (Doctoral dissertation, Ohio University, Ohio.)
- Duncan, Cromwell, Susan. (2006). *Improving the prediction of differential item functioning: A comparison of the use of an effect size for logistic regression DIF and Mantel-Haenszel DIF methods*. (Doctoral dissertation, Texas A&M University).
- Guyer, R. D. (2008). *Effect of early misfit in computerized adaptive testing on the recovery of theta* (Doctoral dissertation). Available from ProQuest Dissertations and Theses database
- Han, K. T. (2007). WinGen: Windows software that generates IRT parameters and item responses. *Applied Psychological Measurement*, 31(5), 457-459.
- Han, K. T., & Hambleton, R. K. (2007). User's Manual: WinGen (Center for Educational Assessment Report No. 642). Amherst, MA: University of Massachusetts, School of Education.
- Harwell, M., Stone, C. A., Hsu, T. C., & Kirisci, L. (1996). Monte Carlo studies in item response theory. *Applied psychological measurement*, 20(2), 101-125.
- Herrera, A. N., & Gómez, J. (2008). Influence of equal or unequal comparison group sample sizes on the detection of differential item functioning using the Mantel-Haenszel and logistic regression techniques. *Quality & Quantity*, 42(6), 739-755.
- Holland, P. W., & Thayer, D. T. (1988). Differential item performance and the Mantel-Haenszel procedure. In H. Wainer, & H. I. Braun (Eds.), *Test validity* (pp. 129-145). Hillsdale, NJ: Lawrence Erlbaum.
- Kabasakal, K. A., Gök, B., Arsan, N., & Kelecioglu, H. (2014). Comparing Performances (Type I Error and Power) of IRT Likelihood Ratio SIBTEST and Mantel-Haenszel Methods in the Determination of Differential Item Functioning. *Kuram ve Uygulamada Eğitim Bilimleri*, 14(6), 2186.
- Kalender, Ilker. (2011). *Effect of Different Computerized Adaptive Testing Strategies on Recovery of Ability*. (Doctoral dissertation, Middle East Technical University).
- Kim, J. (2010). Controlling type 1 error rate in evaluating differential item functioning for four DIF methods: Use of three procedures for adjustment of multiple item testing.
- Kim, S. (2006). A comparative study of IRT fixed parameter calibration methods. *Journal of Educational Measurement*, 43(4), 355-381.

- Lei, P.-W., Chen, S.-Y., & Yu, L. (2006). Comparing methods of assessing differential item functioning in a computerized adaptive testing environment. *Journal of Educational Measurement*, 43, 254-264.
- Li, H.-H., & Stout, W. (1996). A new procedure for detection of crossing DIF. *Psychometrika*, 61(4), 647-677.
- Miller, T. R. (1992). Practical considerations for conducting studies of differential item functioning (DIF) in a CAT environment. Paper presented at the annual meeting of the American Educational Research Association, San Francisco, CA.
- Nandakumar, R., & Roussos, L. (2001). CATSIB: A Modified SIBTEST Procedure To Detect Differential Item Functioning in Computerized Adaptive Tests. Law School Admission Council Computerized Testing Report. LSAC Research Report Series.
- Nandakumar, R., & Roussos, L. (2004). Evaluation of the CATSIB DIF procedure in a pretest setting. *Journal of Educational and Behavioral Statistics*, 29, 177-199.
- Narayanan, P., & Swaminathan, H. (1994). Performance of the Mantel-Haenszel and simultaneous item bias procedures for detecting differential item functioning. *Applied Psychological Measurement*, 18(4), 315-328.
- Näsström, Gunilla. (2003). Differential item functioning for items in the Swedish National test in mathematics, course B. Paper presented at the Pre-ICME Conference in Växjö.
- Nydict, S. W., & Weiss, D. J. (2009). A hybrid simulation procedure for the development of CATs. In *Proceedings of the 2009 GMAC Conference on Computerized Adaptive Testing*. Retrieved from www.psych.umn.edu/psylabs/CATCentral
- Piromsombat, C. (2014). *Differential Item Functioning in Computerized Adaptive Testing: Can CAT Self-Adjust Enough?* (Doctoral dissertation, University of Minnesota).
- Raju, N. S. (1990). Determining the significance of estimated signed and unsigned areas between two item response functions. *Applied Psychological Measurement*, 14, 197-207.
- Shealy, R., & Stout, W. (1993). A model-based standardization approach that separates true bias/DIF from group ability differences and detects test bias/DTF as well as item bias/DIF. *Psychometrika*, 58, 159-194.
- Swaminathan, H., & Rogers, H. J. (1990). Detecting differential item functioning using logistic regression procedures. *Journal of Educational measurement*, 361-370.
- Thissen, D. (2001). IRTLRDIF v.2.0b: Software for the computation of the statistics involved in item response theory likelihood-ratio tests for differential item functioning. Chapel Hill: L.L. Thurstone Psychometric Laboratory, University of North Carolina.
- Thissen, D., Steinberg, L., & Wainer, H. (1993). Detection of differential item functioning using the parameters of item response models.

- Wang, T., & Vispoel, W. P. (1998). Properties of ability estimation methods in computerized adaptive testing. *Journal of Educational Measurement*, 35(2), 109-135.
- Zumbo, B. D. (1999). A handbook on the theory and methods of differential item functioning (DIF): Logistic regression modeling as a unitary framework for binary and Likert-like (ordinal) item scores. Ottawa, Canada: Directorate of Human Resources Research and Evaluation.
- Zwick, R. (2010). The investigation of differential item functioning in adaptive tests. In *Elements of adaptive testing* (pp. 331-352). Springer New York.
- Zwick, R., & Thayer, D. T. (2002). Application of an empirical Bayes enhancement of Mantel-Haenszel differential item functioning analysis to a computerized adaptive test. *Applied Psychological Measurement*, 26, 57-76.
- Zwick, R., Thayer, D. T., & Lewis, C. (1997). An investigation of the validity of an empirical Bayes approach to Mantel-Haenszel DIF analysis (ETS Research Report RR-97-21). Retrieved from Educational Testing Service website: <http://www.ets.org/Media/Research/pdf/RR-97-21.pdf>
- Zwick, R., Thayer, D. T., & Wingersky, M. (1994). A simulation study of methods for assessing differential item functioning in computerized adaptive tests. *Applied Psychological Measurement*, 18(2), 121-140.