

ارزیابی کارکرد افتراقی جنسیتی سؤال‌های آزمون ریاضی با استفاده از دو روش مانتل-هنزل و نظریه سؤال-پاسخ

نگار شریفی یگانه *

چکیده

آزمون‌ها می‌بایست برای تمام افراد جامعه از هر جنس، نژاد، سن، موقعیت اجتماعی و اقتصادی عادلانه باشد. بر این اساس ارزیابی وجود سوگیری و کارکرد افتراقی در سؤال‌های آزمون‌ها بسیار حائز اهمیت است. در مطالعه حاضر کارکرد افتراقی جنسیتی با استفاده از روش مانتل-هنزل و روش مبتنی بر نظریه سؤال-پاسخ مورد بررسی قرار گرفته است. در ابتدا مبانی نظری و روش‌های ارزیابی کارکرد افتراقی معرفی شده است و در ادامه به منظور ارائه نمونه عملی، پاسخ‌های یک نمونه تصافی طبقه‌ای ۴۰۰۰ نفری مشتمل بر ۲۲۰۰ آزمودنی مرد و ۱۸۰۰ زن به ۵۵ سؤال آزمون ریاضی کنکور سراسری مورد تحلیل قرار گرفته است. نتایج تحلیل بیانگر آن است که سؤال‌های آزمون دارای کارکرد افتراقی است. در بررسی کارکرد افتراقی با استفاده از روش مانتل-هنزل ۲۳ سؤال دارای شاخص مانتل-هنزل معنادار بودند. بر اساس رویکرد سؤال-پاسخ نیز ۹ سؤال دارای کارکرد افتراقی جنسیتی بودند که همگی به نفع دختران بودند. سؤال‌های دارای کارکرد افتراقی به نفع دختران بیشتر در حوزه محتوایی توابع و معادلات و سؤال‌های دارای کارکرد افتراقی به نفع پسران بیشتر در حوزه محتوایی مثلثات، هندسه و احتمال هستند.

واژگان کلیدی: نظریه سؤال-پاسخ، سوگیری سؤال، کارکرد افتراقی سؤال، روش مانتل-هنزل

تاریخ دریافت مقاله: ۹۱/۰۲/۰۵

تاریخ پذیرش مقاله: ۹۱/۰۶/۲۲

پژوهشگاه علوم انسانی و مطالعات فرهنگی
پرتال جامع علوم انسانی

* دانشجوی دکتری سنجش و اندازه‌گیری دانشگاه علامه طباطبایی و کارشناس سازمان سنجش آموزش کشور
(مسئول مکاتبات: nsh-yeganeh@yahoo.com)

مقدمه

روش‌های مختلف سنجش به ویژه آزمون‌ها نقش مهمی در سرنوشت افراد ایفا می‌کنند، از این رو، عادلانه^۱ و منصفانه بودن شیوه‌های سنجش اهمیت بسیاری دارد. در واقع روش‌های سنجش از جمله آزمون‌ها می‌بایست برای تمام افراد جامعه از هر جنس، نژاد، سن، موقعیت اجتماعی و اقتصادی منصفانه باشد و آزمودنی‌های متعلق به گروه‌های مختلف با توانایی‌های مشابه، امکان توفیق یکسانی در آزمون داشته باشند. هدف از اجرای آزمون‌ها، تفکیک و تمایز افراد است، اما آزمون‌ها نبایستی منجر به تبعیض میان آزمودنی‌ها شود (ناسستروم^۲، ۲۰۰۳).

در سال‌های اخیر، با توجه به افزایش کاربرد آزمون‌ها در موقعیت‌های مختلف، نگرانی فزاینده‌ای در باره امکان عملکرد متفاوت سؤال‌های آزمون و به عبارت دیگر سوگیری سؤال‌های^۳ آزمون در زیرگروه‌های آزمودنی (به طور مثال زنان و مردان) ایجاد شده است. به دلیل تأثیر سوگیری سؤال‌ها در انحراف نتایج آزمون و بالتبع تصمیم‌های اتخاذ شده بر مبنای این آزمون‌ها، دامنه این نگرانی حتی به سیاست‌گذاران آموزشی هم منتقل شده است. به اعتقاد منتقدان آزمون‌ها، سوگیری سؤال‌های آزمون ممکن است منجر به از دست رفتن فرصت‌های آموزشی و شغلی شود (کونولی^۴، ۲۰۰۳).

سوگیری زمانی مطرح می‌شود که احتمال پاسخ‌گویی درست گروهی از آزمودنی‌ها به دلیل مشخصات سؤال‌ها و شرایط آزمون، از گروه دیگری از آزمودنی‌ها کم‌تر یا بیشتر باشد (زومبو^۵، ۱۹۹۹). در واقع سؤال زمانی سوگیری دارد که ویژگی‌های غیرمرتبط با سازه مورد سنجش را بسنجد. سوگیری سؤال ممکن است به علت ابهام در متن، گزینه‌ها و راهنمای آزمون باشد (همبلتون^۶ و راجرز^۷، ۱۹۹۵). کارکرد افتراقی سؤال‌ها^۸ (DIF) روش آماری تعیین وجود سوگیری در سؤال‌های آزمون است. این اصطلاح توسط هالند^۹ و واینر^{۱۰} (۱۹۹۳) ارائه شده است (اینگ^{۱۱} و هوی^{۱۲}، ۲۰۰۵).

1. Fair
2. Näsström
3. Item Bias
4. Conoley
5. Zumbo
6. Hambleton
7. Rodgers
8. Differential Item Function
9. Holland
10. Wainer
11. Eng
12. Hoe

کارکرد افتراقی سؤال، زمانی مطرح می‌شود که آزمودنی‌ها با توانایی یکسان، متعلق به گروه‌های مختلف، احتمال پاسخ‌گویی درست اما متفاوتی داشته باشند. به عبارت دیگر سؤال در صورتی دارای کارکرد افتراقی است که در زیر گروه‌های مختلف جامعه عملکرد متفاوت داشته باشد. وجود کارکرد افتراقی بیانگر آن است که عوامل مربوط به عضویت در گروه، احتمال پاسخ‌گویی درست را تحت تأثیر قرار می‌دهد. در صورتی که سؤال، کارکرد افتراقی نداشته باشد، فاقد سوگیری نیز هست، ولی با وجود کارکرد افتراقی، وجود سوگیری سؤال نیز، باید به روش‌های مختلفی چون تحلیل محتوا و ارزشیابی تجربی مورد بررسی قرار گیرد (دانکن، ۲۰۰۶).

تاریخچه و رویکردهای ارزیابی کارکرد افتراقی سؤال

چالش اساسی در خصوص آزمون عادلانه، نگرانی در خصوص سوگیری سؤال‌های آزمون و ضرورت بررسی آن به عنوان بخشی از فرایند آزمون‌سازی، در اواخر دهه ۱۹۶۰ و اوایل دهه ۱۹۷۰ در راستای انتظارات مبنی بر یکسانی نتایج آزمون‌ها و هم‌زمان با آغاز عصر حقوق مدنی مطرح شد (دانکن، ۲۰۰۶). در چهل سال اخیر، نگرانی‌های مربوط به عادلانه بودن آزمون‌ها، منجر به ایجاد روش‌های مختلفی برای ارزیابی کارکرد افتراقی سؤال‌ها شده است. به عنوان مثال هیلز^۱ بیش از چهل روش را مشخص کرده است. روش‌های مختلف ارزیابی کارکرد افتراقی از لحاظ مفروضه‌ها، شرایط کاربرد و نتایج حاصل با یکدیگر تفاوت دارند و همین مسئله زمینه‌ای را برای پژوهشگران، به منظور مقایسه روش‌های مختلف ارزیابی کارکرد افتراقی سؤال‌ها، با هدف دستیابی به مناسب‌ترین روش برای هر موقعیت ایجاد کرده است (زومبو، ۱۹۹۹). پژوهشگران روش‌های ارزیابی کارکرد افتراقی را بر اساس الگوها و آماره‌های به کار رفته طبقه‌بندی کرده‌اند که این موجب درک بهتری از روش‌های مختلف می‌شود. در طبقه‌بندی کلی می‌توان روش‌های تحلیل کارکرد افتراقی را در چارچوب نظریه کلاسیک و یا نظریه سؤال - پاسخ^۲ طبقه‌بندی کرد. در نظریه کلاسیک، تفاوت مشخصه‌های کلاسیک سؤال، در بین گروه‌های مختلف آزمودنی با نمره‌های مشابه مورد مقایسه قرار می‌گیرد. سادگی اجرا و کاربرد نمونه‌های کوچک از مزیت‌های روش‌های مبتنی بر رویکرد کلاسیک است (امبرستون^۳ و رایس^۴، ترجمه شریفی و همکاران، ۱۳۸۸). برخی از متداول‌ترین روش‌های بررسی کارکرد افتراقی عبارتند از:

1. Hills
2. Item response theory
3. Embreston
4. Reise

۱- روش‌های دشواری سؤال تبدیل شده^۱: این روش توسط انگاف^۲ و فورد^۳ پیشنهاد شده است و به آن طرح دلتا نیز گفته می‌شود. در این روش پارامتر دشواری سؤال در هر یک از زیر گروه‌های جامعه به طور جداگانه محاسبه می‌شود، سپس مقادیر دشواری سؤال به مقیاس دلتا (میانگین ۱۳ و انحراف استاندارد ۴) تبدیل می‌شوند (باقی^۴ و فرارا^۵، ۱۹۸۹). همبستگی مقادیر دلتا مربوط به دو گروه محاسبه می‌شود، به علاوه برای تمام سؤال‌ها نمودار مقادیر دلتا مربوط به هر زوج رسم می‌شود و با داده‌ها خط مستقیمی برازش داده می‌شود. سؤال‌هایی که مقادیر دلتا در آنها با فاصله زیادی از این خط قرار گرفته‌اند دارای سوگیری هستند (آنیل، ۱۹۹۱). این روش به علت سادگی مورد توجه بسیار واقع شده است، اما نسبت به عدم یکسانی قدرت تشخیص سؤال‌ها حساس است.

۲- روش مانتل-هنزل^۶: این روش توسط مانتل و هنزل (۱۹۵۹) به عنوان روشی برای مطالعه گروه‌های جور شده معرفی شده است. هالند و تایر (۱۹۸۸) این روش را در بررسی کارکرد افتراقی سؤال مورد استفاده قرار دادند (دورانس^۷ و هالند، ۱۹۹۲). روش مانتل-هنزل، روشی ناپارامتریک برای ارزیابی کارکرد افتراقی است. سؤال‌ها به صورت دوازده‌گانه نمره‌دهی می‌شوند و برای تعیین آزمودنی‌ها با سطوح توانایی یکسان در دو گروه مرجع^۸ و کانونی^۹ نمره کل مورد استفاده قرار می‌گیرد (گروه مرجع گروهی از آزمودنی‌ها است که عملکرد آنها به عنوان نقطه مرجع محسوب می‌شود. گروه کانونی گروهی از آزمودنی‌ها است که انتظار می‌رود سؤال‌های آزمون به نفع آنان نباشد). این روش به دلیل ارائه آزمون آماری، برآورد اندازه اثر و کارآمدی آن در نمونه‌های با حجم کم، به عنوان یکی از رایج‌ترین روش‌های ارزیابی کارکرد افتراقی سؤال‌ها مطرح است. البته این روش فقط در تشخیص کارکرد افتراقی یکنواخت^{۱۰} مناسب است و همین موضوع یکی از محدودیت‌های روش مانتل هنزل است. کارکرد افتراقی به دو صورت یکنواخت و غیریکنواخت^{۱۱} است. در کارکرد

1. Transformed Item Difficulty
2. Angoff
3. Ford
4. Baghi
5. ferrara
6. Mantel- Haenszel
7. Dorans
8. Reference
9. Focal
10. Uniform
11. Non -Uniform

افتراقی یکنواخت تفاوت احتمال پاسخ‌گویی درست گروه کانونی و گروه مرجع در تمام سطوح توانایی یکسان است یا به عبارت دیگر در کارکرد افتراقی یکنواخت تعامل بین سطح توانایی و عضویت در گروه وجود ندارد (دریانا^۱، ۲۰۰۷). در کارکرد افتراقی یکنواخت سؤال‌ها از لحاظ سطوح دشواری در دو گروه متفاوت هستند، اما از لحاظ قدرت تشخیص تفاوتی ندارند (شولتز^۲ و ویتنی^۳، ۲۰۰۵). در کارکرد افتراقی افتراقی غیریکنواخت احتمال پاسخ‌گویی درست گروه کانونی و گروه مرجع در تمام سطوح توانایی یکسان نیست و در واقع بین سطح توانایی و عضویت در گروه تعامل وجود دارد (دریانا، ۲۰۰۷) و سؤال‌ها از لحاظ سطح دشواری و قدرت تشخیص در دو گروه متفاوت قرار می‌گیرند (شولتز و ویتنی، ۲۰۰۵).

۳- روش رگرسیون لجستیک: این روش توسط سوامی‌ناتان^۴ و راجرز^۵ معرفی شده است. در این روش پاسخ سؤال‌ها، به عنوان متغیری وابسته در نظر گرفته می‌شود. رگرسیون لجستیک مبتنی بر مدل‌سازی آماری احتمال پاسخ صحیح به سؤال، بر اساس عضویت در گروه و ملاک است که در آن ملاک معمولاً نمره کل آزمون است. وجود کارکرد افتراقی با بررسی بهبود ایجاد شده در برازش مدل رگرسیون، پس از اضافه کردن عضویت در گروه و تعامل بین نمره آزمون و عضویت در گروه در مدل رگرسیون تعیین می‌شود (روسسو^۶، برتراند^۷ و بیتو^۸، ۲۰۰۴). روش رگرسیون لجستیک لجستیک مانند روش مانتل هنزل، آزمون معناداری آماری و اندازه اثر ارائه می‌کند. از جمله مزایای این روش می‌توان به توانایی بررسی کارکرد افتراقی یکنواخت و غیریکنواخت اشاره کرد.

بررسی کارکرد افتراقی بر اساس نظریه سؤال - پاسخ به ابتدای دهه شصت میلادی و کاربرد مدل راش بر می‌گردد. لرد^۹ و نوایک^{۱۰} ظرفیت بالای رویکرد سؤال - پاسخ را در ارزیابی کارکرد افتراقی سؤال مورد تاکید قرار دادند (دانکن، ۲۰۰۶). از لحاظ نظری رویکرد سؤال - پاسخ روش مناسبی برای ارزیابی کارکرد افتراقی سؤال‌ها

1. Driana
2. Shultz
3. Whitney
4. Swaminathan
5. Rogers
6. Rousseau
7. Bertrand
8. Boiteau
9. Lord
10. Novick

است. ویژگی نامتغیر بودن پارامترها در این رویکرد، چارچوبی نظری برای تعریف و تعیین کارکرد افتراقی سؤال‌های آزمون فراهم می‌سازد. اگر سؤال، در زیرگروه‌های جامعه یکسان عمل کند، در این حالت احتمال پاسخ‌گویی درست آزمودنی‌ها با سطح توانایی یکسان می‌بایست مشابه باشد. به این ترتیب بر اساس نظریه سؤال - پاسخ، سؤال زمانی دارای کارکرد افتراقی است که منحنی ویژگی سؤال یا به عبارت دیگر احتمال شرطی پاسخ درست به سؤال در سطوح توانایی یکسان در گروه مرجع و کانونی متفاوت باشند (امبرتسون و رایس، ترجمه شریفی و همکاران، ۱۳۸۸). رویکرد سؤال - پاسخ قادر به تشخیص کارکرد افتراقی یکنواخت و غیریکنواخت است. یکی از محدودیت‌های روش‌های مبتنی بر رویکرد سؤال - پاسخ ضرورت وجود گروه‌های نمونه بزرگ است. زیکی^۱ (۱۹۹۳)، روسو و ستوات^۲ کاربرد گروه نمونه با حجم بالای بالای صد نفر را برای گروه کانونی و بین دویست تا هزار نفر را برای گروه مرجع در ارزیابی کارکرد افتراقی توصیه کردند. از سویی کاربرد نمونه‌های بزرگ منجر به افزایش موارد مثبت کاذب می‌شود (دانکن، ۲۰۰۶).

در چارچوب رویکرد سؤال - پاسخ چندین روش برای ارزیابی کارکرد افتراقی مطرح شده است که این روش‌ها در دو گروه کلی مساحت میان منحنی ویژگی سؤال^۳ گروه مرجع و کانونی و روش آزمون‌های آماری یکسانی پارامترهای سؤال طبقه‌بندی می‌شوند (شولتز و ویتنی، ۲۰۰۵). روش آزمون‌های آماری، شامل آزمون چند متغیره و آزمون t بر روی مقادیر پارامتر دشواری می‌شود. در رویکرد مساحت میان دو منحنی ویژگی سؤال، منحنی ویژگی و نه پارامتر سؤال‌ها مورد توجه قرار می‌گیرد. هر چند که منحنی ویژگی سؤال‌ها بر اساس پارامترهای سؤال‌ها ترسیم می‌شوند، اما روش بررسی مساحت بین دو منحنی ویژگی سؤال، مستلزم نگاه دقیق‌تری به تفاوت بین پارامتر سؤال‌ها است. در این روش منطقه بین دو منحنی ویژگی سؤال مورد توجه قرار می‌گیرد. پس از برآورد پارامتر سؤال‌ها و قرار دادن آن‌ها روی مقیاس مشترک، منحنی سؤال در دو گروه ترسیم می‌شود. در صورتی که فاصله بین دو منحنی ویژگی سؤال صفر باشد و در واقع دو منحنی بر هم منطبق باشند، کارکرد افتراقی وجود ندارد. برعکس هنگامی که مساحت بین دو منحنی سؤال صفر نیست کارکرد افتراقی با درجاتی وجود دارد. هر چه

1. Zieky
2. Stout
3. Item Characteristic Curve

این مساحت بیشتر باشد، کارکرد افتراقی نیز بیشتر است. همبلتون و راجو^۱ فرمول‌هایی برای محاسبه مساحت ارائه کرده‌اند (همبلتون، سوامیناتان و راجرز، ترجمه فلسفی‌نژاد، ۱۳۸۹). راجو نارسایی رویکرد محاسبه مساحت بین دو منحنی ویژگی سؤال را درک کرده و روش‌های مبتنی بر چارچوب کارکرد افتراقی سؤال‌ها و آزمون^۲ (DFIT) را برای ارزیابی کارکرد افتراقی ارائه کرد. شاخص DFIT ویژگی‌هایی دارد که آن را توانمند و منعطف می‌سازد. این ویژگی‌ها عبارتند از: ۱- امکان کاربرد این شاخص برای سؤال‌های دو ارزشی و چندارزشی ۲- امکان کاربرد این شاخص در مدل‌های تک بعدی و چند بعدی ۳- توانایی ارزیابی کارکرد افتراقی سؤال و آزمون (اشیما و موریس، ۲۰۰۸).

مقایسه روش‌های مختلف تعیین کارکرد افتراقی سؤال‌ها

به منظور مقایسه روش‌های مختلف ارزیابی کارکرد افتراقی، مطالعات متعددی با استفاده از داده‌های حقیقی و یا شبیه‌سازی شده انجام شده است که در برخی موارد نتایج متناقضی به همراه داشته‌اند. دو روش رگرسیون لجستیک و روش مانتل هنزل به عنوان رایج‌ترین روش‌های ارزیابی کارکرد افتراقی در بسیاری از مطالعات مورد مقایسه قرار گرفته‌اند. به عنوان مثال هیدالگو^۳ و لویز^۴ در مطالعه‌ای به بررسی این دو رویکرد در ارزیابی کارکرد افتراقی سؤال‌های آزمون پرداختند. نتایج مطالعه آنان بیانگر آن بود که روش رگرسیون لجستیک سؤال‌های بیشتری را با کارکرد افتراقی مشخص می‌سازد. سوامیناتان و راجرز (۱۹۹۰) نیز با استفاده از مطالعات شبیه‌سازی شده به نتیجه مشابه دست یافتند (دانکن، ۲۰۰۶). گیرل^۵ و همکاران (۱۹۹۹) در مطالعه‌ای کارکرد افتراقی جنسیتی سؤال‌های آزمون ریاضی و علوم را با استفاده از دو روش مانتل هنزل و رگرسیون لجستیک مورد بررسی قرار دادند. در حالی که نتایج این دو روش در بررسی کارکرد افتراقی یکنواخت در آزمون ریاضی مشابه بود، نتایج آزمون علوم ثبات کم‌تری داشت. روش مانتل هنزل در مقایسه با روش رگرسیون لجستیک تعداد کم‌تری از سؤال‌ها را دارای کارکرد افتراقی تشخیص داد و بنابر این روشی محافظه‌کارانه‌تر می‌باشد. با توجه به مطالعات ارائه شده می‌توان نتیجه گرفت

1. Raju
2. Differential Functioning of Items and Tests
3. Hidalgo
4. Lopez
5. Gierl

که تعداد سؤال‌های دارای کارکرد افتراقی، به روش به‌کار رفته وابسته است و از این رو سازندگان آزمون‌ها و سیاست‌گذاران آموزشی می‌بایست تفاوت‌های روش‌های مختلف ارزیابی کارکرد افتراقی سؤال‌ها را همواره در نظر بگیرند (احمدی، ۱۳۸۷).

باقی و فرارا (۱۹۸۹) در مطالعه‌ای سه روش نمودار دلتا، روش مانتل هنزل و رویکرد مبتنی بر نظریه سؤال - پاسخ را با استفاده از نمونه‌هایی با حجم متفاوت مورد مقایسه قرار دادند. روش‌ها از لحاظ ثبات (مطابقت نتایج هر روش در نمونه‌ها با حجم‌های متفاوت)، توافق (مطابقت نتایج روش‌های مختلف ارزیابی کارکرد افتراقی با نتایج روش مبتنی بر رویکرد سؤال - پاسخ) و عملی بودن (نرم‌افزارهای مورد نیاز، حجم گروه نمونه مورد نیاز) با یکدیگر مقایسه شدند. ثبات روش نمودار دلتا و رویکرد سؤال - پاسخ کم تا متوسط بود. در نمونه‌های بزرگ ثبات روش سؤال - پاسخ، نمودار دلتا و روش مانتل هنزل در حد متوسط بود. توافق روش‌ها با محاسبه همبستگی رتبه‌ای اسپیرمن میان شاخص‌های کارکرد افتراقی روش‌ها و مشخص کردن تعداد سؤال‌ها دارای کارکرد افتراقی و فاقد کارکرد افتراقی در هر یک از روش‌ها تعیین شد. بر اساس ضرایب همبستگی رتبه‌ای اسپیرمن، توافق میان روش راش و نمودار دلتا در تمام نمونه‌ها در حد بسیار بالا بود، توافق روش راش با مدل سه پارامتری در حد متوسط بود و توافق سایر روش‌ها با روش سه پارامتری کم بود. هریس^۱ و کولن^۲ (۱۹۸۶)، هریس و هوور^۳ (۱۹۸۶) و اسکاگز^۴ و لیسیتز^۵ (۱۹۸۶) اظهار داشتند که میزان توافق بین روش‌های مختلف ارزیابی کارکرد افتراقی خیلی زیاد نیست (باقی و فرارا، ۱۹۸۹).

عبدالعزیز^۶ (۲۰۱۰) کارکرد افتراقی جنسیتی سؤال‌ها را در آزمون ریاضی با استفاده از روش مانتل هنزل، تفاوت پارامتر دشواری سؤال‌های آزمون و روش پارامتر دشواری تبدیل شده (TID) را بررسی و میزان توافق این سه روش را تعیین کرد. نتایج تحقیق وی بیانگر آن بود که میزان توافق بین این سه روش پایین است. بیشترین توافق، بین روش مانتل هنزل و روش پارامتر دشواری تبدیل شده (TID) به دست آمد که این مسئله ممکن است به این دلیل باشد که هر دو این روش‌ها در چارچوب

-
1. Harris
 2. Kolen
 3. Hoover
 4. Skaggs
 5. Lissitz
 6. Abedalaziz

نظریه کلاسیک هستند. کم‌ترین توافق میان روش پارامتر دشواری تبدیل شده (TID) و تفاوت پارامتر دشواری سئوال‌ها است.

دوران^۱ و پارک^۲ (۲۰۰۶) کارکرد افتراقی را در سئوال‌های آزمون زبان با استفاده از سه روش SIBTEST، مانتل هنزل و رویکرد سئوال - پاسخ بررسی کردند. نتایج نشان داد که دو روش SIBTEST و مانتل هنزل تعداد زیادی سئوال دارای کارکرد افتراقی را مشخص ساختند، در حالی‌که روش سئوال - پاسخ تعداد سئوال‌های کم‌تری را به عنوان سئوال‌های دارای کارکرد افتراقی مشخص کرده است. بر اساس این نتایج آنان مطالعات بیشتر و دقیق‌تری را در باره روش‌های مختلف ارزیابی کارکرد افتراقی توصیه کردند.

یکی از مسایل مهم در مطالعات کارکرد افتراقی تعیین عوامل به وجود آورنده آن است. در باره سئوال‌های دو ارزشی مطالعاتی صورت گرفته است که از آن جمله می‌توان به مطالعه اونیل و مک‌پیک^۳ اشاره کرد که تأثیر عواملی چون محتوا و شکل سئوال را در آزمون استعداد تحصیلی^۴ (SAT)، آزمون‌های ورودی تحصیلات تکمیلی^۵ (GRE) و سایر آزمون‌های پذیرش مورد بررسی قرار دادند. نتایج مطالعه آنان بیانگر آن بود که در آزمون ریاضی آزمون‌های دختر در مقایسه با آزمون‌های پسر در سئوال‌های جبر عملکرد بهتری دارند، در حالی‌که آزمون‌های پسر در حل مسایل، عملکرد بهتری دارند (وانگ و لان، ۱۹۹۴). پژوهشگران هم‌چنین نشان دادند که کارکرد افتراقی جنسیتی به میزان شباهت سئوال با سئوال‌های کتاب‌های درسی نیز بستگی دارد. آزمون‌های دختر در مقایسه با پسران در سئوال‌هایی که مشابه سئوال‌های کتاب‌های درسی هستند، عملکرد بهتری دارند (ویلینگهام و کول، ۱۹۹۷). در مورد عوامل به وجود آورنده کارکرد افتراقی در سئوال‌های چند ارزشی مطالعات اندکی انجام شده است (وانگ و لان، ۱۹۹۴). دودین^۶ و الدرابی^۷ (۲۰۰۳) نیز کارکرد افتراقی سئوال‌های آزمون ریاضی را با استفاده از روش مانتل هنزل مورد بررسی قرار دادند. آزمون متشکل از پنجاه سئوال چندگزینه‌ای بود که با نمونه چهارصد نفری از دانش‌آموزان اجرا شد (۲۰۰ دانش‌آموزان دختر و ۲۰۰ دانش‌آموزان پسر). از این پنجاه

1. Durand
2. Park
3. McPeck
4. Scholastic Aptitude Test
5. Graduate Record Examinations (GRE)
6. Doudin
7. Al-Darabee

سؤال ۱۷ سؤال (۳۴ درصد) کارکرد افتراقی داشتند که هشت سؤال به نفع دانش‌آموزان پسر و نه سؤال به نفع دانش‌آموزان دختر بود. نتایج مطالعه آنان بیانگر برتری دانش‌آموزان پسر در سؤال‌های حساب و حل مسئله و همچنین برتری دانش‌آموزان دختر در سؤال‌ها جبر و هندسه بود. بربراولو^۱ (۱۹۹۶) وجود کارکرد افتراقی جنسیتی در سؤال‌های ریاضی آزمون سراسری ورود به دانشگاه را با استفاده از مدل دو پارامتری مورد بررسی قرار داد. نتایج تحلیل مبین آن بود که از ۳۲ سؤال چندگزینه‌ای آزمون، ۵۳ درصد سؤال‌ها به نفع مردان و ۴۷ درصد به نفع زنان بوده است (دریانا، ۲۰۰۷).

مندز^۲، بارنت^۳ و اریکان^۴ (۲۰۰۶) منابع کارکرد افتراقی سؤال‌ها با توجه به عامل عامل جنسیت را مورد بررسی قرار دادند. بر اساس فرضیه آنان، محتوا، سطوح شناختی و متن سؤال، منابع احتمالی کارکرد افتراقی سؤال محسوب می‌شوند. آنان دریافتند که در حوزه محتوایی، مجموعه سؤال‌های حل مسئله و مجموعه سؤال‌های مربوط به حوزه‌های لگاریتم و توان به نفع دانش‌آموزان پسر و مجموعه سؤال‌های چندجمله‌ای‌ها و روابط درجه دوم به نفع دانش‌آموزان دختر است. بر اساس سطوح شناختی سؤال‌ها، دانش‌آموزان پسر در مقایسه با دختران در سؤال‌های سطوح بالای شناختی عملکرد بهتری داشتند. در سؤال‌های مربوط به سطوح پایین شناختی تفاوتی مشاهده نشد. مطالعه آنان در باره متن سؤال‌ها نشان داد سؤال‌هایی که شامل شکل و نمودار هستند به نفع پسران است. نتایج این پژوهش موافق با نتایج سایر مطالعات انجام شده در این زمینه است (دریانا، ۲۰۰۷).

اهداف مطالعه

اتخاذ تصمیمات درست و منطقی برپایه نتایج آزمون‌ها مستلزم کسب اطمینان از عادلانه بودن آزمون‌ها است. بنابراین ارزیابی کارکرد افتراقی سؤال‌های آزمون‌ها اهمیت فراوانی دارد و می‌بایست به عنوان بخشی از فرایند تحلیل آزمون در نظر

-
1. Berberoglu
 2. Mendes
 3. Barnett
 4. Ercikan

گرفته شود. با توجه به اهمیت آزمون‌های ورودی دانشگاه‌ها و نقش آنها در سرنوشت تحصیلی و شغلی داوطلبان، لازم است سؤال‌های آزمون‌های ورودی از نظر وجود کارکرد مورد بررسی دقیق و موشکافانه قرار گیرد. در این راستا مطالعه حاضر به ارزیابی کارکرد افتراقی جنسیتی در سؤال‌های آزمون ریاضی گروه آزمایشی ریاضی و فنی با استفاده از روش مانتل - هنزل و روش مبتنی بر نظریه سؤال- پاسخ پرداخته است. به علاوه میزان توافق بین این دو روش در تعیین کارکرد افتراقی نیز مورد بررسی قرار گرفته است. در مطالعه حاضر آزمودنی‌های دختر به عنوان گروه کانونی و آزمودنی‌های پسر به عنوان گروه مرجع در نظر گرفته شده‌اند.

سؤال تحقیق

آیا عملکرد آزمودنی‌های دختر و پسر با سطوح توانایی مشابه در آزمون ریاضی کنکور سراسری گروه آزمایشی ریاضی و فنی متفاوت است؟

داده‌های تحلیل

تفاوت عملکرد در آزمون‌های ریاضی با توجه به ویژگی‌ها نظیر وضعیت اقتصادی - اجتماعی، جنسیت و نژاد از جمله مسائلی است که مورد توجه و پژوهش فراوان قرار گرفته است. در مورد آزمون‌های ریاضی، کارکرد افتراقی دغدغه‌ای اساسی است. در مطالعه حاضر کارکرد افتراقی جنسیتی سؤال‌های آزمون ریاضی مورد بررسی قرار گرفته است. در تحلیل حاضر از پاسخ‌های آزمودنی‌های شرکت کننده در کنکور سراسری گروه آزمایشی علوم ریاضی و فنی سال ۱۳۸۲ استفاده شده است. ۳۹۹۱۹۹ آزمودنی در این آزمون شرکت کرده‌اند که ۷۸۷۵ آزمودنی به تمام سؤال‌های آزمون ریاضی پاسخ غلط داده‌اند و در نتیجه از روند تحلیل کنار گذاشته شدند. به این ترتیب تحلیل براساس پاسخ‌های ۳۹۱۳۲۴ آزمودنی انجام شده است. ۲۱۳۶۲۵ آزمودنی پسر و ۱۷۷۶۹۹ نفر دختر بودند. با استفاده از روش نمونه‌گیری تصادفی طبقه‌ای گروه نمونه چهارهزار نفری از آزمودنی‌ها انتخاب شد که شامل ۲۲۰۰ آزمودنی مرد و ۱۸۰۰ آزمودنی بود. آزمون شامل ۵۵ سؤال چهارگزینه‌ای بود که برای پاسخ‌های نادرست نمره منفی لحاظ شده است.

روش تحلیل

در مطالعه حاضر بررسی کارکرد افتراقی سؤال‌ها براساس دو روش سؤال - پاسخ و مانند - هزل صورت گرفته است. اولین گام در تحلیل مبتنی بر نظریه سؤال - پاسخ بررسی برقراری مفروضه‌های مدل است، چون کاربرد این نظریه و کسب نتایج دقیق از آن مستلزم برقراری مفروضه بنیادی آن است. هر چند برخی شواهد بیانگر مقاوم بودن نظریه سؤال - پاسخ به تخطی از مفروضه‌های بنیادی‌اش است (همبلتون، ۱۹۸۹). یکی از مفروضه‌های مدل، تک‌بعدی بودن^۱ است، به این معنا که یک خصیصه مکنون با مجموعه سؤال‌های آزمون اندازه‌گیری شود. برقراری مطلق این مفروضه به دلیل دخالت عوامل گوناگون (اضطراب، انگیزش، مهارت‌ها و فرایندهای شناختی و وضعیت جسمانی و روانی آزمودنی‌ها) بر عملکرد آزمودنی‌ها غیرممکن است، اما آن چه اهمیت دارد وجود عاملی بارز در آزمون است که معرف عملکرد آزمودنی‌ها در آزمون باشد. روش‌های مختلفی برای بررسی برقراری مفروضه تک‌بعدی هست که یکی از مناسب‌ترین آنها تحلیل عاملی^۲ است. در مطالعه حاضر تحلیل عاملی سؤال‌ها با استفاده از نرم‌افزار TESTFACT و بر اساس کل داده‌ها انجام شده است. ارزش ویژه نخستین عامل استخراج شده در تحلیل عاملی برابر با ۱۵/۶۳ است و ۲۸/۳۶ درصد واریانس را تبیین می‌کند. دومین ارزش ویژه برابر با ۲/۰۱ است که ۳/۰۵ درصد واریانس را تبیین می‌کند. طبق اظهار ریکازی^۳ هنگامی که عامل نخست حداقل ۲۰ درصد واریانس کل را تبیین کند، می‌توان نتیجه گرفت آزمون تک‌بعدی است (احمدی، ۲۰۰۸). بر این اساس در مطالعه حاضر آزمون تک‌بعدی است و تحلیل آزمون بر اساس نظریه سؤال - پاسخ تک‌بعدی امکان پذیر است.

جدول (۱) عامل‌های استخراج شده آزمون ریاضی

عامل	ارزش ویژه	درصد واریانس تبیین شده
۱	۱۵/۶۳	۲۸/۳۶
۲	۲/۰۱	۳/۰۵
۳	۱/۳۵	۲/۴۴
۴	۱/۱۷	۲/۱۱
۵	۱/۰۹	۱/۹۸

1. Unidimensionality
2. Factor analysis
3. Reakase

بررسی وجود کارکرد افتراقی سئوال‌ها با استفاده از نرم‌افزار BILOG-MG صورت گرفته است. این نرم‌افزار توانایی مدل‌سازی چندگروهی نظریه سئوال - پاسخ را دارد. بنابراین می‌توان از آن برای ارزیابی کارکرد افتراقی سئوال استفاده کرد. نرم‌افزار BILOG-MG کارکرد افتراقی سئوال را از نظر پارامترهای دشواری سئوال بررسی می‌کند و قادر به بررسی تفاوت‌های گروهی در پارامتر تشخیص نیست (امبرتسون و رایس، ترجمه شریفی و همکاران، ۱۳۸۸). نکته قابل تأمل در کاربرد روش سئوال-پاسخ، برازش مدل با داده‌ها است. در واقع، برتری‌های رویکرد سئوال-پاسخ تنها در صورت برازش داده‌ها با مدل برقرار است. برازش مدل انتخابی با استفاده از شاخص لگاریتم درست‌نمایی صورت گرفت. تفاوت مقدار این شاخص در مدل‌های مختلف دارای توزیع مجذور کای با درجه آزادی برابر با تعداد پارامترهای افزوده برای مدل پیچیده‌تر است (امبرتسون و رایس، ترجمه شریفی و همکاران، ۱۳۸۸). مقادیر این شاخص در جدول (۲) ارائه شده است. با توجه به جدول، مدل دو پارامتری با داده‌ها برازش دارد.

جدول (۲) مقادیر شاخص برازش لگاریتم درست‌نمایی مدل‌های مختلف

مدل سئوال - پاسخ	لگاریتم درست‌نمایی
تک پارامتری	۴۷۵۲۵/۴۸۹
دو پارامتری	۴۶۸۰۶/۸۷۷۹
سه پارامتری	۴۶۸۱۱/۱۴۳۹

پس از کسب اطمینان از برقراری مفروضه‌های مدل و برازش مدل با داده‌ها، در گام اول به منظور تعیین وجود کارکرد افتراقی در سئوال‌ها، داده‌ها دوبار تحلیل می‌شوند. در نوبت اول داده‌ها به طور کلی تحلیل می‌شوند به گونه‌ای که از جامعه یکسانی انتخاب شده‌اند. در نوبت دوم داده‌ها در زیرگروه‌ها با استفاده از تحلیل کارکرد افتراقی با فرض صفر عدم وجود کارکرد افتراقی، تحلیل می‌شوند. تفاوت لگاریتم درست‌نمایی دو مرحله دارای توزیع مجذور کای با درجه آزادی $(n-1)(m-1)$ است که n معرف تعداد سئوال‌ها و m معرف تعداد گروه‌ها است. هنگامی که مجذور کای معنادار است، شواهدی از کارکرد افتراقی وجود دارد. بر این اساس در مطالعه حاضر داده‌ها یک‌بار به طور کلی و یک‌بار در دو گروه داوطلبان دختر و پسر با استفاده از تحلیل کارکرد افتراقی تحت فرضیه صفر مبنی بر عدم وجود کارکرد افتراقی

تحلیل شدند. مقدار لگاریتم درستنمایی در مرحله اول برابر با $۱۸۸۰۳۴/۸۶$ و در مرحله دوم $۱۸۷۷۸۴/۰۲$ است که تفاوت آن $۲۵۰/۸۴$ است و در سطح $۰/۰۵$ با درجه آزادی ۵۴ معنادار است که این نتیجه بیانگر وجود اثرات کارکرد افتراقی در آزمون ریاضی است.

پس از تأیید وجود کارکرد افتراقی در آزمون، در گام بعدی تحلیل کارکرد افتراقی سؤال‌ها انجام شد. نرم‌افزار کارکرد افتراقی را فقط بر اساس پارامتر دشواری بررسی می‌کند. کارکرد افتراقی در سه مرحله مورد بررسی قرار می‌گیرد. در مرحله نخست پارامتر سؤال‌ها برای هر یک از دو گروه آزمودنی به طور جداگانه برآورد می‌شود. البته پارامتر تشخیص سؤال‌ها برای دو گروه یکسان در نظر گرفته می‌شود. برنامه BILOG-MG میانگین و انحراف استاندارد گروه‌ها را یکسان فرض نمی‌کند. میانگین و انحراف استاندارد گروه مرجع به ترتیب صفر و یک است، میانگین و انحراف استاندارد گروه کانونی پارامترهایی هستند که برآورد می‌شوند. در مرحله دوم میانگین دشواری سؤال‌ها برای گروه مرجع و کانونی محاسبه می‌شود و تفاوت میانگین دشواری سؤال‌های دو گروه کانونی و مرجع محاسبه می‌شود. پارامتر دشواری سؤال‌های گروه کانونی با کم کردن تفاوت میانگین دشواری دو گروه از پارامتر دشواری سؤال‌های گروه کانونی، تعدیل می‌شوند. در مرحله سوم تفاوت میان پارامترهای دشواری سؤال‌های مربوط به گروه مرجع و کانونی و خطای استاندارد محاسبه می‌شود، تفاوت زیاد نشان‌دهنده کارکرد افتراقی سؤال می‌باشد. در واقع سؤال‌هایی که تفاوت میان پارامتر دشواری آنها در دو گروه تقریباً بیش از دو برابر خطای استاندارد ($S. E \times 1/96$) است به عنوان سؤال‌هایی با کارکرد افتراقی در نظر گرفته می‌شوند. در تحلیل حاضر با توجه به این‌که داوطلبان دختر گروه کانونی و داوطلبان پسر، گروه مرجع را تشکیل می‌دهند، مقادیر منفی تفاوت میان پارامترهای دشواری سؤال‌ها در گروه مرجع و کانونی بیانگر آن است که سؤال به نفع داوطلبان دختر است و برعکس. نتایج تحلیل سؤال‌ها آزمون در جدول (۳) ارائه شده است. همان گونه که در جدول مشخص شده است نه سؤال آزمون ($۱۶/۴$ درصد سؤال‌ها) کارکرد افتراقی دارند که تمام این سؤال‌ها به نفع آزمودنی‌های دختر سوگیری دارند (سؤال‌های شماره ۵، ۷، ۹، ۱۸، ۲۷، ۳۲، ۳۶، ۴۱ و ۴۹).

جدول (۳) نتایج تحلیل مبنی بر رویکرد سئوال - پاسخ

DIF	تفاوت دشواری گروه‌ها	سئوال	DIF	تفاوت دشواری گروه‌ها	سئوال
	داوطلبان پسر - داوطلبان دختر			داوطلبان پسر - داوطلبان دختر	
	۰/۰۵۴ *۰/۱۵۹	۲۹		-۰/۰۳ *۰/۰۵۴	۱
	۰/۰۷۱ *۰/۲۹۷	۳۰		۰/۰۵۴ *۰/۲۰۸	۲
	-۰/۰۲۴ *۰/۱۹۴	۳۱		-۰/۰۴۲ *۰/۰۶۵	۳
+	-۰/۳۴۶ *۰/۰۶۴	۳۲		۰/۰۱۲ *۰/۳۶۷	۴
	۰/۳۱ *۰/۴۲۱	۳۳	+	-۰/۲۱۳ *۰/۰۶	۵
	۰/۲۵۱ *۰/۲۵۸	۳۴		-۰/۰۳۲ *۰/۰۶۹	۶
	۰/۱۵۹ *۰/۰۹۸	۳۵	+	-۰/۲۸۵ *۰/۱۱۱	۷
+	-۰/۱۸۸ *۰/۰۹۴	۳۶		۰/۱۵۲ *۰/۰۹۱	۸
	-۰/۲۰۱ *۰/۱۱۶	۳۷	+	-۰/۱۸ *۰/۰۵۲	۹
	۰/۰۱ *۰/۲۰۵	۳۸		۰/۶۲۲ *۰/۴۲	۱۰
	-۰/۰۱۸ *۰/۲۷۲	۳۹		-۰/۱۴۵ *۰/۱۰۷	۱۱
	-۰/۸ *۰/۱۸۴	۴۰		۰/۲۰۸ *۰/۵۸۸	۱۲
+	-۰/۲۰۲ *۰/۰۷۷	۴۱		۰/۲۸۴ *۰/۳۳۳	۱۳
	-۰/۰۲۹ *۰/۰۹۹	۴۲		۰/۳۱۱ *۰/۳۴۸	۱۴
	۰/۲۸۱ *۰/۳۲۲	۴۳		-۰/۲۴۴ *۰/۲۸۹	۱۵
	۰/۳۹۴ *۰/۳۰۶	۴۴		۰/۳۷۴ *۰/۵۲	۱۶

DIF	تفاوت دشواری گروه‌ها	سئوال	DIF	تفاوت دشواری گروه‌ها	سئوال
	داوطلبان پسر - داوطلبان دختر			داوطلبان پسر - داوطلبان دختر	
	-۰/۳۰۹ *۰/۱۶۵	۴۵		۰/۰۳۸ *۰/۱۲۴	۱۷
	۰/۰۶۳ *۰/۱۵۹	۴۶	+	-۰/۳۸ *۰/۰۷۶	۱۸
	۰/۰۵ *۰/۱۰۵	۴۷		۰/۱۱۶ *۰/۲۴۸	۱۹
	-۰/۲۲۵ *۰/۱۴۶	۴۸		۰/۱۲۷ *۰/۴۵۳	۲۰
+	-۲/۱۲۶ *۰/۶۸۴	۴۹		۰/۰۰۷ *۰/۱۱۲	۲۱
	۰/۰۸۵ *۰/۳۳۹	۵۰		-۰/۰۰۴ *۰/۱۱	۲۲
	۰/۰۲۳ *۱/۱۰۸	۵۱		-۰/۱۸۷ *۰/۱۹۷	۲۳
	۰/۲۰۶ *۰/۳۱۸	۵۲		-۰/۱۳۸ *۰/۰۷۲	۲۴
	۰/۲۵۴ *۰/۱۷۱	۵۳		۰/۰۵۲ *۰/۱۶	۲۵
	-۰/۰۲۹ *۰/۱۷۱	۵۴		-۰/۱۱۷ *۰/۰۸۲	۲۶
	۱/۳۰۶ *۱/۰۳۹	۵۵	+	-۰/۱۹۷ *۰/۰۷۵	۲۷
				۰/۱۸۸ *۰/۱۳۱	۲۸

(* خطای استاندارد، (+) وجود کارکرد افتراقی

در ادامه تحلیل کارکرد افتراقی سئوال‌ها با استفاده از روش مانتل هنزل مورد ارزیابی قرار گرفت. در مطالعه حاضر تحلیل مانتل هنزل با استفاده از نرم‌افزار Leartap انجام شده است. برای هر سئوال در هر سطح نمره، تعداد و نسبت افراد هریک از گروه‌های مرجع و کانونی که به سئوال، پاسخ درست داده‌اند، شاخص مجذور کای مانتل هنزل،

آلفای مانندل هنزل^۱، دلتای مانندل هنزل^۲، نسبت بخت‌ها، معناداری و اندازه اثر محاسبه می‌شود. نسبت بخت، اندازه‌ای نسبی است که احتمال این‌که فردی از گروه مرجع در مقایسه با گروه کانونی به سئوال پاسخ درست بدهد را نشان می‌دهد. مقادیر بزرگ‌تر از یک نشان‌دهنده این است که احتمال پاسخ‌گویی درست بیشتر به نفع گروه مرجع است و مقادیر کم‌تر از یک بیانگر آن است که احتمال پاسخ‌گویی درست بیشتر به نفع گروه کانونی است. آلفای مانندل هنزل عبارت است از میانگین نسبت بخت‌ها در تمام سطوح نمره که براساس تعداد افراد هر یک از گروه‌ها در هر سطح نمره وزن‌دهی شده‌اند. بخت، به عنوان شاخصی از احتمال وقوع یک حادثه، با این محدودیت همراه است که دامنه آن همواره بین صفر تا مثبت بی‌نهایت است یا به عبارتی ناقرینه^۳ است. به همین علت از آن لگاریتم گرفته و به این ترتیب لگاریتم آلفای مانندل هنزل به دست می‌آید که دامنه آن در بازه منفی بی‌نهایت تا مثبت بی‌نهایت است. سئوال‌ها با مقادیر مثبت به سود گروه مرجع و سئوال‌ها با مقادیر منفی به نفع گروه کانونی سوگیری دارند. هالند و تایر برای نشان دادن مقدار کارکرد افتراقی شاخص دیگری را پیشنهاد کردند که در واقع تبدیل آلفای مانندل هنزل به مقیاسی دیگر است که دلتای مانندل هنزل نامیده می‌شود و از رابطه زیر به دست می‌آید:

$$MH\ D-DIF = -2.35 * Ln\ MH\ alpha$$

سئوال‌ها با مقادیر دلتای مانندل هنزل مثبت به نفع گروه کانونی و سئوال‌ها با مقادیر منفی به نفع گروه مرجع هستند. خدمات سنجش آموزشی^۴ طبقه‌بندی سه سطحی را برای کارکرد افتراقی سئوال ارائه کرده است (زویک^۵، ۲۰۱۲) که عبارتند از:

۱- سطح A: قدر مطلق دلتای مانندل هنزل سئوال‌های این طبقه کم‌تر از ۱ است. سئوال‌های این طبقه فاقد کارکرد افتراقی یا کارکرد افتراقی ناچیزی دارند.

۲- سطح B: سئوال‌ها با قدر مطلق دلتای مانندل هنزل بین ۱ تا ۱/۵ در این طبقه قرار می‌گیرند. این سئوال‌ها کارکرد افتراقی متوسطی دارند.

-
1. MH alpha
 2. MH D-DIF
 3. Asymmetry
 4. Educational Testing Service (ETS)
 5. Zwick

۳- سطح C: سئوال‌های این طبقه با شاخص مجذور کای مانتل هنزل معنادار و قدر مطلق دلتای مانتل هنزل بیش از ۱/۵ مشخص می‌شوند. کارکرد افتراقی سئوال‌ها زیاد است و چنین سئوال‌هایی از آزمون حذف می‌شوند (زویک، ۲۰۱۲).
نتایج تحلیل مانتل هنزل سئوال‌های آزمون در جدول (۵) ارائه شده است.

جدول (۵) نتایج تحلیل مانتل هنزل

سئوال	شاخص مانتل هنزل	معناداری	آلفای مانتل هنزل	دلتای مانتل هنزل	طبقه‌بندی ETS
۱	۰/۶۱	۰/۴۴	۱/۰۶	-۰/۱۳	A
۲	۰/۶۵	۰/۴۲	۱/۱۰	-۰/۲۲	A
۳	۰/۲	۰/۶۶	۱/۰۳	-۰/۰۸	A
۴	۰/۱۴	۰/۷۱	۱/۰۴	-۰/۱۰	A
۵	۷/۸۱	*۰/۰۱	۰/۸۰	۰/۵۲	A
۶	۰/۲۹	۰/۵۹	۱/۰۴	-۰/۱۰	A
۷	۷/۷۲	*۰/۰۱	۰/۸۱	۰/۵۰	A
۸	۹/۲۴	*۰/۰۰	۱/۳۱	-۰/۶۳	A
۹	۴/۲۲	*۰/۰۴	۰/۸۶	۰/۳۶	A
۱۰	۱۴/۳۶	*۰/۰۰	۱/۵۹	-۱/۰۹	B
۱۱	۰/۷۷	۰/۳۸	۰/۹۴	۰/۱۵	A
۱۲	۲/۱۶	۰/۱۴	۱/۱۳	-۰/۲۸	A
۱۳	۴/۵۹	*۰/۰۳	۱/۲۴	-۰/۵۰	A
۱۴	۵/۲۶	*۰/۰۲	۱/۳۶	-۰/۷۲	A
۱۵	۱/۴۹	۰/۲۲	۰/۹۱	۰/۲۱	A
۱۶	۴/۵۷	*۰/۰۳	۱/۳۲	-۰/۶۵	A
۱۷	۰/۰۱	۰/۹۳	۰/۹۹	۰/۰۲	A
۱۸	۲۳/۵۵	*۰/۰۰	۰/۷۰	۰/۸۵	A
۱۹	۱/۰۴	۰/۳۱	۱/۱۳	-۰/۲۹	A
۲۰	۱/۰۷	۰/۳۰	۱/۱۱	-۰/۲۵	A
۲۱	۰/۴۴	۰/۵۱	۱/۰۶	-۰/۱۴	A
۲۲	۰/۲۳	۰/۶۳	۱/۰۴	-۰/۰۹	A
۲۳	۱/۵۶	۰/۲۱	۰/۸۸	۰/۳۱	A
۲۴	۱/۳۷	۰/۲۴	۰/۹۱	۰/۲۳	A
۲۵	۱/۴۹	۰/۲۲	۱/۱۴	-۰/۳۰	A
۲۶	۱/۱۴	۰/۲۹	۰/۹۰	۰/۲۵	A

سئوال	شاخص مانتل هنزل	معناداری	آلفای مانتل هنزل	دلتای مانتل هنزل	طبقه‌بندی ETS
۲۷	۶/۱۲	*۰/۰۱	۰/۷۵	۰/۶۷	A
۲۸	۵/۶۲	*۰/۰۲	۱/۲۰	-۰/۴۳	A
۲۹	۰/۵۶	۰/۴۵	۱/۰۸	-۰/۱۸	A
۳۰	۰/۳۹	۰/۵۳	۱/۰۸	-۰/۱۷	A
۳۱	۰/۰۲	۰/۹۰	۱/۰۱	-۰/۰۳	A
۳۲	۲۳/۰۹	*۰/۰۰	۰/۶۹	۰/۸۶	A
۳۳	۴/۳۲	*۰/۰۴	۱/۳۷	-۰/۷۴	A
۳۴	۴/۰۶	*۰/۰۴	۱/۳۰	-۰/۶۲	A
۳۵	۶/۲۸	*۰/۰۱	۱/۱۹	-۰/۴۲	A
۳۶	۳/۴۱	۰/۰۶	۰/۸۶	۰/۳۴	A
۳۷	۲/۵۳	۰/۱۱	۰/۸۸	۰/۳۰	A
۳۸	۰/۰۶	۰/۸۱	۱/۰۳	-۰/۰۸	A
۳۹	۰/۰۲	۰/۹۰	۰/۹۸	۰/۰۵	A
۴۰	۰/۰	۰/۹۵	۰/۹۹	۰/۰۲	A
۴۱	۵/۰۶	*۰/۰۲	۰/۸۱	۰/۵۱	A
۴۲	۰/۱۱	۰/۷۴	۱/۰۳	۰/۰۷	A
۴۳	۵/۰	*۰/۰۳	۱/۳۸	-۰/۷۵	A
۴۴	۹/۰۶	*۰/۰۰	۱/۳۶	-۰/۷۳	A
۴۵	۴/۷۹	*۰/۰۳	۰/۸۶	۰/۳۶	A
۴۶	۱/۲۷	۰/۲۳	۱/۱۰	-۰/۲۳	A
۴۷	۲/۰۹	۰/۱۵	۱/۱۴	-۰/۳۰	A
۴۸	۴/۴۲	*۰/۰۴	۰/۸۴	۰/۴۰	A
۴۹	۲۱/۶۴	*۰/۰۰	۰/۷۴	۰/۷۱	A
۵۰	۰/۳۰	۰/۵۸	۱/۰۵	-۰/۱۳	A
۵۱	۰/۰۵	۰/۸۲	۱/۰۳	-۰/۰۸	A
۵۲	۲/۴۷	۰/۱۲	۱/۲۷	-۰/۵۷	A
۵۳	۶/۸۴	*۰/۰۱	۱/۳۲	-۰/۶۶	A
۵۴	۰/۰۴	۰/۸۴	۱/۰۲	-۰/۰۴	A
۵۵	۱۸/۸۸	*۰/۰۰	۱/۴۰	-۰/۷۸	A

همان‌گونه که در جدول مشخص شده است، شاخص مانتل هنزل ۲۳ سئوال (۴۱/۸ درصد سئوال‌ها) در سطح ۰/۰۵ معنادار است. از این ۲۳ سئوال با توجه به مقادیر دلتای مانتل هنزل ۲۲ سئوال در سطح A هستند، به این معنا که مقدار دلتای مانتل هنزل آنها

کم‌تر از یک است و در واقع با وجود معنادار بودن شاخص مانتل هنزل، کارکرد افتراقی ناچیزی دارند و می‌توان این سئوال‌ها را در آزمون حفظ کرد. سئوال شماره ۱۰ در سطح B است و دارای کارکرد افتراقی متوسط است. هیچ یک از سئوال‌های آزمون کارکرد افتراقی زیاد از خود نشان نداده‌اند (سطح C). از این ۲۳ سئوال که شاخص مانتل هنزل معنادار دارند، ۱۰ سئوال به نفع داوطلبان دختر (سئوال‌های شماره ۵، ۷، ۹، ۱۸، ۲۷، ۳۲، ۴۱، ۴۵، ۴۸ و ۴۹) و ۱۳ سئوال (سئوال‌های شماره ۸، ۱۰، ۱۳، ۱۴، ۱۶، ۲۸، ۳۳، ۳۴، ۳۵، ۴۲، ۴۴، ۵۳ و ۵۵) به نفع داوطلبان پسر است.

همان‌گونه که مشاهده می‌شود روش مبتنی بر نظریه سئوال - پاسخ تعداد سئوال‌های کم‌تری را به عنوان سئوال‌های دارای کارکرد افتراقی مشخص کرده است (۹ سئوال در مقابل ۲۳ سئوال) که این یافته با نتایج مطالعه دورانند و پارک همسو است. یک دلیل تفاوت حساسیت روش‌های مختلف تعیین کارکرد افتراقی ممکن است مربوط به تعدیلات ریاضی خطای نوع اول باشد که در هر برنامه‌ای به گونه متفاوتی اجرا می‌شود. بررسی سئوال‌هایی که دارای کارکرد افتراقی هستند بیانگر آن است که سئوال‌هایی که به نفع دختران کارکرد افتراقی دارند در حوزه محتوایی توابع و معادلات هستند و سئوال‌هایی که به نفع پسران کارکرد افتراقی دارند بیشتر در حوزه محتوایی مثلثات، هندسه و احتمال می‌باشند. البته بررسی دقیق‌تر سئوال‌های آزمون توسط کارشناسان موضوعی به منظور روشن ساختن منبع و دلیل کارکرد افتراقی ضروری می‌باشد.

به منظور بررسی توافق و هم‌خوانی میان دو روش مانتل هنزل و روش مبتنی بر رویکرد سئوال - پاسخ در تشخیص کارکرد افتراقی، ضریب کاپا^۱ محاسبه شد و به این ترتیب میزان هم‌خوانی روش‌ها در تعیین سئوال‌های دارای کارکرد افتراقی مشخص شد. در جدول (۶) فراوانی سئوال‌های دارای کارکرد افتراقی و فاقد کارکرد افتراقی براساس دو رویکرد ارائه شده است. همان‌گونه که در جدول مشخص شده است، دو روش در تعیین هشت سئوال دارای کارکرد افتراقی و در ۳۱ سئوال فاقد کارکرد افتراقی مطابقت داشتند. مقادیر ضریب کاپا در مطالعه حاضر برابر با ۰/۳۴۶ و $p=۰/۰۰۲$ و فاصله اطمینان ۹۵ درصد (۰/۵۶۵، ۰/۱۲۶) است. پیشنهاد لاندیس^۲ و کخ^۳ (۱۹۹۷) این است که کاپای بیش از ۰/۷۵ حاکی از توافق عالی، کاپای کمتر از

1. Kappa coefficient

2. Landis

3. Kokh

۰/۴ نشان‌دهنده توافق ضعیف و بالاخره کاپای بین ۰/۴ تا ۰/۷۵ نشان از توافق نسبی خوب دارد. در مطالعه حاضر میزان توافق میان دو روش، ضعیف است. بیشتر مطالعات انجام شده در خصوص توافق روش‌های مختلف ارزیابی کارکرد افتراقی بیانگر توافق کم میان روش‌های مختلف است. همگرایی روش‌های بررسی کارکرد افتراقی به شدت تحت تأثیر ناپایایی شاخص‌های ارزیابی کارکرد افتراقی است (عبدالعزیز، ۲۰۱۰).

جدول (۶) فراوانی سئوال‌های دارای کارکرد افتراقی و فاقد کارکرد افتراقی

تعداد کل	رویکرد سئوال - پاسخ		دارای کارکرد افتراقی	مانتل هنزل
	فاقد کارکرد افتراقی	دارای کارکرد افتراقی		
۲۳	۱۵	۸	دارای کارکرد افتراقی	مانتل هنزل
۳۲	۳۱	۱	فاقد کارکرد افتراقی	
۵۵	۴۶	۹	تعداد کل	

نتیجه‌گیری

مطالعه حاضر به ارزیابی کارکرد افتراقی سئوال‌های آزمون ریاضی با استفاده از دو روش مانتل هنزل و روش مبتنی بر نظریه سئوال - پاسخ پرداخته است. تحلیل مبتنی بر رویکرد سئوال - پاسخ با استفاده از نرم افزار BILOG-MG انجام شده است. این برنامه کارکرد افتراقی را صرفاً براساس پارامتر دشواری بررسی می‌کند. نظریه سئوال - پاسخ روش مناسبی برای ارزیابی کارکرد افتراقی سئوال‌ها محسوب می‌شود که این به واسطه ویژگی نامتغیر بودن پارامترها در این رویکرد است. نتایج تحلیل، بیانگر وجود کارکرد افتراقی در سئوال‌های آزمون ریاضی است. براساس رویکرد سئوال - پاسخ، نه سئوال کارکرد افتراقی دارند که تمامی آنها نیز به نفع آزمودنی‌های دختر هستند. در تحلیل آزمون براساس روش مانتل هنزل، شاخص مانتل هنزل ۲۳ سئوال آزمون در سطح ۰/۰۵ معنادار بود که ده سئوال به نفع آزمودنی‌های دختر و ۱۳ سئوال به نفع آزمودنی‌های پسر است. با توجه به طبقه‌بندی خدمات سنجش آموزشی ۲۲ سئوال دارای کارکرد افتراقی ناچیز می‌باشند (سطح A) و یک سئوال کارکرد افتراقی متوسط دارد (سطح B). شایان ذکر است که از نه سئوالی که بر اساس رویکرد سئوال - پاسخ دارای کارکرد افتراقی هستند، هشت مورد آن نیز در روش مانتل هنزل دارای کارکرد

افتراقی می‌باشند. تنها سؤال شماره ۳۶ است که براساس روش مبتنی بر رویکرد سؤال - پاسخ دارای کارکرد افتراقی است، اما در روش مانتل هنزل فاقد کارکرد افتراقی است. میزان توافق و هم‌خوانی میان دو روش مانتل هنزل و روش مبتنی بر رویکرد سؤال - پاسخ در تشخیص کارکرد افتراقی با توجه به ضریب کاپا $0/۳۴۶$ ضعیف است. رویکرد سؤال - پاسخ در مقایسه با روش مانتل هنزل تعداد کمتری از سؤال‌ها را دارای کارکرد افتراقی تشخیص داده است، این یافته مشابه نتایج مطالعه دورانند و پارک (۲۰۰۶) است. بر این اساس سازندگان آزمون‌ها می‌بایست همواره به تفاوت‌های روش‌های مختلف ارزیابی کارکرد افتراقی سؤال‌ها توجه داشته باشند و سؤال‌هایی را که بر اساس یک روش دارای کارکرد افتراقی هستند با سایر روش‌های معتبر نیز مورد بررسی قرار دهند.

بررسی سؤال‌های دارای کارکرد افتراقی در آزمون ریاضی بیانگر آن است که سؤال‌های دارای کارکرد افتراقی به نفع دختران بیشتر در حوزه محتوایی توابع و معادلات و سؤال‌های دارای کارکرد افتراقی به نفع پسران در حوزه محتوایی مثلثات، هندسه و احتمال هستند. پس از تعیین کارکرد افتراقی، کارشناسان موضوعی می‌بایست محتوای سؤال‌ها را به منظور روش ساختن منبع و دلیل کارکرد افتراقی مورد بازبینی کامل قرار دهند. از این رو لازم است در مطالعات بعدی سؤال‌های آزمون که دارای کارکرد افتراقی بودند از لحاظ محتوایی توسط کارشناسان موضوعی مورد بررسی قرار گیرند. وجود سؤال‌های دارای کارکرد افتراقی در آزمون، ضرورت بررسی و توجه بیشتر به فرایند طراحی سؤال‌های آزمون و آموزش جامع طراحان سؤال را به منظور تهیه آزمون‌های عادلانه و منصفانه مطرح می‌سازد.

منابع

- امیرتسون، سوزان ای و رایس، استیون پی (۲۰۰۰). نظریه‌های جدید روانسنجی برای روانشناسان (IRT). ترجمه: دکتر حسن پاشا شریفی، دکتر ولی‌الله فرزاد، مجتبی حیبی عسگرآباد و بلال ایزانلو (۱۳۸۸). تهران: انتشارات رشد.
- همبلتون، رونالد ک، سوامینانان، اچ و راجرز، جین. (۱۹۹۱). مبانی نظریه پرسش-پاسخ. ترجمه: دکتر محمدرضا فلسفی نژاد (۱۳۸۹). تهران: انتشارات دانشگاه علامه طباطبایی.
- Abedalaziz, Nabeel. (2010). A gender - related differential item functioning of mathematics test item. *The International Journal of Educational and Psychological Assessment*, 5,101-116 .
- Ahmadi, Alireza. (2008). *Differential Item Functioning in High-stakes Tests: the Effect of Gender and Field of Study*. Doctoral dissertation, University of Isfahan, Faculty of Foreign Languages, Department of English .
- Baghi, Heibatollah & Ferrara, Steven. (1989). *A comparison of IRT, Delta Plot and Mantel- Haenszel techniques for detecting differential item functioning across subpopulations in the Maryland test of citizenship skills*. Paper presented at the annual meeting of the American Educational Research Association. San Francisco, CA: March 27-31 (Eric database, Report: ED324364) .
- Conoley, C. Adele (2003). *Differential item functioning in the Peabody picture vocabulary test - third edition: partial correlation versus expert judgment*. Doctoral dissertation. Texas A&M University .
- Dorans, Neil J. & Holland, Paul W. (1992). *DIF Detection and Description: Mantel - Haenszel and Standardization*. Paper presented at the Educational Testing Service/AFHRL Conference (Princeton. NJ. October)
- 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 .
- Durand, Jeffrey & Park, Siwo. (2006). *A Study of Gender and Academic Major - Based Differential Item Functioning (DIF) In KEPT 2006, Mexico* .

- Eng. L. S. & Hoe. L. S. (2005). Detecting Differential Item Functioning (DIF) in Standardized Multiple-Choice Test: An Application of Item Response Theory (IRT). [ONLINE] Available at: <http://www.ipbl.edu.my/inter/penyelidikan/seminarpapers/2005/linguitm.pdf>.
- Hambleton, R. , & Rodgers, J. (1995). Item bias review. *Practical Assessment, Research, and valuation*, 4(6). Retrieved November 18, 2006, from <http://PAREonline.net/getvn.asp?v=4&n=6>
- Hambleton, R. K. (1989). Principles and selected applications of item response theory. In R. L. Linn (Ed.) *Educational measurement* (3rd ed ,pp. 147-200). New York NY: American Council on Education & Macmillan Publishing .
- Landis, J. R. , Koch, G. G. (1977). The measurement of observer agreement for categorical data. *Biometrics* 33:159-174.
- 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 .
- O'Neal, Marcia R. (1991). A Comparison of Method for Detecting Item Bias. Paper presented at the annual meeting of the Mid-South Educational Research Association (20th, Lexington, KY, November 12-15 .
- Oshima, T. C. and Morris, S. B. (2008), Raju's Differential Functioning of Items and Tests (DFIT). *Educational Measurement: Issues and Practice*, 27: 43-50. doi: 10. 1111/j. 1745-3992. 2008. 00127 .
- Rousseau, M. , Bertrand, R. , & Boiteau, N. (2004). *Impact of missing data on robustness of DIF IRT-based and Non-IRT-based methods*. Paper presented at the 2004 AERA annual meeting .
- Shultz, S. Kenneth & Whitney, David. (2005). *Measurement Theory in Action. Case Studies and Exercises*. Sage publication .
- Willingham, W. W. & Cole, N. S; (1997). *Gender and fair assessment*. New Jersey, U. S. A: Lawrence Erlbaum associate
- 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 .Rebecca (2012). A Review of ETS Differential Item Functioning Assessment Procedures: Flagging Rules, Minimum Sample Size Requirements, and Criterion Refinement. [ONLINE] Available at: <http://www.ets.org/Media/Research/pdf/RR-12-08.pdf>.