

کنش افتراقی سؤال (DIF) و آزمون (DTF) مرتبط با جنسیت در آزمون‌های کنکور سراسری سازمان سنجش آموزش کشور^۱

مسعود گرامی‌پور*
اعظم رضانی صدر**
احمد رضایی***
لیلا نوروزی****
فرانک مختاریان*****

چکیده

عملکرد آزمودنی با کنترل آماری توانایی ایشان در گروه‌های مختلف جنسیتی ممکن است متفاوت باشد. در این صورت، وقوع کنش افتراقی سؤال (DIF) / آزمون (DTF) مرتبط با جنسیت و به دنبال آن سوگیری سؤال/آزمون محتمل است. هدف اصلی از اجرای این پژوهش، مطالعه DTF/DIF و سوگیری جنسیتی در آزمون‌های کنکور سراسری ایران است. آزمون‌های تخصصی یک دفترچه آزمون در پنج گروه آزمایشی شرکت‌کننده در کنکورهای سراسری سال‌های ۱۳۸۷ تا ۱۳۹۰ به صورت خوشه‌ای تک مرحله‌ای انتخاب شدند. سپس برای مطالعه DIF از تحلیل رگرسیون لوجستیک و برای مطالعه DTF از رویکرد مبتنی بر IRT استفاده شد. نتایج نشان داد که به‌طور متوسط، حدود ۱۴ درصد از پرسش‌های آزمون‌های مورد مطالعه دارای DIF جنسیتی با اندازه اثر ناچیز ($EF < 0.001$) هستند و حداکثر ۲ درصد آنها بر اساس نقطه نظرات گروه‌های خبرگان نسبت به جنس مؤنث یا مذکر سوگیری دارند. همچنین یافته‌های تحلیل DTF نشان داد که به جز آزمون خلاقیت‌های نمایشی گروه هنر که نسبت به زنان دارای سوگیری کمی است، سایر آزمون‌ها فاقد DTF هستند. واژگان کلیدی: کنش افتراقی سؤال/آزمون، DTF/DIF مرتبط با جنسیت، سوگیری سؤال/آزمون، آزمون‌های ورودی دانشگاه، کنکور سراسری

تاریخ دریافت مقاله: ۹۵/۰۷/۱۲

تاریخ پذیرش مقاله: ۹۶/۰۳/۲۷

^۱. این پژوهش با حمایت سازمان سنجش آموزش کشور اجرا شده است.

* استادیار دانشکده روان‌شناسی و علوم تربیتی دانشگاه خوارزمی (نویسنده مسئول):
mgramipour@yahoo.com

** کارشناسی ارشد تحقیقات آموزشی، دانشگاه خوارزمی
*** دانشجوی دکتری رشته سنجش و اندازه‌گیری دانشگاه علامه طباطبایی
**** کارشناسی ارشد تحقیقات آموزشی، دانشگاه خوارزمی
***** کارشناس پژوهشی سازمان سنجش آموزش کشور

مقدمه

آزمون‌های بزرگ‌مقیاسی مانند کنکور سراسری که سازمان سنجش آموزش کشور در ایران برگزار می‌کند، نقش کلیدی در آینده جوانان دارند. دختران و پسران بعد از گذراندن دوره متوسطه برای ورود به دانشگاه و رشته مورد علاقه خود، باید با یکدیگر در کنکور سراسری رقابت کنند. رقابت در چنین آزمونی مشکل بوده و شاید بتوان گفت مهم‌ترین آزمون در ایران است که به صورت سالیانه برگزار می‌شود. والدین دانش‌آموزان نیز زمان و هزینه بسیاری متحمل می‌شوند تا فرزندان آنها در کنکور سراسری پذیرفته شوند. یکی از مهم‌ترین متغیرهایی که در چنین آزمونی باید مد نظر قرار گیرد، جنسیت است. به طوری که این اطمینان حاصل شود که آزمون برای دختران و پسران به یک اندازه منصفانه است. بنابراین مطالعه روایی^۱ آزمون بر اساس جنسیت بسیار اهمیت دارد. در دهه‌های گذشته، متخصصان اندازه‌گیری و پژوهشگران درباره کارکرد متفاوت آزمون‌های روانی در گروه‌های مختلف (به طور مثال زنان و مردان) به طور فزاینده‌ای نگران شده‌اند. به همین دلیل، پژوهش‌های زیادی درباره سوگیری سؤال در آزمون‌های روانی و تربیتی انجام شده است (امبرتسون^۲، ۲۰۰۰). در چنین پژوهش‌هایی، پرسش این چنین مطرح می‌شود: آیا نمره‌های زنان و مردان در جامعه‌های مختلف در یک مقیاس اندازه‌گیری قابل مقایسه است؟ زمانی گفته می‌شود سؤالی دارای کنش افتراقی (DIF^۳) است که ویژگی‌های روان‌سنجی آن در گروه مورد نظر (گروه کانونی^۴) در مقایسه با گروه مورد مقایسه (مرجع^۵) متغیر باشد (دراسگو^۶، ۱۹۸۴، ۱۹۸۷). کارکرد افتراقی آزمون (DTF^۷) نیز زمانی اتفاق می‌افتد که توابع نمره‌های کل آزمون در گروه‌ها (مثلاً جنسیت) متفاوت باشد. در آزمون‌های دارای DTF نمره‌های واقعی مورد انتظار آزمون در دو گروه آزمودنی‌ها یکسان نیست. بنابراین می‌توان گفت که کنش افتراقی در دو سطح

-
1. Validity
 2. Embretson
 3. Differential Item Functioning
 4. Focal group
 5. Reference
 6. Drasgow
 7. Differential Test Functioning

اتفاق می‌افتد: ۱- سطح سؤال (DIF)؛ ۲- مجموع سؤال‌های آزمون (DTF) (رادنر^۱، جتسون^۲ و نایت^۳، ۱۹۸۰).

در پیشرفت‌های اخیر حوزه DTF دو نوع آزمون کلی^۴ آن شامل sDTF و uDTF که مبتنی بر نظریه پاسخ سؤال (IRT^۵) است، توسعه یافته است. در این روش‌ها، کنش افتراقی آزمون با رویه‌های انتگرال‌گیری سطوح میان منحنی‌های کل آزمون در دو گروه آزمودنی (مثلاً جنسیت مذکر و مؤنث) محاسبه می‌شود. در این رویکرد sDTF با بررسی اختلاف دو گروه مرجع و کانونی، مقداری علامت‌دار مبنی بر ارجحیت گروه کانونی (مثبت) یا گروه مرجع (منفی) فراهم می‌کند، در حالی که uDTF متوسط این اختلاف را بدون علامت مثبت یا منفی نشان می‌دهد. مقادیر منفی sDTF بر این دلالت دارد که نمره گروه مرجع به‌طور متوسط کمتر از گروه کانونی است، در حالی که مقادیر مثبت آن دلالت بر توفیق نمره‌های گروه کانونی دارد. مقدار uDTF بر قدر مطلق انحراف‌های مقادیر پرسش دلالت دارد که روی کل آزمون متراکم شده و نشانگر متوسط مساحت میان دو منحنی آزمون است. دامنه نمره‌های uDTF بین ۰ تا نمره آزمون است و اگر شکل تابعی منحنی‌های آزمون دقیقاً شکل هم باشد، مساحت میان منحنی‌ها صفر است. هنگامی که DTF کلی به‌واسطه sDTF و uDTF بررسی می‌شود، سه سناریوی پایانی ممکن است مشاهده شود: ۱- هنگامی که sDTF و uDTF هر دو کوچک است، در این حالت، مقدار DTF در کل دامنه توانایی بسیار کوچک یا صفر است؛ ۲- هنگامی که sDTF کوچک و uDTF بزرگ است، منحنی‌های آزمون در یک یا چند نقطه یکدیگر را قطع می‌کنند تا یک نمره کل متعادل را بسازند اما در این حالت، سوگیری غیر قابل اغمازی در نقاط به‌خصوصی از سطوح توانایی وجود دارد؛ ۳- هنگامی که sDTF و uDTF هر دو بزرگ است، DTF کلی در مجموع نمره‌های آزمون وجود

1 . Rudner
2 . Getson
3 . Knight
4 . Omnibus test
5 . Item Response Theory

دارد و سوگیری و تفاوت منحنی‌های آزمون نظام‌دار است (چالمرز، کانسِل و فلورا، ۲۰۱۵).

در مورد ارتباط میان DIF و DTF گمانه‌زنی‌هایی وجود دارد مبنی بر اینکه DTF ناشی از تجمع پرسش‌های DIF آزمون است. راجو، وندرلیندن و فلییر^۲ (۱۹۹۵) دو نوع DIF جبرانی و غیر جبرانی را مفهوم‌سازی کرده‌اند که به نظر می‌رسد مجموع پرسش‌های DIF نوع جبرانی آن به DTF در کل آزمون منجر شود. بنابراین حتی اگر تک‌تک پرسش‌های آزمون دارای DIF معنی‌داری نباشند، ممکن است مجموع آن مقدار DTF معنی‌داری را نشان دهد. البته وجود پرسش‌های DIF غیر جبرانی در آزمون ممکن است نوع جبرانی آن را خنثی کرده و آشکارسازی DTF را تحت تأثیر قرار دهد (فلورا، کارن، هوسانگ و ادواردز^۳، ۲۰۰۸). از نقطه‌نظر آماری نیز از آنجایی که کل آزمون نسبت به تک‌تک پرسش‌های آن آگاهی بیشتری فراهم می‌کند، DTF می‌تواند بسیار آسیب‌زننده باشد. از نقطه‌نظر غیر آماری و اجرایی نیز از آنجایی که همه تصمیم‌های سرنوشت‌ساز برای آزمودنی‌ها بر اساس نمره کل آزمون گرفته می‌شود، DTF نقش مخرب‌تری در تصمیم‌گیری‌های نهایی خواهد داشت (راسل^۴، ۲۰۰۵؛ پای و پارک^۵، ۲۰۰۶). در چنین وضعیتی نمره‌های نهایی آزمون، اندازه‌گیری یکسانی در جامعه آزمودنی‌ها ندارند. در این صورت، آزمون دارای «تغییرناپذیری اندازه‌گیری»^۶ نیست و تفاوت میانگین‌ها در گروه‌های مختلف یا همبستگی نمره‌های آزمون با نمره یک آزمون ملاک، گمراه‌کننده خواهد بود (راجو، وندرلیندن و فلییر، ۱۹۹۵). علی‌رغم DTF/DIF که ویژگی منفی برای آزمون است، تغییرناپذیری اندازه‌گیری، ویژگی مثبتی برای آزمون‌های روانی و تربیتی است. این ویژگی آزمون با روایی سازه‌ای آزمون در ارتباط است. علاوه بر این صرف وجود DTF/DIF آزمون به این معنی نیست که آزمون، سوگیری دارد بلکه مراحل تحلیل محتوای پرسش‌ها و آزمون‌های دارای DIF و DTF توسط متخصصان محتوایی و موضوع درسی به تشخیص نهایی پرسش‌ها و آزمون‌های سودار منجر می‌شود. بنابراین وجود DTF/DIF در آزمون، شرط لازم برای

-
1. Chalmers, Counsell & Flora
 2. Raju, Van der Linden & Fleer
 3. Flora, Curran, Hussong & Edwards
 4. Russell
 5. Pae & Park
 6. Measurement invariance

سوگیری بوده ولی کافی نیست (گرامی پور و فلسفی نژاد، ۱۳۹۲). در این رابطه، پژوهش حاضر نیز به منظور مطالعه این گونه روایی در آزمون‌های کنکور سراسری سازمان سنجش آموزش کشور، اجرا می‌شود و بر آن است که پرسش‌ها و آزمون‌هایی را که نسبت به جنسیت دارای سوگیری هستند، مشخص کند.

در مورد روایی آزمون‌های کنکور سراسری با استفاده از روش‌های مبتنی بر نظریه‌های کلاسیک و پیشرفته اندازه‌گیری در داخل کشور مطالعات متعددی انجام گرفته است (معلمی اوره، ۱۳۸۷؛ آبکار، ۱۳۹۱؛ حبیبی، ۱۳۹۲؛ نژادنجف، ۱۳۹۳؛ ترکشوند، ۱۳۹۴ و میری، ۱۳۹۴)، اما درباره روایی آزمون‌های کنکور سراسری ایران از منظر DIF جنسیتی، چند مطالعه بیشتر انجام نگرفته است (براتی و احمدی^۱، ۲۰۱۰؛ امیریان، علوی و فیدالگو^۲، ۲۰۱۴؛ فلاحی سرشت، ۱۳۹۴؛ آسیابی، ۱۳۹۱). مطالعه براتی و احمدی (۲۰۱۰) در مورد DIF مرتبط با جنسیت در آزمون زبان تخصصی انگلیسی کنکور سراسری انجام گرفته است. بررسی آنها نشان داد که حدود ۴۳ درصد پرسش‌های آزمون زبان انگلیسی دارای DIF مرتبط با جنسیت هستند. به علاوه، آنها نشان دادند که فرمت پرسش‌های آزمون و محتوای آن با DIF جنسیتی در ارتباط است. در پژوهش امیریان، علوی و فیدالگو (۲۰۱۴) درباره DIF مرتبط با جنسیت اجرا شده و در آن آزمون زبان انگلیسی دانشگاه تهران، مطالعه شده است. بررسی آنها نشان داد که حدود ۲۸ درصد پرسش‌های آزمون زبان انگلیسی دانشگاه تهران دارای DIF مرتبط با جنسیت هستند، در حالی که اندازه‌های اثر مربوط به DIF بسیار کوچک بودند. مرحله تحلیل محتوا در این مطالعه نشان داد که پرسش‌هایی که دارای مضامین علوم انسانی هستند، نسبت به مردان و پرسش‌های دارای مضامین علوم سخت نسبت به زنان سوگیری دارند. فلاحی سرشت (۱۳۹۴) و آسیابی (۱۳۹۱) نیز DIF مرتبط با جنسیت را به ترتیب در کنکور دکتری و کارشناسی ارشد بررسی کرده‌اند. این در حالی است که تنها یک مطالعه DTF داخلی توسط مینایی (۱۳۹۲) در مورد آزمون علوم پایه هشتم برای مقایسه عملکرد آزمون برای دانش‌آموزان ایرانی و آمریکایی انجام گرفته و تاکنون در ایران هیچ‌گونه مطالعه DTF در مورد آزمون‌های کنکور سراسری انجام نگرفته است.

^۱ . Barati & Ahmadi

^۲ . Amirian, Alavi & Fidalgo

پژوهش‌های خارجی بسیاری نیز در مورد DIF مرتبط با جنسیت در آزمون‌های بزرگ‌مقیاس، اجرا شده اما برخلاف DIF مطالعات کمی در مورد DTF مرتبط با جنسیت انجام گرفته است. این امر می‌تواند به واسطه نو بودن نسبی DTF و روش‌شناسی آن نسبت به روش‌شناسی DIF باشد که به‌تازگی توسعه و تکامل پیدا کرده و کاستی‌های پیشین آن تا حدودی برطرف شده است (چالمرز، کانسل و فلورا، ۲۰۱۵). برای مثال، پای (۲۰۱۱) در مطالعه DTF مرتبط با جنسیت آزمون مهارت‌های تحصیلی زبان انگلیسی پی‌رسون^۱ نشان داد که به‌طور کلی و فارغ از DIF موجود در پرسش‌های آزمون، این آزمون بدون DTF است. دولیتل و کلییری^۲ (۱۹۸۷) در مطالعه DIF آزمون کالج آمریکایی (ACT^۳) نشان دادند که پرسش‌های هندسه و استدلال ریاضی برای جنس مؤنث، دشوارتر و پرسش‌هایی که بیشتر جنبه الگوریتمی و محاسباتی دارند، آسان‌تر می‌نمایند. کارلتون و هریس^۴ (۱۹۹۲) در مطالعه DIF آزمون استعداد تحصیلی (SAT^۵) نشان دادند که پرسش‌های بخش شباهت‌ها و تفاوت‌های آزمون که دارای محتوای روابط انسانی، زیبایی‌شناسی یا فلسفی هستند برای زنان در مقایسه با مردان آسان‌تر هستند. لی^۶ (۲۰۰۶) در مطالعه DIF برنامه سنجش بین‌المللی دانش‌آموزان (PISA^۷) نشان داد که به‌طور کلی پرسش‌های چندگزینه‌ای به سود پسران است. وی همچنین به این نتیجه رسید که DIF در PISA به حوزه‌های محتوایی آن وابسته است. دودین و انابی^۸ (۲۰۰۸) در مطالعه آزمون بین‌المللی علوم و ریاضیات (TIMSS^۹) نشان دادند که به‌طور متوسط حدود ۲۰ درصد پرسش‌های این آزمون در هر دفترچه، دارای DIF مرتبط با جنسیت هستند. آنها همچنین نشان دادند که پرسش‌های DIF در میان دانش‌آموزان دختر و پسر تقریباً به‌طور مساوی تقسیم شده است. آریادوست، گو و کیم^{۱۰} (۲۰۱۱) در مطالعه DIF مجموعه آزمون شنیدن زبان

^۱ . Pearson Test of English Academic

^۲ . Doolittle & Cleary

^۳ . American College Testing

^۴ . Carlton & Harris

^۵ . Scholastic Aptitude Test

^۶ . Le

^۷ . Programme for International Student Assessment

^۸ . Doudeen & Annabi

^۹ . Trends in International Mathematics and Science Study

^{۱۰} . Aryadoust, Goh & Kim

انگلیسی‌میشیگان^۱ دریافتند که بسیاری از پرسش‌های آزمون در سطوح پایین توانایی، به نفع آزمودنی‌های مذکر هستند. بریلند، لی، نجاریان و موراکی^۲ (۲۰۰۴) در مطالعه DIF آزمون انگلیسی به‌عنوان زبان دوم (TOEFL^۳) نشان دادند که پرسش‌های تکلیف نوشتن آزمون TOEFL به نفع آزمودنی‌های مؤنث عمل می‌کنند. پارک^۴ (۲۰۰۸) در پژوهشی DIF پرسش‌های آزمون شنیدن زبان انگلیسی را بررسی کرد. او به این نتیجه رسید که چهار متغیر در ساخت آزمون شامل محتوای متن شنیداری، ارائه عکس در پرسش ارائه شده، نوع زبان به‌کار گرفته شده و فرمت پرسش با DIF جنسیتی در ارتباط هستند. برای مثال، او نشان داد که پرسش‌هایی که محتوای خرید دارند از زنان و پرسش‌هایی که محتوای ورزشی دارند از مردان حمایت می‌کنند. وانگ و لین^۵ (۱۹۹۶) DIF مرتبط با جنسیت را در سنجش عملکرد ریاضی بررسی کردند. آنها نشان دادند تکالیفی که در آن آزمودنی با یک مسئله کاربردی مواجه می‌شود، به نفع جنس مؤنث عمل می‌کنند. براون و کانیونگو^۶ (۲۰۰۷) وجود DIF را در آزمون‌های ریاضی بزرگ‌مقیاس کشور ترینیداد و توباگو بررسی کردند. آنها نشان دادند که حدود ۱۷ درصد از پرسش‌های آزمون به‌صورت معنی‌داری دارای DIF مرتبط با جنسیت هستند، اما شدت DIF ناچیز بود. بربراکلو^۷ (۱۹۹۵) در مطالعه‌ای تقریباً شبیه به مطالعه حاضر در کشور ترکیه، DIF مرتبط با جنسیت را در برخی خرده‌آزمون‌های آزمون ورودی دانشگاه شامل حساب، لغات و هندسه مطالعه کرد. او نشان داد که تقریباً در همه مقایسه‌ها، پرسش‌های حساب به نفع مردان و پرسش‌های لغات و هندسه به نفع زنان است. این نتیجه به‌وضوح نشان می‌دهد که جنس مؤنث در توانایی کلامی و فضایی و جنس مذکر در مهارت‌های حساب برتری دارند. از سوی دیگر، این نتیجه با یافته‌های

1. Michigan English Language Assessment Battery listening test

2. Breland, Lee, Najarian & Muraki

3. Test of English as a Foreign Language

4. Park

5. Wang & Lane

6. Brown & Kanyongo

7. Berberoglu

بسیاری از پژوهشگران تناقض دارد (فنما^۱، ۱۹۸۰؛ فنما و کارپنتر^۲، ۱۹۸۱؛ پاتیسون^۳ و گریو^۴، ۱۹۸۴؛ وود^۵، ۱۹۷۶) که در آنها جنس مذکر در توانایی فضایی برتری داشته، در حالی که جنس مؤنث در مهارت‌های حساب مثل مجموعه‌های جبری برتری خود را نشان می‌دهد. برخی پژوهشگران مثل هوسن^۶ (۱۹۶۷)، هانا^۷، ۱۹۸۹ و اتینگتون^۸ (۱۹۹۰) معتقدند که یافته‌های متناقض DIF در مورد تفاوت‌های محتوای سوگیرانه نسبت به جنس مؤنث و مذکر ممکن است به واسطه اثر کشوری باشد که مطالعه در آن انجام می‌شود. بنابراین با توجه به نتایج متفاوت و در برخی موارد متناقض مطالعات DIF و DTF مرتبط با جنسیت و ضرورت بررسی بیشتر در این حوزه و نقصان چنین مطالعه‌ای در آزمون‌های کنکور سراسری سازمان سنجش آموزش کشور، نتایج این پژوهش می‌تواند برای متخصصان این حوزه بسیار راهگشا بوده و برای تصمیم‌گیران و آزمون‌سازان داخل کشور بسیار مفید به فایده باشد. با توجه به این موارد و نقش حیاتی که کنکور سراسری در زندگی دانش‌آموزان و ادامه تحصیل آنها در دانشگاه دارد، کمترین مقدار DTF/DIF در آزمون‌های سرنوشت‌ساز می‌تواند پیامدهای ناخواسته گسترده‌ای در جامعه به دنبال داشته باشد (پای و پارک، ۲۰۰۶). بنابراین پرسش اصلی مطالعه حاضر این است که آیا آزمون‌های کنکور سراسری سازمان سنجش آموزش کشور نسبت به گروه‌های جنسیتی دارای کنش افتراقی سؤال و آزمون (DTF/DIF) و به دنبال آن دارای سوگیری هستند؟

روش

پژوهش حاضر از نظر هدف، کاربردی و از حیث روش، از نوع ارزشیابی است. جامعه و نمونه تحقیق شامل همه آزمودنی‌های یک دفترچه آزمون و درس‌های امتحانی تخصصی تمام گروه‌های آزمایشی کنکور سراسری از سال ۱۳۸۷ تا ۱۳۹۰ بود. از میان درس‌ها برای هر گروه آزمایشی درس‌های تخصصی انتخاب شدند. درس‌های انتخاب

-
1. Fennema
 2. Carpenter
 3. Pattison
 4. Grieve
 5. Wood
 6. Husen
 7. Hanna
 8. Ethington

شده شامل: درس ریاضی (۵۰ پرسش)، فیزیک (۴۵ پرسش) و شیمی (۳۵ پرسش) از گروه ریاضی و فنی؛ زیست‌شناسی (۵۰ پرسش) و شیمی (۳۵ پرسش) از گروه علوم تجربی؛ ادبیات فارسی (۳۰ پرسش) از گروه علوم انسانی؛ درک عمومی هنر (۳۰ پرسش)، خلاقیت تصویری و تجسمی (۲۰ پرسش) و خلاقیت نمایشی (۲۰ پرسش) از گروه هنر و زبان انگلیسی تخصصی (۷۰ پرسش) از گروه زبان بودند. خلاصه اطلاعات حجم نمونه آزمودنی‌ها بر اساس جنسیت به تفکیک سال‌های برگزاری کنکور سراسری در جدول (۱) ملاحظه می‌شود.

جدول (۱) حجم نمونه مورد مطالعه برای DTF/DIF به تفکیک سال‌های برگزاری کنکور

سراسری و جنسیت

جنسیت		سال و رشته
مؤنث	مذکر	
۱۰۳۴۵	۱۲۶۳۴	ریاضی
۲۱۶۶۷	۱۱۶۷۵	تجربی
۲۶۷۷۵	۱۳۴۴۵	انسانی
۲۹۷۷	۱۱۱۲	هنر
۷۳۲۴	۳۱۱۲	زبان
۱۰۵۵۳	۱۳۴۶۶	ریاضی
۲۵۴۴۲	۹۹۶۳	تجربی
۲۴۵۵۷	۱۳۶۷۵	انسانی
۳۱۱۶	۱۰۴۸	هنر
۶۲۴۶	۳۱۰۴	زبان
۱۱۲۹۹	۱۳۴۶۹	ریاضی
۲۴۷۴۰	۱۱۲۷۴	تجربی
۲۳۴۴۶	۱۳۸۹۹	انسانی
۲۳۸۹	۱۲۸۶	هنر
۷۸۹۹	۳۸۷۷	زبان
۱۰۱۳۵۶	۱۳۱۶۷۷	ریاضی
۳۱۲۳۵۳	۸۹۸۴۶	تجربی
۱۹۹۵۶۶	۱۱۴۹۰۰	انسانی
۲۴۷۲۴	۱۰۹۹۱	هنر
۶۷۵۴۴	۲۹۸۵۵	زبان

برای تشخیص کنش افتراقی پرسش‌های آزمون (DIF) از رویکرد رگرسیون لجستیک دو وجهی^۱ و بسته‌های نرم‌افزاری difR در نرم‌افزار R و SPSS استفاده شد. رویه‌های محاسباتی در این رویکرد شامل سه گام در قالب سه معادله بود که در آن معادله‌ها، پاسخ (درست یا غلط) به سؤال متغیر تابع است: ۱- در معادله اول متغیر پیش‌بین (نمره کل آزمون) وارد معادله شد؛ ۲- در معادله دوم نشانگر جنسیت وارد معادله شد و ۳- در معادله سوم، متغیر تعامل دو متغیر وارد شده قبلی وارد معادله رگرسیون لجستیک شد. سپس با استفاده از لگاریتم نسبت درست‌نمایی در آزمونی با $df=2$ ، مدل مرحله ۱ و ۲ و مدل مرحله ۲ و ۳ مقایسه شد. این آزمون، دو آزمون دارای $df=1$ را ترکیب می‌کند. آزمون اول، مدل مرحله ۱ را با مدل مرحله ۲ مقایسه می‌کند که در این آزمون DIF هماهنگ بررسی شد. رد فرض صفر به معنی وجود DIF هماهنگ است. آزمون دوم نیز با استفاده از لگاریتم نسبت درست‌نمایی، مدل مرحله ۲ را با مدل مرحله ۳ مقایسه می‌کند. این مقایسه، DIF ناهماهنگ را ارزیابی می‌کند (گرامی‌پور، ۱۳۹۳). در نرم‌افزار R برای بررسی معنی‌داری DIF، صرفاً مقادیر اختلاف لگاریتم‌های درست‌نمایی، سطح معنی‌داری و اندازه‌های اثر DIF فراهم شد. همچنین برای تشخیص کنش افتراقی آزمون (DTF) از رویکرد مبتنی بر نظریه IRT چالمرز، کانسل و فلورا (۲۰۱۵) و بسته نرم‌افزاری mirt در نرم‌افزار R استفاده شد. به‌علاوه، یکی از پیش‌فرض‌های تحلیل داده‌ها پیش از آشکارسازی DTF/DIF، تک‌بعدی بودن داده‌ها است (فریزر و مک‌دونالد^۲، ۱۹۸۸). برای آزمون تک‌بعدی بودن داده‌ها از نرم‌افزار NOHARM استفاده شد و برای تأیید آن، سه شاخص مجموع مجذورات باقی‌مانده‌ها، ریشه دوم مجذور میانگین باقی‌مانده‌ها و شاخص تاناکا و تاناکا^۳ بنا بر پیشنهاد متخصصان (فینچ و هابینگ^۴، ۲۰۰۷) مورد استناد قرار گرفت.

یافته‌ها

نخست، نتایج سه شاخص مورد استناد در NOHARM نشان‌دهنده تک‌بعدی بودن همه آزمون‌های مورد مطالعه بود، بدین ترتیب که: ۱- مجموع مجذورات باقی‌مانده برای همه آزمون‌های مورد مطالعه کوچک‌تر یا مساوی با ۰/۰۰۵ بودند؛ ۲- مجذور

1. Binary

2. Fraser & McDonald

3. Tanaka & Tanaka

4. Finch & Habing

میانگین ریشه باقی مانده برای همه آزمون‌ها کوچک‌تر مساوی با ۰/۰۱ بودند و ۳- مقدار شاخص تاناکا و تاناکا که یک شاخص بیشینه درست‌نمایی است، برای همه آزمون‌ها بزرگ‌تر مساوی با ۰/۹۵ بود.

در ادامه، یافته‌های رویکرد رگرسیون لجستیک برای آشکارسازی DIF جنسیتی برای یک نمونه پرسش آزمون درس زیست‌شناسی در گروه علوم تجربی سال ۱۳۹۰ ارائه می‌شود.

پرسش ۳، درس زیست‌شناسی از گروه علوم تجربی سال ۱۳۹۰ به‌طور معمول در باکتری‌هایی که کروموزوم‌های کمی دارند، به تعداد مولکول‌های DNA وجود دارد.

۱- جایگاه شروع همانندسازی ۲- ژن مقاومت نسبت به آنتی‌بیوتیک ۳- دوراهی همانندسازی ۴- جایگاه تشخیص آنزیم محدودکننده

همان‌طور که در جدول (۲) ملاحظه می‌شود، نتایج تحلیل رگرسیون لجستیک پرسش ۳ آزمون درس زیست‌شناسی در گروه تجربی سال ۹۰، نشان می‌دهد که جنس مؤنث احتمال بیشتری برای دادن پاسخ درست به پرسش ۳ آزمون را دارد ($b = 0/176, P < 0/01$). همچنین مقدار ثابت معادله رگرسیون لجستیک ($\alpha = -3/135, P < 0/01$) و نمره کل آزمون زیست‌شناسی ($b = 0/083, P < 0/01$) با بیش از ۹۹ درصد اطمینان معنی‌دار است. نمره کل آزمون از جمع نمره‌های ۵۰ سؤال آزمون، تشکیل شده است و در معادله رگرسیون لجستیک تأثیر آن بر احتمال پاسخ به پرسش کنترل می‌شود. مقدار ثابت نیز در معادله رگرسیون لجستیک نشان‌دهنده برآورد خام احتمال پاسخ به پرسش است.

جدول (۲) مدل مرحله ۲ رگرسیون لجستیک

منبع	آماره	مقدار برآورد	خطای استاندارد برآورد	مقدار والد	درجه آزادی	سطح معنی‌داری	توان نمایی برآورد
جنسیت	۰/۱۷۶	۰/۱۱	۰/۰۱۱	۲۷۵/۹۳۰	۱	۰/۰۰۰	۱/۱۹۲
نمره کل زیست	۰/۰۸۳	۰/۰۰۱	۰/۰۰۱	۲۴۹۶۱/۵۲۸	۱	۰/۰۰۰	۱/۰۸۷
مقدار ثابت	-۳/۱۳۵	۰/۰۱۷	۰/۰۱۷	۳۵۶۶۸/۶۰۸	۱	۰/۰۰۰	۰/۰۴۴

معادله مدل مرحله دوم رگرسیون لجستیک به شرح زیر خواهد بود:
 مدل کامل لوجیت: $P[\text{پرسش ۳ آزمون} = \text{پاسخ درست}] = [-3/135 - ++ \text{جنسیت } 0/176]$
 نمره کل آزمون ۰/۰۸۳

جدول (۳) نشان‌دهنده درصد پیش‌بینی پاسخ‌ها با مدل رگرسیون لجستیک است. هر چقدر درصد پیش‌بینی پاسخ‌ها برای پرسش ۳ آزمون بیشتر باشد مدل رگرسیون لجستیک قدرت بیشتری دارد. جدول (۳) نشان می‌دهد که مدل رگرسیون لجستیک توانسته است حدود ۸۷ درصد پاسخ‌ها برای پرسش ۳ آزمون را پیش‌بینی کند.

جدول (۳) طبقه‌بندی پیش‌بینی‌های درست رگرسیون لجستیک در مدل مرحله ۲

پیش‌بینی شده		مشاهده شده	
درصد درست	پرسش ۳		پرسش ۳
	پاسخ درست	پاسخ غلط	
۹۹/۲	۲۷۱۱	۳۳۹۵۶۰	پاسخ غلط
۸/۱	۴۱۲۱	۴۶۷۸۳	پاسخ درست
۸۷/۴			درصد کل پیش‌بینی درست

همچنین جدول (۴) خلاصه مدل مقدار واریانس پیش‌بینی شده مدل توسط متغیر جنسیت و نمره کل آزمون را نشان می‌دهد.

جدول (۴) خلاصه مدل رگرسیون لجستیک برای پیش‌بینی پاسخ پرسش ۳ آزمون در مرحله ۲

مقدار ۲- برابر لگاریتم بخت	ضریب تبیین کاکس و اسنل	ضریب تبیین ناگلکرک
۲۷۸۰۰۲/۵۵۴	۰/۰۶۲	۰/۱۱۵

نتایج جدول (۴) نشان می‌دهد که حدود ۱۱ درصد احتمال پاسخ برای پرسش ۳ آزمون درس زیست‌شناسی توسط متغیر جنسیت و نمره کل آزمون زیست‌شناسی پیش‌بینی می‌شود. تفاوت ضرایب تبیین ملاحظه شده در جدول (۴) در این است که ضریب تبیین کاکس و اسنل ضریب تبیین نمونه و ضریب ناگلکرک مقدار تبیین شده آن برای جامعه است.

معادله محاسبه شده برای مدل کاهش یافته (مدل مرحله اول) در مثال حاضر عبارت است از:
 مدل کاهش یافته: $P = [P(3 \text{ پرسش} = \text{آزمون} = \text{پاسخ درست}) - 2/904] +$
 نمره کل آزمون $0/083$

همان‌طور که ملاحظه می‌شود در مدل کاهش یافته (مدل مرحله ۱) متغیر جنسیت از معادله رگرسیون لجستیک حذف شده و محاسبه دوباره انجام شده است. مقدار (۲-) برابر لگاریتم بخت برای معادله کاهش یافته معادل با $78274/977$ است که اگر مقدار (۲-) برابر لگاریتم بخت مدل کامل از آن کم شود مقدار $272/423$ به دست می‌آید، این مقدار با ۲ درجه آزادی نشان‌دهنده DIF هماهنگ پرسش بررسی شده در سطح معنی‌داری کوچک‌تر از یک‌صدم معنی‌دار است. بنابراین پرسش ۳ درس زیست‌شناسی از گروه علوم تجربی سال ۹۰ نسبت به جنسیت مرد سوگیری آماری دارد.

جدول (۵) مدل مرحله ۳ رگرسیون لجستیک

منبع	آماره	مقدار برآورد	خطای استاندارد برآورد	مقدار والد	درجه آزادی	سطح معنی داری	توان نمایی برآورد
جنسیت	۰/۲۱۳	۰/۰۱۸	۱۳۶/۲۵۵	۱	۰/۰۰۰	۱/۲۳۸	
نمره کل آزمون زیست	۰/۰۸۷	۰/۰۰۲	۳۰۲۷/۴۷۹	۱	۰/۰۰۰	۱/۰۹۱	
جنسیت ×	-۰/۰۰۳	۰/۰۰۱	۶/۴۰۴	۱	۰/۰۱۱	۰/۹۹۷	
نمره کل آزمون مقدار ثابت	-۳/۱۸۶	۰/۰۲۶	۱۴۷۷۱/۷۲	۱	۰/۰۰۰	۰/۰۴۱	

از سوی دیگر، همان‌طور که در جدول (۵) ملاحظه می‌شود، مدل مرحله سوم با در نظر گرفتن متغیر تعامل، نتایج تحلیل رگرسیون لجستیک نشان می‌دهد که جنسیت احتمال بیشتری برای دادن پاسخ درست به پرسش ۳ آزمون را دارد ($p < 0/01, b = 0/213$). مقدار ثابت معادله رگرسیون لجستیک ($\alpha = -3/186, P < 0/01$) و نمره کل آزمون ($b = 0/087, P < 0/01$) با بیش از ۹۹ درصد اطمینان، معنی دار است. معادله رگرسیون لجستیک به شرح زیر خواهد بود:

کامل مدل (لوجیت): $P = [0/087 \times \text{نمره کل آزمون} - 3/186 + \text{جنسیت} \times 0/213] = \text{پاسخ درست}$

برای DIF ناهماهنگ معادله فوق به این معنی است که لوجیت احتمال پاسخ درست به پرسش ۳ آزمون عبارت است از تأثیر تعامل جنسیت و نمره کل آزمون بر لگاریتم بخت پاسخ درست به پرسش ۳ آزمون در حالی که تأثیر جنسیت و نمره کل آزمون بر احتمال پاسخ درست به پرسش کنترل شده است.

جدول (۶) نشان‌دهنده درصد پیش‌بینی پاسخ‌ها با مدل رگرسیون لجستیک است. هر چقدر درصد پیش‌بینی پاسخ‌ها برای پرسش ۳ آزمون بیشتر باشد، مدل رگرسیون لجستیک قدرت بیشتری دارد. جدول (۶) نشان می‌دهد که مدل رگرسیون لجستیک توانسته است حدود ۸۷ درصد پاسخ‌ها برای پرسش ۳ درس زیست‌شناسی را پیش‌بینی کند.

جدول (۶) طبقه‌بندی پیش‌بینی‌های درست مدل رگرسیون لجستیک مرحله سوم

درصد درست	پیش‌بینی شده		مشاهده شده	
	پرسش ۳			
	پاسخ درست	پاسخ غلط	پاسخ غلط	پرسش ۳
۹۹/۲	۲۷۱۱	۳۳۹۵۶۰	پاسخ غلط	پرسش ۳
۸/۱	۴۱۲۱	۴۶۷۸۳	پاسخ درست	
	۸۷/۴		درصد کل پیش‌بینی درست	

همچنین جدول (۷) خلاصه مدل مقدار واریانس پیش‌بینی شده مدل توسط متغیر جنسیت و نمره کل آزمون درس ریاضی و تعامل جنسیت با نمره کل آزمون درس زیست‌شناسی را نشان می‌دهد.

جدول (۷) خلاصه مدل رگرسیون لجستیک مرحله سوم برای پیش‌بینی پاسخ پرسش ۳ درس زیست‌شناسی

مقدار ۲- برابر لگاریتم بخت	ضریب تبیین کاکس و اسنل	ضریب تبیین ناگلکرک
۲۷۷۹۹۹/۱۵۶	۰/۰۶۲	۰/۱۱۵

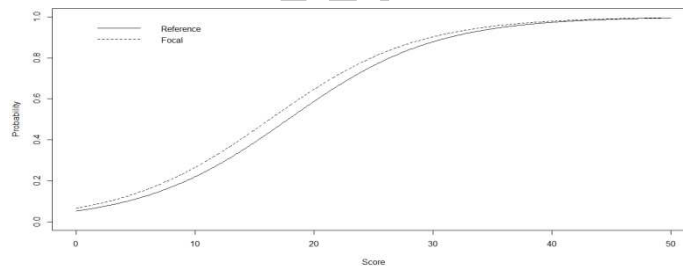
نتایج جدول (۷) نشان می‌دهد که در مدل مرحله سوم نیز حدود ۱۱ درصد احتمال پاسخ برای پرسش ۳ آزمون درس زیست‌شناسی توسط متغیر جنسیت و نمره کل درس زیست‌شناسی و تعامل این دو پیش‌بینی می‌شود.

معادله محاسبه شده برای مدل کاهش یافته DIF ناهماهنگ در پرسش ۳ درس زیست‌شناسی کنکور تجربی سال ۱۳۹۰ حاضر عبارت است از:

مدل کاهش یافته (مرحله ۲): لوجیت $P[\text{پرسش ۳ آزمون} = \text{پاسخ درست}] = -3/135 - \text{جنسیت} + 0/176 \times \text{نمره کل آزمون} + 0/083$

همان‌طور که ملاحظه می‌شود در مدل کاهش یافته DIF ناهماهنگ متغیر تعامل جنسیت با نمره کل آزمون از معادله رگرسیون لجستیک حذف شده و محاسبه دوباره انجام شده است. مقدار (۲-) برابر لگاریتم بخت برای معادله کاهش یافته معادل با $278002/554$ است که اگر مقدار (۲-) برابر لگاریتم بخت مدل کامل از آن کم شود مقدار $3/398$ به دست می‌آید. این مقدار با ۲ درجه آزادی، مقدار معنی‌داری برای نشان دادن DIF ناهماهنگ پرسش ۳ آزمون نیست.

نمودار (۱) نشان‌دهنده منحنی‌های ویژگی پرسش (ICC) مورد مطالعه به تفکیک جنسیت است.



نمودار (۱) منحنی‌های ویژگی پرسش (ICC) مورد مطالعه به تفکیک جنسیت برای پرسش ۳ آزمون زیست‌شناسی سال ۱۳۹۰

¹ . Item Characteristic Curve

همان‌طور که در نمودار (۱) ملاحظه می‌شود گروه کانونی^۱ (جنسیت مؤنث) نسبت به گروه مرجع^۲ تا حدود زیادی در سراسر نمره‌های توانایی (به جز سطوح بالای توانایی) احتمال بیشتری برای پاسخ درست به پرسش مورد بررسی را دارند. همچنین، اندازه اثر DIF مبتنی بر R^2 ناگلکرک برای این پرسش کمتر از ۰/۰۰۰۱ بود که بسیار ناچیز و قابل اغماض تلقی می‌شود.

سایر نتایج تحلیل DIF مبتنی بر رویکرد تحلیل رگرسیون هیچ‌گونه DIF ناهماهنگ معنی‌داری را بر اساس آماره آزمون نسبت درست‌نمایی در پرسش‌های آزمون‌های کنکور سراسری نشان نداد ($LRT^3 < 1, P > 0/05$)، اما برخی از پرسش‌های آزمون‌های سراسری، بر اساس آماره آزمون نسبت درست‌نمایی، دارای DIF هماهنگ معنی‌داری بودند ($LRT > 4, P < 0/05$) که نرخ آن در آزمون‌های درس‌های مورد مطالعه برحسب سال‌های برگزاری آزمون (در سطح معنی‌داری ۰/۰۵) در جدول (۸) ملاحظه می‌شود.

جدول (۸) تعداد و درصد پرسش‌های DIF در آزمون‌های درس‌های مورد مطالعه برحسب سال‌های برگزاری آزمون

درصد متوسط در گروه	متوسط آزمون		۱۳۹۰		۱۳۸۹		۱۳۸۸		۱۳۸۷		سال‌های برگزاری آزمون	
	درصد	تعداد	درصد	تعداد	درصد	تعداد	درصد	تعداد	درصد	تعداد	گروه داوطلبی/درس	
۱۴/۶۹	۱۴/۰۰	۷	۱۶/۰۰	۸	۱۰/۰۰	۵	۱۲/۰۰	۶	۱۸/۰۰	۹	ریاضی	ریاضی و فیزیک
	۱۷/۷۸	۸	۲۲/۲۲	۱۰	۱۷/۷۸	۸	۲۰	۹	۱۱/۱۱	۵	فیزیک	
	۱۳/۵۷	۴/۷۵	۱۷/۱۴	۶	۱۱/۴۳	۴	۱۴/۲۹	۵	۱۱/۴۳	۴	شیمی	
۱۶/۱۸	۱۴/۵	۷/۲۵	۱۶	۸	۱۶	۸	۱۴	۷	۱۲	۶	زیست‌شناسی	علوم تجربی
	۱۷/۸۶	۶/۲۵	۱۴/۲۹	۵	۲۰	۷	۱۷/۱۴	۶	۲۰	۷	شیمی	
۱۵/۸۳	۱۵/۸۳	۴/۷۵	۲۰	۶	۱۶/۶۷	۵	۱۳/۳۳	۴	۱۳/۳۳	۴	ادبیات فارسی	علوم انسانی
۱۲/۹۲	۷/۵	۲/۲۵	۲۰	۶	۰	۰	۰	۰	۱۰	۳	درک عمومی هنر	هنر

1. Focal
2. References
3. Item Characteristic

درصد	متوسط آزمون		۱۳۹۰		۱۳۸۹		۱۳۸۸		۱۳۸۷		سال‌های برگزاری آزمون	گروه داوطلبی/درس
	درصد	تعداد	درصد	تعداد	درصد	تعداد	درصد	تعداد	درصد	تعداد		
	۷/۵	۱/۵	۰	۰	۲۰	۴	۰	۰	۱۰	۲	خلاقیت تصویری و تجسمی	
	۲۳/۷۵	۴/۷۵	۳۵	۷	۲۵	۵	۱۵	۳	۲۰	۴	خلاقیت نمایشی	
۱۲/۱۴	۱۲/۱۴	۸/۵	۱۰	۷	۱۱/۴۳	۸	۱۴/۲۹	۱۰	۱۲/۸۶	۹	زبان انگلیسی تخصصی	زبان
	۱۴/۳۲		۱۸/۸۲		۱۲/۷۴		۱۲/۹۰		۱۲/۷۱		درصد متوسط برحسب سال	

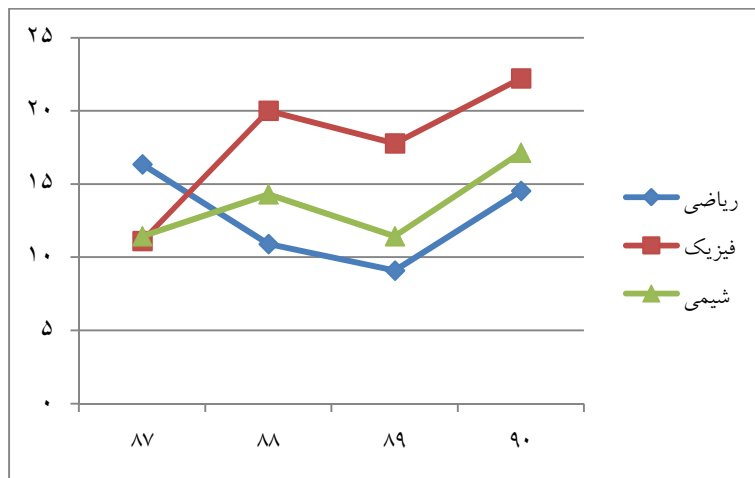
جدول (۹) نیز تعداد و درصد سال‌های DIF علیه جنس مذکر یا مؤنث را نشان می‌دهد.

جدول (۹) تعداد و درصد سال‌های DIF آزمون علیه جنسیت مذکر و مؤنث برحسب سال‌های برگزاری آزمون

متوسط آزمون	۱۳۹۰		۱۳۸۹		۱۳۸۸		۱۳۸۷		سال‌های برگزاری آزمون	گروه داوطلبی/درس/جنسیت	
	درصد	تعداد	درصد	تعداد	درصد	تعداد	درصد	تعداد			
۲/۵۰	۱/۲۵	۶	۳	۰	۰	۴	۲	۰	۰	مذکر	ریاضی و فیزیک
۱۱/۵۰	۵/۷۵	۱۰	۵	۱۰	۵	۸	۴	۱۸	۹	مؤنث	
۶/۱۱	۲/۷۵	۱۳/۳۳	۶	۰	۰	۰	۱۱/۱۱	۵	۵	مذکر	فیزیک
۱۱/۶۷	۵/۲۵	۸/۸۹	۴	۱۷/۷۸	۸	۲۰	۹	۰	۰	مؤنث	
۷/۸۶	۲/۷۵	۲/۸۶	۱	۱۱/۴۳	۴	۵/۷۱	۲	۱۱/۴۳	۴	مذکر	شیمی
۵/۷۲	۲	۱۴/۲۹	۵	۰	۰	۸/۵۷	۳	۰	۰	مؤنث	
۷/۵۰	۳/۷۵	۱۰	۵	۸	۴	۸	۴	۴	۲	مذکر	علوم تجربی
۷	۳/۵۰	۶	۳	۸	۴	۶	۳	۸	۴	مؤنث	
۱۳/۵۷	۴/۷۵	۲/۸۶	۱	۱۴/۲۹	۵	۱۷/۱۴	۶	۲۰	۷	مذکر	شیمی
۴/۲۹	۱/۵۰	۱۱/۴۳	۴	۵/۷۱	۲	۰	۰	۰	۰	مؤنث	
۱۲/۵۰	۳/۷۵	۱۰	۳	۱۳/۳۳	۴	۱۳/۳۳	۴	۱۳/۳۳	۴	مذکر	علوم انسانی
۳/۳۳	۱	۱۰	۳	۳/۳۳	۱	۰	۰	۰	۰	مؤنث	
۳/۳۸	۱	۶/۶۷	۲	۰	۰	۰	۰	۶/۶۷	۲	مذکر	هنر
۴/۱۷	۱/۲۵	۱۳/۳۳	۴	۰	۰	۰	۰	۳/۳۳	۱	مؤنث	

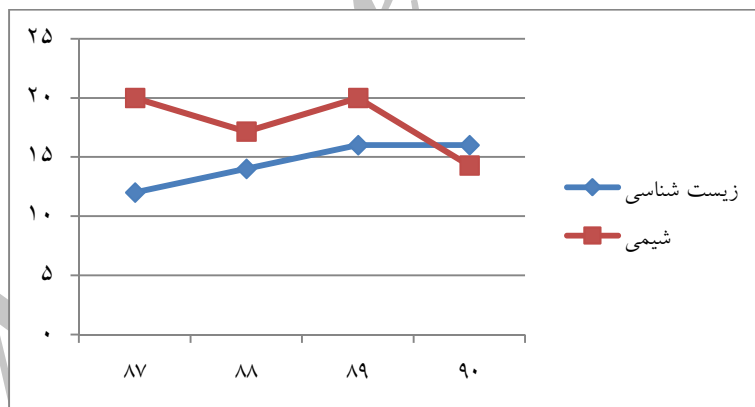
متوسط آزمون		۱۳۹۰		۱۳۸۹		۱۳۸۸		۱۳۸۷		سال‌های برگزاری آزمون	
درصد	تعداد	درصد	تعداد	درصد	تعداد	درصد	تعداد	درصد	تعداد	گروه داوطلبی/درس/جنسیت	
۲/۵۰	۰/۵۰	۰	۰	۱۰	۲	۰	۰	۰	۰	مذکر	خلاقیات
۵	۱	۰	۰	۱۰	۲	۰	۰	۱۰	۲	مؤنث	تصویری و تجسمی
۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	مذکر	خلاقیات
۲۳/۷۵	۴/۷۵	۳۵	۷	۲۵	۵	۱۵	۳	۲۰	۴	مؤنث	نمایشی
۷/۸۶	۵/۵۰	۵/۷۱	۴	۴/۲۹	۳	۸/۵۷	۶	۱۲/۸۶	۹	مذکر	زبان انگلیسی
۴/۲۹	۳	۴/۲۹	۳	۷/۱۴	۵	۵/۷۱	۴	۰	۰	مؤنث	تخصصی
											زبان

همان‌طور که در جدول‌های (۸ و ۹) ملاحظه می‌شود در گروه ریاضی و فیزیک، نرخ پرسش‌های DIF علیه جنس مؤنث در سال‌های ۸۷ و ۹۰ نسبت به سال‌های ۸۸ و ۸۹ بیشتر است. این در حالی است که نرخ DIF علیه جنس مؤنث در درس ریاضی مشهود است. نرخ پرسش‌های DIF در آزمون فیزیک سال ۸۷ نسبت به سایر سال‌های برگزاری آزمون کمتر است، در حالی که پرسش‌های DIF علیه جنس مذکر کمی بیشتر است. به‌علاوه، آزمون فیزیک سال ۹۰ به نسبت دارای بیشترین درصد پرسش‌های DIF علیه هر دو جنس مؤنث و مذکر است، اما به‌طور کلی نرخ پرسش‌های DIF علیه جنس مؤنث کمی بیشتر است. اما در آزمون شیمی گروه ریاضی و فیزیک، پرسش‌های DIF کمی بیشتر علیه جنس مذکر است. به‌طور متوسط می‌توان مشاهده کرد که آزمون فیزیک در این گروه در مقایسه با سایر درس‌های مورد مطالعه دارای پرسش‌های DIF بیشتری است. روند پرسش‌های DIF در آزمون‌های مورد مطالعه گروه ریاضی و فیزیک در نمودار (۲) ملاحظه می‌شود.



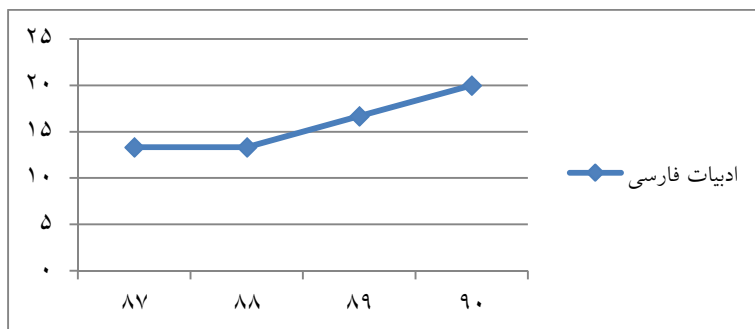
نمودار (۲) درصد پرسش‌های دارای DIF در آزمون‌های مورد مطالعه گروه ریاضی و فیزیک

در گروه علوم تجربی، برخی پرسش‌های آزمون زیست‌شناسی نسبت جنس مؤنث و برخی نسبت به جنس مذکر دارای DIF است. درس شیمی در این گروه نیز شبیه به گروه ریاضی و فیزیک دارد، اما نرخ پرسش‌های DIF در آزمون شیمی به‌طور متوسط از درس زیست‌شناسی بیشتر است. نمودار (۳) نشان‌دهنده نرخ پرسش‌های DIF در درس‌های مورد مطالعه گروه علوم تجربی است.



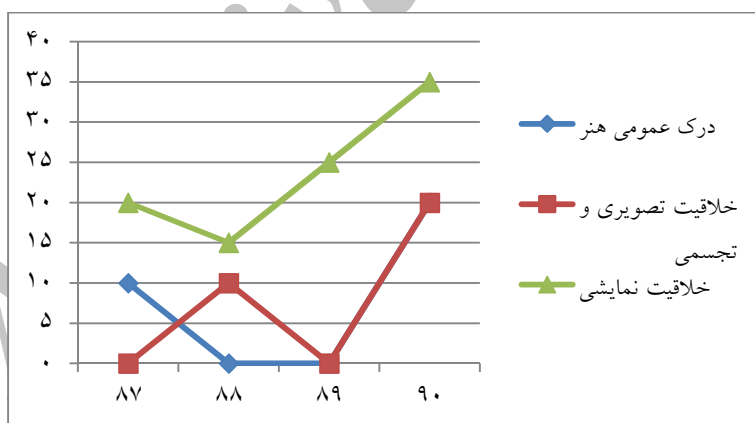
نمودار (۳) درصد پرسش‌های DIF در درس‌های مورد مطالعه گروه علوم تجربی

تحلیل آزمون ادبیات فارسی گروه علوم انسانی نشان می‌دهد که پرسش‌های DIF به نفع جنس مؤنث عمل می‌کنند، در حالی که نرخ آنها از سال ۸۷ تا ۹۰ رو به افزایش است. چنین روندی در نمودار (۴) ملاحظه می‌شود.



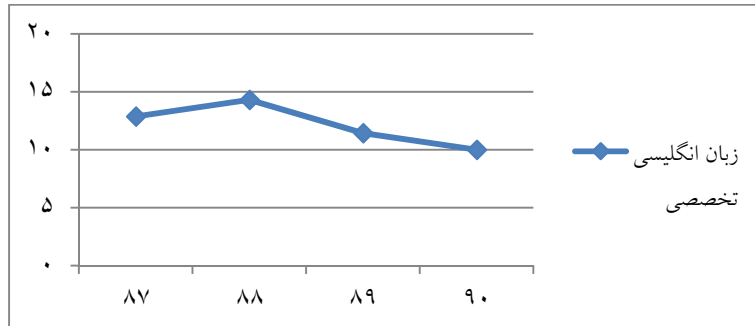
نمودار (۴) درصد پرسش‌های DIF در آزمون ادبیات فارسی گروه علوم انسانی

نمودار (۴) نشان می‌دهد که بیشترین نوسان در گروه هنر است، بدین ترتیب که کمترین نرخ پرسش‌های DIF در درس‌های درک عمومی هنر و خلاقیت تصویری و تجسمی (۷/۵۰ درصد) وجود دارد و در عین حال بیشترین نرخ پرسش‌های DIF (۲۳/۷۵ درصد) نیز در همین گروه و در آزمون درس خلاقیت نمایشی است. در این درس به‌طور مشهودی پرسش‌های DIF علیه جنس مؤنث عمل می‌کنند.



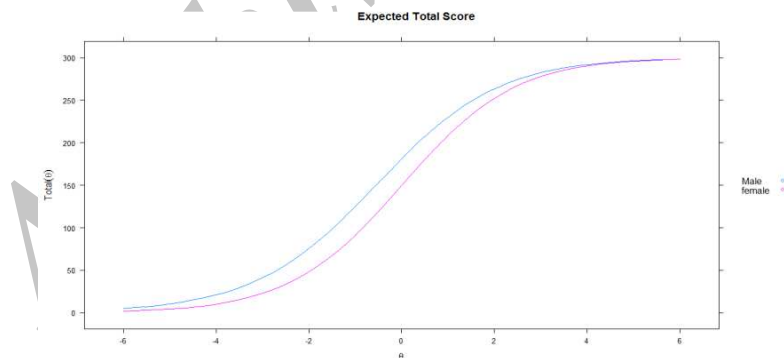
نمودار (۵) نرخ سؤال‌های دارای DIF در آزمون‌های تخصصی گروه هنر

سرانجام نمودار (۶) نشان می‌دهد که آزمون زبان تخصصی گروه زبان دارای کمترین نرخ پرسش‌های DIF است، در حالی که نرخ پرسش‌های DIF از سال ۸۸ تا ۹۰ سیر نزولی داشته است.



نمودار (۶) نرخ پرسش‌های DIF در آزمون زبان تخصصی گروه زبان

همچنین در جدول (۹) می‌توان دید که در درس زبان تخصصی گروه زبان، به جز سال ۸۷ که پرسش‌ها علیه جنس مذکر عمل می‌کنند، در سایر سال‌ها برخی پرسش‌های DIF علیه جنس مؤنث و برخی علیه جنس مذکر عمل می‌کنند. سرانجام، نتایج تحلیل DTF نشان داد که به غیر از آزمون درس خلاقیت نمایشی گروه هنر، سایر آزمون‌های مورد مطالعه فاقد DTF جنسیتی هستند، به طوری که مقادیر قدر مطلق sDTF و مقادیر uDTF محاسبه شده برای آنها کمتر از ۰/۰۶ و غیر معنی‌دار بود ($P > 0/05$). نمودار (۷) نشان‌دهنده DTF جنسیتی برای آزمون خلاقیت نمایشی است.



نمودار (۷) DTF مرتبط با جنسیت برای آزمون درس خلاقیت نمایشی گروه هنر

مقادیر محاسبه شده DTF آزمون خلاقیت نمایشی شامل $sDTF = -3/12$ و $sDTF = 3/34$ است. نمودار (۷) نشان می‌دهد که این آزمون حدود ۳ نمره خام از گروه مرجع (جنس مذکر) طرفداری می‌کند. بنابراین فرض صفر نبود DTF مرتبط با جنسیت در آزمون خلاقیت نمایشی رد شد ($P < 0/01$). همچنین همان‌طور که در نمودار (۷) دیده می‌شود، در سطوح بسیار بالای توانایی سوگیری آزمون خلاقیت نمایشی نسبت به جنس مؤنث از بین می‌رود و این سوگیری در توانایی پایین‌تر از آن سطح مشاهده می‌شود. البته مقدار اندازه اثر استاندارد شده برای این مقدار DTF معادل با ۱۶ درصد محاسبه شده و همچنان در حد «کم» ارزیابی می‌شود.

بحث و نتیجه‌گیری

وجود پرسش‌های DIF در آزمون‌های بزرگ‌مقیاس، اجتناب‌ناپذیر است، اما وجود پرسش‌های DIF لزوماً به معنای سوگیری پرسش‌های آزمون نیست (لی، ۲۰۰۶). با این وجود، تحلیل DIF پرسش‌های آزمون سرنوشت‌سازی مانند کنکور سراسری بر اساس جنسیت (به‌عنوان اساسی‌ترین متغیر گروهی) بسیار ضروری است؛ زیرا وجود DIF مرتبط با جنسیت به‌صورت منفی بر منصفانه بودن آزمون برای دختران و پسران تأثیر می‌گذارد و دقت نمره‌ها و کاربرد آزمون برای گزینش مبتنی بر جنسیت را کاهش می‌دهد.

یافته‌های این پژوهش نشان می‌دهد که به‌طور کلی حدود ۱۴ درصد پرسش‌های آزمون‌های تخصصی مورد مطالعه در کنکور سراسری سال‌های ۸۷ تا ۹۰ دارای DIF هستند. البته شدت DIF محاسبه شده برای این پرسش‌ها بسیار ناچیز هستند. نرخ کلی آشکارسازی پرسش‌های DIF در آزمون تقریباً با یافته‌های دودین و انابی (۲۰۰۸) و براون و کانینگو (۲۰۰۷) همخوانی دارد. همچنین به غیر از آزمون درس خلاقیت نمایشی گروه هنر، سایر آزمون‌های مورد مطالعه بدون DTF جنسیتی بودند، در حالی که شدت DTF برای آزمون این درس نیز «کم» ارزیابی شد. به‌طور کلی، این نرخ و به‌ویژه شدت DIF و DTF می‌تواند تا حدود زیادی نشانگر نبود سوگیری جنسیتی در آزمون‌های کنکور سراسری باشد.

همان‌طور که اشاره شد DIF/DTF مرتبط با جنسیت لزوماً به معنای سوگیری جنسیتی نیست، زیرا برای سوگیری پرسش و آزمون وجود DIF/DTF لازم است

ولی کافی نیست (زومبو^۱، ۱۹۹۹). بنابراین برای قضاوت نهایی در مورد بود یا نبود سوگیری مرتبط با جنسیت، کمیته‌های تخصصی موضوعی تشکیل شد. این کمیته‌ها شامل ۶ تا ۱۲ تن از متخصصان و استادان درس‌های کنکور (بسته به آزمون درس تخصصی) شامل متخصصان زن و مرد بودند که با رهبری یک متخصص سنجش و اندازه‌گیری هدایت می‌شدند. معیارهای اصلی قضاوت در مورد سوگیری جنسیتی پرسش‌های آزمون در درجه اول یافته‌های مربوط به DIF و اندازه اثر DIF پرسش‌های آزمون و سپس محتوای پرسش‌های آزمون بود که از نظر میزان تطابق با تجربیات زیسته دختران و پسران یا زنان و مردان از سوی متخصصان موضوع درسی مورد بررسی و مذاقه قرار گرفت و سرانجام پیشینه پژوهشی به‌عنوان یکی از ملاک‌های اصلی برای قضاوت در مورد سوگیری پرسش‌های آزمون مد نظر گروه متخصصان قرار گرفت. به‌طور مثال در بررسی DIF و DTF آزمون درس خلاقیت نمایشی گروه هنر در کنکور سراسری سال ۱۳۹۰، پرسش ۲۰۳ دفترچه گروه هنر که دارای DIF آماری در درس مذکور بود، به‌صورت زیر مطرح شده است:

۲۰۳- این فیلم مستند داستان ندارد، گزارش روزانه زندگی یک پسر بچه است.

۱) سیاه و سفید ۲) خاک ۳) یک اتفاق ساده ۴) طبیعت بی‌جان

که بر اساس قضاوت‌های کیفی متخصصان موضوع درسی مذکور (اعم از متخصصان زن یا مرد) دارای سوگیری علیه جنس مؤنث بود. نتایج مرور و تشخیص نرخ سوگیری پرسش‌های DIF در درس‌های تخصصی بنا به قضاوت کمیته‌های موضوع درسی در جدول (۱۰) ملاحظه می‌شوند.

جدول (۱۰) نرخ سوگیری پرسش‌های DIF در درس‌های تخصصی بنا به قضاوت کمیته‌های

متخصص موضوع درسی

درصد متوسط در گروه	متوسط آزمون		۱۳۹۰		۱۳۸۹		۱۳۸۸		۱۳۸۷		سال‌های برگزاری آزمون	
	تعداد	درصد	تعداد	درصد	تعداد	درصد	تعداد	درصد	تعداد	درصد	گروه داوطلبی/درس	
۲/۶۸	۲/۷۳	۱/۵	۳/۶۴	۲	۰	۰	۱/۸۲	۱	۵/۴۵	۳	ریاضی	ریاضی و فیزیک
	۳/۸۹	۱/۷۵	۶/۶۷	۳	۲/۲۲	۱	۴/۴۴	۲	۲/۲۲	۱	فیزیک	
	۱/۴۳	۰/۵	۲/۸۶	۱	۰	۰	۲/۸۶	۱	۰	۰	شیمی	
۲/۶۸	۲/۵۰	۱/۲۵	۲	۱	۲	۱	۴	۲	۲	۱	زیست‌شناسی	

^۱. Zumbo

درصد	متوسط آزمون		۱۳۹۰		۱۳۸۹		۱۳۸۸		۱۳۸۷		سال‌های برگزاری آزمون	
	درصد	تعداد	درصد	تعداد	درصد	تعداد	درصد	تعداد	درصد	تعداد	گروه داوطلبی/درس	
	۲/۸۶	۱	۰	۰	۵/۷۱	۲	۲/۸۶	۱	۲/۸۶	۱	شیمی	علوم تجربی
۱/۶۷	۱/۶۷	۰/۵	۳/۳۳	۱	۳/۳۳	۱	۰	۰	۰	۰	ادبیات فارسی	علوم انسانی
۱/۶۷	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	درک عمومی هنر	هنر
	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	خلاقیت تصویری و تجسمی	
	۵	۱	۱۰	۲	۵	۱	۰	۰	۵	۱	خلاقیت نمایشی	
۲/۱۴	۲/۱۴	۱/۵	۱/۴۳	۱	۱/۴۳	۱	۲/۸۶	۲	۲/۸۶	۲	زبان انگلیسی تخصصی	زبان
	۲/۲۲		۲/۹۹		۱/۹۷		۱/۸۸		۲/۰۴		درصد متوسط برحسب سال	

همان‌طور که در جدول (۱۰) ملاحظه می‌شود، بر اساس قضاوت کمیته‌های تخصصی به‌طور متوسط حدود ۲ درصد پرسش‌های DIF در آزمون‌های مورد مطالعه کنکور سراسری دارای سوگیری که با توجه به اندازه‌های اثر بسیار کوچک و ناچیزی که به دست آمد، قابل توجیه است. در گروه ریاضی و فیزیک، پرسش‌های سودار بیشتر به نفع جنس مذکر عمل می‌کنند. به نظر می‌رسد پرسش‌هایی که برای محاسبه پاسخ نهایی به مراحل بیشتری نیاز دارند، به نفع جنس مذکر عمل می‌کنند. این نتیجه با یافته‌های دولیتل و کلییری (۱۹۸۷)، انابی و دودین (۲۰۰۶) و کارلتون و هریس (۱۹۹۲) همخوانی دارد، آنها نیز به این نتیجه رسیدند که پرسش‌هایی که در پاسخ به آن راهبردهای متعددی مورد نیاز است، به ضرر دانش‌آموزان دختر عمل می‌کند. آنها همچنین به این نتیجه رسیدند که پرسش‌های دارای نمودار یا جدول نیز به نفع دانش‌آموزان مذکر عمل می‌کنند. همچنین نرخ پرسش‌های سودار تشخیص داده شده در آزمون درس شیمی نسبت به سایر درس‌های انتخاب شده از گروه ریاضی و فیزیک کمتر است، به‌طوری که تنها ۲ پرسش در همه آزمون‌های برگزار شده، به دلایل محتوایی

فوق‌الذکر، نسبت به جنس مؤنث سوگیری داشتند. این در حالی بود که نرخ پرسش‌های DIF آماری علیه جنس مذکر در درس شیمی بیشتر بود.

همچنین نرخ پرسش‌های سودار تشخیص داده شده در جدول (۱۰) نمایانگر این است که پرسش‌های سودار درس‌های منتخب گروه‌های علوم تجربی و علوم انسانی بیشتر به نفع جنس مؤنث عمل می‌کنند. به نظر می‌رسد که پاسخگویی به این پرسش‌ها وابستگی بیشتری به حافظه و یادآوری دارند. این نتیجه تا حدودی با یافته‌های انابی (۲۰۰۸) همخوانی دارد. همچنین این یافته با پژوهش‌هایی همخوانی دارد که نشان می‌دهند دانش‌آموزان مؤنث در پرسش‌هایی که یادآوری مفاهیم تدریس شده، اجتناب از تفکر چالشی و ترجمه مفاهیم را می‌سازند، به‌صورت سوگیرانه‌ای عملکرد بهتری دارند (گالاکر، ۱۹۹۸؛ گالاکر و دلیسی، ۱۹۹۴ و هریس و کارلتون، ۱۹۹۳).

همان‌طور که در جدول (۱۰) ملاحظه می‌شود، بر اساس قضاوت کمیته متخصصان، در آزمون‌های درس‌های درک عمومی هنر و خلاقیت تصویری و تجسمی گروه هنر هیچ پرسش سوداری یافت نشد. اما در آزمون درس خلاقیت نمایشی، پرسش‌های سودار بر علیه جنس مؤنث عمل می‌کردند. بر اساس نظرات کیفی کمیته خبرگان این گروه، دلیل احتمالی DIF و DTF در این آزمون می‌تواند این باشد که محتوای پرسش‌های سودار این درس و به‌طور کلی این آزمون، بیشتر با تجارب جنس مذکر متناسب است، بدین معنی که بافت پرسش برای جنس مذکر، ملموس‌تر بوده و با تجربیات واقعی زندگی آنها بیشتر قابل پیوند است.

سرانجام، یافته‌های تحلیلی کمیته متخصصان درس زبان انگلیسی، نشان‌دهنده آن بود که پرسش‌های سوداری که محتوای علوم اجتماعی دارند بیشتر به نفع جنس مؤنث و پرسش‌های سوداری که محتوای علوم سخت دارند بیشتر به نفع جنس مذکر عمل می‌کنند. این نتیجه با یافته‌های امیریان، علوی و فیدالگو^۴ (۲۰۱۴) همخوانی دارد. همچنین این نتیجه با یافته‌های کارلتون و هریس (۱۹۹۲) و اونیل و مک‌پیک^۵ (۱۹۹۳) از این جهت همخوانی دارد که آنها نیز به این نتیجه رسیدند که پرسش‌های دارای مفاهیم انتزاعی به نفع جنس مؤنث عمل می‌کنند.

1. Gallagher

2. Gallagher & DeLisi

3. Harries & Carlton

4. Amirian, Alavi & Fidalgo

5. O'Neill & McPeck

یافته‌های مطالعه حاضر دلالت‌های کاربردی معینی در مورد DIF و DTF مرتبط با جنسیت در کنکور سراسری سازمان سنجش آموزش کشور دارد. از این مطالعه درباره تأثیر جنسیت بر عملکرد آزمودنی‌ها در کنکور سراسری اطلاعات سودمندی به دست آمد و ویژگی‌های محتوایی پرسش‌هایی که ممکن است نسبت به جنس مؤنث یا مذکر غیرمنصفانه (حتی به میزان خیلی کم و ناچیز) عمل کنند، مشخص شد. بنابراین تعیین نرخ پرسش‌های DIF و آزمون‌های DTF و میزان سوگیری آنها و تعیین ویژگی‌های محتوایی پرسش‌های سوادار، از جمله هدف‌هایی بوده است که در پژوهش حاضر، محقق شد.

برای پژوهش‌های آینده پیشنهاد می‌شود که بدون در نظر گرفتن اطلاعات آماری مربوط به DIF مرتبط با جنسیت، پرسش‌های کنکور صرفاً به لحاظ محتوایی و به صورت کیفی توسط گروه خبرگان بررسی و نتایج حاصل با یافته‌های این پژوهش، مقایسه شوند. انگل‌هارد، هانش و روتلج^۱ (۱۹۹۰) چنین رویکردی را برای مطالعات DIF پیشنهاد کرده است. همچنین یافته‌های پژوهش نشان داد که تنها آزمون خلاقیت نمایشی که دارای بیشترین نرخ پرسش‌های DIF است، دارای DTF معنی‌داری بود. اگرچه چندین مطالعه تجربی در زمینه رابطه DIF و DTF آزمون انجام گرفته است (تاکالا و کاftاندیجوا^۲، ۲۰۰۰؛ زومبو، ۲۰۰۳)، اما هنوز مشخص نیست که ۱- وجود DTF در آزمون ناشی از تأثیر تراکمی پرسش‌های DIF آزمون است، کما اینکه چنین پدیده‌ای در مطالعه حاضر به صورت موردی مشاهده شد، ولی هنوز نمی‌توان آن را تعمیم داد و یا ۲- آزمون دارای پرسش‌های DIF به دلیل لغو DIF در سطح تحلیل آزمون، DTF را نشان نمی‌دهند و یا اینکه ۳- آزمونی با پرسش‌های DIF، DTF را نشان نمی‌دهد، زیرا DIF از DTF مستقل است. بنابراین پیشنهاد می‌شود که رابطه میان وجود پرسش‌های DIF در آزمون و کنش افتراقی آزمون (DTF) داده‌های شبیه‌سازی شده و داده‌های تجربی کنکور سراسری با حجم نمونه‌های بیشتر، به منظور کاربرد در تحلیل داده‌های کنکور سراسری، بررسی و مطالعه شود.

1. Engelhard, Hansche & Rutledge

2. Takala & Kaftandjieva

منابع

- آبکار، کبری (۱۳۹۱). بررسی ویژگی‌های روان‌سنجی سؤالات کنکور سراسری در رشته علوم تجربی سال ۱۳۸۹ از نظر تئوری سؤال و پاسخ (IRT). پایان‌نامه کارشناسی ارشد، دانشگاه آزاد اسلامی واحد تهران مرکز.
- آسیابی، مینا (۱۳۹۱). ارزشیابی آزمون کارشناسی ارشد رشته جغرافیای سیاسی با استفاده از مدل‌های جدید اندازه‌گیری و تعیین سوگیری جنسیتی در آن. پایان‌نامه کارشناسی ارشد، دانشگاه علامه طباطبایی.
- ترکاشوند، علی (۱۳۹۴). بررسی ویژگی‌های روان‌سنجی آزمون سراسری درس زیست‌شناسی بر اساس مدل چندگزینه‌ای IRT. پایان‌نامه کارشناسی ارشد، دانشگاه خوارزمی.
- حبیبی، مجتبی (۱۳۹۲). بررسی عوامل مؤثر بر پیشرفت تحصیلی دانشجویان مقطع کارشناسی و پیش‌بینی آن بر اساس نمرات تراز کنکور: اعتباریابی بیرونی نمرات تراز کنکور با مطالعه موردی دانشگاه شهید بهشتی. طرح پژوهشی، وزارت علوم، تحقیقات و فناوری.
- فلاحی سرشت، شیوا (۱۳۹۴). بررسی کارکرد افتراقی سؤالات (DIF) استعداد تحصیلی آزمون نیمه‌متمرکز دکتری سال ۹۳ با کاربرد نظریه سؤال-پاسخ (IRT) و رگرسیون لجستیک. پایان‌نامه کارشناسی ارشد، دانشگاه علامه طباطبایی.
- گرامی‌پور، مسعود (۱۳۹۳). مبانی نظری و کاربرد نظریه‌های اندازه‌گیری در علوم رفتاری. تهران: انتشارات تمدن علمی.
- گرامی‌پور، مسعود و فلسفی‌نژاد، محمدرضا (۱۳۹۲). روش‌های آماری بررسی کنش افتراقی سؤال (DIF) در آزمون‌های سرنوشت‌ساز. تهران: انتشارات جهاد دانشگاهی، واحد تربیت معلم.
- معلمی‌اوره، مهرناز (۱۳۸۷). مقایسه دقت برآورد توانایی در سؤالات چندگزینه‌ای با به‌کارگیری مدل‌های سؤال-پاسخ دو و چندارزشی. پایان‌نامه کارشناسی ارشد، دانشگاه علامه طباطبایی.
- میری، محمد (۱۳۹۴). بررسی و مقایسه ویژگی‌های روان‌سنجی بخش فیزیک آزمون سراسری ورود به دانشگاه بر اساس مدل‌های دو ارزشی IRT. پایان‌نامه کارشناسی ارشد، دانشگاه خوارزمی.

مینایی، اصغر (۱۳۹۲). سنجش مقایسه پذیری سازه و تحلیل کارکرد افتراقی سؤال‌ها (DIF) و بلوک‌های (DTF) آزمون علوم پایه هشتم تیمز ۲۰۰۷ در بین دانش آموزان ایران و آمریکا. فصلنامه اندازه‌گیری تربیتی، ۴ (۱۱)، ۱۰۹-۱۴۶.

نژادنجف، فیروز (۱۳۹۳). نقد و بررسی سؤالات کنکور سراسری درس دین و زندگی. رشد آموزش معارف اسلامی، ۲۶، ۴۸-۵۳.

- Amirian, S. M.R.; Alavi, S. M. & Fidalgo, A. M. (2014). Analyzing Gender Differences with an English Proficiency Test in EFL Context. *Iranian Journal of Language Testing*.
- Aryadoust, V.; Goh, C. C. M. & Kim, L. O. (2011). An investigation of differential item functioning in the MELAB listening test. *Language Assessment Quarterly*, 8 (4), 361-385.
- Barati, H. & Ahmadi, A. R. (2010). Gender-based DIF across the subject area: A study of the Iranian National University Entrance Exam. *The Journal of Teaching Language Skills (JTLS)*, 2 (3), 1-22.
- Berberoglu, G. (1995). Differential item functioning (DIF) analysis of computation, word problem and geometry questions across gender and SES groups. *Studies in Educational Evaluation*, 21 (4), 439-456.
- Breland, H.; Lee, Y. W.; Najarian, M. & Muraki, E. (2004). *An analysis of the TOEFL CBT writing prompt difficulty and comparability of different gender groups* (TOEFL Research Report No. 76). Princeton, NJ: Educational Testing Service.
- Brown, I. & Kanyongo, Y. (2007). Differential Item Functioning and male-female differences in a large-scale mathematics assessment in Trinidad and Tobago. *Caribbean Curriculum*, 14, 49-71.
- Carlton, S. T. & Harris, A. M. (1992). *Characteristics associated with differential item functioning on the Scholastic Aptitude Test: Gender and majority/minority group comparisons*. Princeton, NJ: Educational Testing Service.
- Chalmers, R. P.; Counsell, A. & Flora, D. B. (2015). It might not make a big DIF: Improved Differential Test Functioning statistics that account for sampling variability. *Educational and Psychological Measurement*, 1-27.
- Doolittle, A. E. & Cleary, T. A. (1987). Gender-based differential item performance in mathematics achievement items. *Journal of Educational Measurement*, 24, 157-166.
- Doudeen, Hamzah M. & Annabi, Hanan A. (2008). Sex-Related Differential Item Functioning (DIF) Analysis of TIMSS. *Dirasat, Educational Sciences*, Vol. 35.

- Drasgow, F. (1984). Scrutinizing psychological tests: Measurement equivalence and equivalent relations with external variables are central issues. *Psychological Bulletin*, 95, 134-135.
- Drasgow, F. (1987). Study of the measurement bias of two standardized psychological tests. *Journal of Applied Psychology*, 72, 19-29.
- Embretson, S. E. & Reise, S. P. (2000). *Item response theory for psychologists*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Engelhard, G.; Hansche, L. & Rutledge, K. (1990). Accuracy of bias review judges in identifying differential item functioning on teacher certification tests. *Applied Measurement in Education*, 3, 347-360.
- Ethington, A. (1990). Gender differences in mathematics: An international perspective. *Journal for Research in Mathematics Education*. 21 (1), 74-80.
- Fennema, E. (1980). Sex-related differences in mathematics achievement: Where and why. In Fox, L. H.; Brody, L. & Tobin, D. (Eds.). *Women and the mathematic mystique*, (pp. 76-93). Baltimore: Johns Hopkins University Press.
- Fennema, E. & Carpenter, T. P. (1981). Sex-related differences in mathematics: Results from national assessment. *Mathematics Teacher*. 74, 554-559.
- Finch, H. & Habing, B. (2007). Performance of DIMTEST- and NOHARM based statistics for testing unidimensionality. *Applied Psychological Measurement*, 31, 292-307.
- Flora, D., Curran, P., Hussong, A., & Edwards, M. (2008). Incorporating measurement Nonequivalence in a cross-study latent growth curve analysis. *Structural Equation Modeling*, 15, 676-704.
- Fraser, C., & McDonald, R. P. (1988). NOHARM: Least squares item factor analysis. *Multivariate Behavioral Research*, 23, 267-269.
- Gallagher, A. (1998). Gender and antecedents of performance in mathematics testing. *Teachers College Record*, 100 (2), 297-314.
- Gallagher, A. M., & DeLisi, R. (1994). Gender differences in scholastic aptitude tests mathematics problem solving among high-ability students. *Journal of Educational Psychology*, 86, 204-211.
- Hanna, G. (1989). Mathematics achievement of girls and boys in grade eight: Results from twenty countries. *Educational Studies in Mathematics*, 20, 225-232.
- Harries, A. & Carlton, S. (1993). Patterns of gender difference on mathematics items on the scholastic aptitude test. *Applied Measurement in Education*, 6 (2), 151- 173.
- Husen, T. (1967). *International study of achievement in mathematics: A comparison of twelve countries*. Vol. 11. Stockholm: Almqvist & Wiksell.

- Innabi, H., & Dodeen, H. (2006). Content Analysis of Gender-related Differential Item Functioning of TIMSS Items in Mathematics in Jordan. *School Science and Mathematics*, 106 (8), 328-337.
- Le, Luc T. (2006). Investigating gender differential item functioning across Countries and Test Languages for PISA science items. *International Journal of Testing*, 9, 2, 122-133.
- O'Neill, K. A. & McPeck, W. M. (1993). Item and test characteristics that are associated with differential item functioning. In Holland, P. W. & Wainer, H. (Eds.), *Differential item functioning*, (pp. 255-276). Hillsdale, N J: Lawrence Erlbaum.
- Pae, H. K. (2011). Differential item functioning and unidimensionality in the Pearson Test of English Academic. <http://pearsonpte.com/research/Documents/Pae.pdf>.
- Pae, T. & Park, G. P. (2006). Examining the relationship between differential item functioning and differential test functioning. *Language testing*, 23 (4), 475-496.
- Park, G. P. (2008). Differential item functioning on an English listening test across gender. *TESOL Quarterly*, 42 (1), 115-123.
- Pattison, P. & Grieve, N. (1984). Do spatial skills contribute to sex differences in different types of mathematical problems? *Journal of Educational Psychology*, 76 (4). 677-689.
- Raju, N. S.; van der Linden, W. J. & Fleer, P. F. (1995). IRT-based internal measures of differential functioning of items and tests. *Applied Psychological Measurement*, 19, 353-368.
- Rudner, L.; Getson, P. & Knight, D. (1980). Biased item detection techniques. *Journal of Educational Statistics*, 5, 213-233.
- Russell, S. S. (2005). *Estimates of Type I error and power for indices of differential bundle and test functioning*. Ph.D. dissertation, Bowling Green State University, United States -- Ohio.
- Takala, S. & Kaftandjieva, F. (2000). Test fairness: A DIF analysis of an L2 vocabulary test. *Language Testing*, 17, 323-340.
- Wang, N. & Lane, S. (1996). Detection of gender-related differential item functioning in a mathematics performance assessment. *Applied Measurement in Education*, 9 (2), 175-199.
- Wood, R. (1976). Sex differences in mathematics attainment at GCE ordinary level. *Educational Studies*, 2. 141- 160.
- 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-Type (Ordinal) Item Scores*. Ottawa, ON: Directorate of Human Resources Research and Evaluation, Department of National Defense.
- Zumbo, B. (2003). Does item-level DIF manifest itself in scale-level analysis? Implications for translating language tests. *Language Testing*, 20, 136-147.

Archive of SID