



## پژوهش‌های زبان‌شناسختی در زبان‌های خارجی

شایانی چاپی: ۴۱۲۳-۲۵۸۸ شایانی الکترونیکی: ۷۵۲۱-۲۵۸۸

www.jflr.ut.ac.ir



# بررسی کنش افتراقی پرسش‌ها و عملکرد در آزمون: مقایسه رگرسیون لجستیک، مدل رش و منتل-هنزل

حسین کرمی\*

(نویسنده مسئول)

استادیار، گروه زبان و ادبیات انگلیسی، دانشگاه تهران،  
تهران، ایران

Email: hkarami@ut.ac.ir



علی خودی\*\*

دکتری گروه زبان و ادبیات انگلیسی، پردیس کیش، دانشگاه تهران،

کیش، ایران

Email: Ali.khodi@ut.ac.ir



## اطلاعات مقاله

تاریخ ارسال: ۱۳۹۹/۰۹/۱۹

تاریخ پذیرش: ۱۳۹۹/۱۰/۱۳

تاریخ انتشار: زمستان ۱۳۹۹

نوع مقاله: علمی پژوهشی

## کلید واژگان:

کنش افتراقی، روایی، عدالت  
آزمون، تبعیض

پژوهش‌های زبان‌شناسختی در زبان‌های خارجی، دوره ۱۰، شماره ۴، زمستان ۱۳۹۹، از صفحه ۸۴۲ تا ۸۵۳

## چکیده

یکی از ابزارهای بررسی عملکرد آزمون، «کنش افتراقی پرسش‌ها» است (Differential item functioning). این روش، می‌تواند عوامل تاثیرگذار بر عملکرد آزمودنی‌ها را پیدا کرده و از بروز سوگیری در آزمون جلوگیری کند. در دو دهه گذشته، روش‌های بسیاری برای تشخیص پیشنهاد شده است. شمار روش‌های تشخیص عملکرد افتراقی پرسش، گاهی باعث سردرگمی پژوهشگران می‌شود. از سوی دیگر، امکان مقایسه یافته‌های پژوهش‌هایی که با روش‌های مختلف به بررسی کنش افتراقی پرسش را پرداخته‌اند، دشوار می‌سازد. پژوهش حاضر به بررسی و سنجه‌ش نتایج بدست آمده از سه روش تشخیص کنش افتراقی پرسش پرداخته است: مدل رش، رگرسیون لجستیک و منتل-هنزل (MH). داده‌های استفاده شده در بررسی‌ها، برگرفته از آزمون توانش انگلیسی دانشگاه تهران (UTEPT) هنزل (MH)، داده‌های استفاده شده در بررسی‌ها، برگرفته از آزمون توانش انگلیسی دانشگاه تهران (UTEPT) است که آزمون با اهمیت ویژه‌ای است و سالانه برای داوطلبان دکترا برگزار می‌شود. تجزیه و بررسی کنش افتراقی یکنواخت با سه روش برگفته شان داد که پرسش‌ها در عملکرد خود تفاوت‌های زیادی ندارند. نتایج تحلیل رگرسیون لجستیک، دو پرسش را برای وجود کنش افتراقی پیدا کرد که مشابه روش منتل-هنزل است. همچنین پرسش‌هایی به عنوان نشانگرهای کنش افتراقی قوی در مدل رش شناسایی شده بودند، همان پرسش‌ها بودند که در دو مدل دیگر نیز معرفی شده بودند. نتایج پژوهش حاضر نشان می‌دهد که استفاده از روش‌های مختلف برای بررسی وجود کنش افتراقی پرسش، الزاماً نتایج متفاوت ناگزیری را در پی ندارد و می‌توان از هر یک از روش‌های استفاده شده در این پژوهش بهره گرفت.

کلیه حقوق محفوظ است ۱۳۹۹

DOI: 10.22059/jflr.2021.315079.783

کرمی، حسین، خودی، علی. (۱۳۹۹). بررسی کنش افتراقی پرسش‌ها و عملکرد در آزمون: مقایسه رگرسیون لجستیک، مدل رش و منتل-هنزل. پژوهش‌های زبان‌شناسختی در زبان‌های خارجی، ۱۰(۴)، ۸۴۲-۸۵۳.

Khodi, Ali, Karami, Hossein (2021). Differential Item Functioning and Test Performance: a Comparison Between the Rasch Model, Logistic Regression and Mantel-Haenszel. *Journal of Foreign Language Research*, 10 (4), 842-853.  
DOI: 10.22059/jflr.2021.315079.783

\* دکتر کرمی متخصص رشته آموزش زبان انگلیسی بوده و به طور تخصصی در زمینه بررسی روایی آزمون‌ها و سنجه‌ش و ارزیابی پژوهش می‌نماید.

\*\* علی خودی دانش‌آموخته رشته آموزش زبان انگلیسی از دانشگاه تهران بوده و در حیطه پژوهشی آزمون‌سازی فعالیت می‌نماید.



# Differential Item Functioning and Test Performance: a Comparison Between the Rasch Model, Logistic Regression and Mantel-Haenszel



**Hossein Karami\***  
(corresponding author)  
Assistant Professor, University of Tehran,  
Tehran, Iran.  
Email: [hkarami@ut.ac.ir](mailto:hkarami@ut.ac.ir)



**Ali Khodi\*\***  
PhD, University of Tehran, Kish international Campus,  
Kish, Iran  
Email: [Ali.khodi@ut.ac.ir](mailto:Ali.khodi@ut.ac.ir)

## ABSTRACT

Differential item functioning (DIF) is considered to be one of the tools for the examination of test fairness. This method is capable of finding the factors affecting the subjects' performance and prevent the occurrence of bias in the test. A plethora of methods for detecting Differential Item Functioning (DIF) has been suggested during the last couple of decades. The multiplicity of methods for diagnosing DIF might be a confusing issue for applied researchers and might lead to complications in the comparability of the findings of various DIF studies which have utilized different DIF detection techniques. This study aimed to investigate the comparability of results from three widely used DIF detection techniques: the Rasch model, Logistic Regression, and Mantel-Haenszel (MH). The data comes from an administration of the University of Tehran English Proficiency Test (UTEPT) which is a high-stakes test administered annually to PhD candidates. DIF analysis through the three techniques indicated that the three methods did not have significant differences in their performance. The Mantel-Hansel model flagged two items having DIF just similar to the findings of logistic regression model. Likewise, the items that were detected as strong-DIF items in Rasch model were the same as items detected by the two aforementioned models. Therefore, it might be concluded that use of different DIF detection techniques does not necessarily lead to flagging different items.

## ARTICLE INFO

Article history:  
Received: 9th, December, 2020  
Accepted: 2nd, January, 2021  
Available online: Winter 2021

## Keywords:

Differential Item Functioning, fairness, bias, validity

DOI: [10.22059/jflr.2021.315079.783](https://doi.org/10.22059/jflr.2021.315079.783)

© 2021 All rights reserved.

Khodi, Ali, Karami, Hossein (2021). Differential Item Functioning and Test Performance: a Comparison Between the Rasch Model, Logistic Regression and Mantel-Haenszel. *Journal of Foreign Language Research*, 10 (4), 842-853.  
DOI: [10.22059/jflr.2021.315079.783](https://doi.org/10.22059/jflr.2021.315079.783)

\* Dr. Karami is an assistant professor of Applied Linguistics. His main research interests include various aspects of language testing and assessment.

\*\* Ali Khodi is a PhD holder of applied linguistics who is interested in language testing.

## ۱. مقدمہ

بررسی ویژگی‌های روانشناسی دانش آموزان در زمینه های مختلف، مانند آموزش؛ بیشتر با استفاده از آزمون انجام می‌شود. هدف از این آزمایش‌ها ارزیابی افراد از نظر توانایی مورد نظر است که سازه‌ی مد نظر سنجش نام دارد (آکار و کله سی اوغلو، *Construct of measurement*) (Acar & Kelecioglu, ۲۰۱۰). از آنجا که مشخصات توانایی‌های افراد باید از طریق یک ابزار اندازه‌گیری واجد شرایط (به عنوان مثال آزمون‌ها) انجام شود، وجود هر عامل یا فاکتور غیر مرتبط سازه‌ای ممکن است کاربرد و قابلیت آزمون را کاهش دهد. بنابراین، آزمون‌ها به طور کلی و موارد آزمایشی به طور خاص باید بتوانند توانایی افراد را اندازه‌گیری کنند، بدون اینکه ویژگی‌های متفرقه‌ای از شرکت کنندگان، کارکرد آنان را در سنجش سازه مورد نظر تحت تاثیر قرار دهد (یوآر، کلیسیاوغلو، و دوغان، Uyar, Kelecioğlu, & Doğan, ۲۰۱۷).

وجود این حساسیت به این دلیل است که افراد با توانایی برابر اما از گروه‌های مختلف بتوانند با همان درجه و شانس برابر به سوال‌های آزمون پاسخ دهنند. اگر در یک آزمون پرسش‌هایی وجود داشته باشد که احتمال میزان پاسخ‌گویی به آن توسط افراد از جنس، سن یا گروه قومی خاص بیش از سایرین باشد، احتمال وجود تبعیض یا کنش افتراقی غیرقابل تصور نیست (کامرون و همکاران، Camero et al., ۲۰۱۴) با این حال باید به ایده داشت که وجود کنش افتراقی لزوماً به معنای وجود تبعیض در سنجش نیست و در حقیقت شرط لازم، اما ناکافی برای آن است (مکنامارا و روور، McNamara, & Roever, ۲۰۰۶).

فرایندهای تعیین وجود تبعیض در عملکرد پرسش‌ها را می‌توان از راه بررسی پاسخ‌های افراد با توانایی یکسان، اما گروه‌های مختلف (به عنوان مثال نژاد، جنسیت و غیره) انجام داد. بنابراین، در این بررسی احتمال پاسخ‌های صحیح داده شده به یک پرسش به عنوان کنش افتراقی آن (DIF) تعریف شده است (استینبرگ و تیسن، Steinberg, & Thissen, 2006). کنش افتراقی به عنوان احتمال مشروط برای دستیابی به پاسخ‌های صحیح تعریف می‌شود و زمانی اتفاق می‌افتد که یک پرسش موارد بیشتری از آن چیزی

این انگاره که اعتبار (Validity) مهمترین ملاحظه در گسترش و بهره‌گیری از آزمون است، به واقعیتی بدینه تبدیل شده است (چپل، Chapelle, ۲۰۲۰). بنابراین، اطمینان از معتبر بودن عملکرد یک آزمون بر عهده طراحان و توسعه‌دهندگان آزمون و کاربران آن است (بکمن، Bachman, ۱۹۹۰). یکی از تهدیدات و موانع اصلی، علیه اعتبار آزمون وجود واریانس برگرفته از عوامل بی‌ربط و سازه‌های نامرتب (construct-irrelevant factors) است، بدین معنا که عملکرد در آزمون، نباید تحت تأثیر عواملی قرار گیرد که متفاوت از سازه مورد نظر آزمون است. در غیر این صورت، آزمون، داری سوگیری خواهد بود (کرمی، Karami, ۲۰۱۳). در این راستا از روشی آماری که به طور گسترده برای تشخیص سوگیری پرسش‌ها در آزمون استفاده می‌شود، عملکرد کنش افتراقی پرسش‌ها (DIF) است.

کنش افتراقی هنگامی رخ می‌دهد که شرکت کنندگان در آزمون با سطح توانایی برابر، اما از دو گروه مختلف، شانس متفاوتی در پاسخ‌گویی به یک پرسش را داشته باشند. با این حال، نمی‌توان عنوان کرد که کنش افتراقی برایر و همسان، به معنای وجود تبعیض (bias) در سنجش است. بلکه می‌توان متصور شد که پیش‌شرطی برای به وجود آمدن تبعیض در سنجش است. بنابراین، وجود کنش افتراقی، شرط اساسی، اما ناکافی برای سوگیری در سنجش است. شیوه‌های مختلفی برای تشخیص کنش افتراقی در پژوهش‌های پیشین ارائه شده است (کلاوزر و مازور، Clauser & Mazor, ۱۹۹۸؛ کرمی، Karami, ۲۰۱۲). آنچه مطرح است، تفاوت‌های عملکردی این شیوه‌های سنجش کنش افتراقی در بررسی‌های انجام شده است. بنابراین، بایسته است که چگونگی و امکان مقایسه این شیوه‌های واکاوی را در پژوهش‌های مختلف بسنجیم و بررسی کنیم. در همین راستا، پژوهش می‌کوشد به بررسی عملکرد سه روش تشخیص کنش افتراقی پرسش‌ها بپردازد: مدل رش، رگرسیون لجستیک و متنل-هنزل. همچنین در پژوهش حاضر اثرگذاری این شیوه‌های سنجش کنش افتراقی بر روایی و پایایی داده‌ها و یافته‌های آزمون مورد ارزیابی قرار می‌گیرد.

۴ پرسش آزمون، میزان آن قابل چشم پوشی است. با بررسی بانک پرسش‌های که برای آزمون ورودی موسسه مورد استفاده قرار می‌گرفت، اثنا اولیوری و همکاران (Elena Oliveri et al., 2018) به بررسی وجود کنش افتراقی برگرفته از اثرات میزان سطح و توانش زبانی، سن، ملیت و سطح و طبقه اجتماعی پرداختند. بدین منظور از دو روش متنل هنزل و نظریه پرسش و پاسخ بهره گرفتند. یافته‌های آنان حاکی از آن بود که از باب ملیت افراد، پرسش‌ها دارای کنش افتراقی‌اند. در پژوهشی دیگر نیز توسط چن، لیو و زومبو (Chen, Liu, & Zumbo, 2020) روش جدید از سنجش کنش افتراقی پیشنهاد شده است که بر اساس نمره کل آزمون عملکرد و نیاز به سنجش و ارزیابی بیشتری دارد.

### ۳. روش انجام پژوهش

#### شرکت‌کنندگان

شرکت‌کنندگان این پژوهش ۳۰۰۰ تن از مقاضیانی بودند(از هر دو جنسیت مرد و زن) که از بین کل شرکت‌کنندگان آزمون توانش مهارت زبان انگلیسی دانشگاه تهران (UTEPT) انتخاب شدند. شرکت‌کنندگان آزمون دانشجویان مشغول به تحصیل در مقطع دکتری رشته‌های مختلف در دانشگاه تهران بودند که بازه سنی بیشتر؛ در رده سنی ۲۵ تا ۴۰ ساله بوده است. شرکت‌کنندگان در پژوهش حاضر بر اساس سوابق تحصیلی خود به دو گروه علوم انسانی و علوم و فناوری تقسیم شدند(هر گروه ۱۵۰۰ نفر). لزوم انتخاب تعداد برابر شرکت‌کنندگان برای هر دو گروه بهدلیل پیشگیری از اثرگذاری اندازه نمونه بر یافته‌های پژوهش است.

#### ابزار پژوهش

مقاضیان دوره‌های دکتری دانشگاه تهران بایسته‌است که در آزمون سنجش توانش و مهارت زبان‌های خارجی به نام "آزمون مهارت‌انگلیسی دانشگاه تهران" (UTEPT) شرکت‌نموده و نتایج آن را به دپارتمان مربوطه ارائه‌دهند. هدف اصلی این آزمون شناسایی افرادی است که سطح مهارت انگلیسی مناسب و قابل قبولی دارند. این آزمون از سه بخش شامل دستور زبان، خواندن و درک مفاهیم و واژگان

را که باید بسنجد اندازه‌گیری می‌کند که آن متغیر اضافی اندازه‌گیری تعریف شده است و به عنوان "پارامتر مداخله گر" (nuisance) (Ackerman, ۱۹۹۲) نامیده می‌شود. کنش افتراقی پرسش‌ها به دو دسته یکنواخت (uniform) و غیر یکنواخت (non-uniform) طبقه بندی می‌شود و به طور سنتی با مقایسه پاسخ‌های داده شده به پرسش‌ها توسط افراد با توانایی مشابه اما متعلق به گروه‌های متفاوت به نام گروه مرجع (reference group) یا گروه آزمایش (focal group) بررسی می‌شود (McNamara و روور، ۲۰۰۶). کنش افتراقی یکنواخت نیز به صورت احتمال پاسخ صحیح به پرسشی تعریف شده است که به طور یکنواخت برای همگی افراد با سطح توانایی متفاوت در گروه الف از افراد در گروه ب بیشتر باشد (Zumbo, ۲۰۰۳) در حالی که کنش افتراقی غیر یکنواخت که گاهی به عنوان کنش متقاطع (crossing) نیز شناخته می‌شود به حالتی اشاره می‌کند که یک پرسش برای افراد با سطح توانایی متفاوت در یک گروه، عملکرد متفاوتی را نشان می‌دهد.

برای تشخیص و بررسی کنش افتراقی شیوه‌های متفاوتی بررسی و به کار برده می‌شود که برخی از آنان (Classical Test Theory) براساس نظریه کلاسیک آزمون (Lord's chi-square test) یا نسبت احتمال (likelihood) که بر اساس نظریه پرسش و پاسخ (IRT) است شناخته می‌شوند (Camili و Shepard, ۱۹۹۴).

برای مثال، در پژوهشی انجام شده ژو و آریادوست (Zhu & Aryadoust, ۲۰۲۰) تاثیر زبان مادری شرکت‌کنندگان آزمون و ارتباط آن را با کنش افتراقی مورد بررسی قرار گرفته است. یافته‌های آن نشان داد پرسش‌های آزمون، متاثر از زبان مادری زبان‌آموزان نبوده است. در پژوهشی دیگر، اثر کنش افتراقی پرسش‌ها بین شرکت‌کنندگان آلمانی و انگلیسی سنجیده شده است (Fischer و همکاران, ۲۰۱۶ al.). یافته‌ها حاکی از آن است که با وجود این اثر در

سوم باید از طریق آزمون خی<sup>۲</sup> (Chi-square) ارزیابی شود.

در بحث شیوه سنجش کنش افتراقی پرسش‌ها با نظریه رگرسیون لجستیک می‌توان عنوان داشت که در سه حالت قابل تبیین است: در حالت اول، مدل فقط در بردارنده نمره واقعی (true score) است در مدل دوم، بررسی در بردارنده نمره واقعی و متغیر دسته‌بندی کننده است و در مدل سوم متغیری به نام همکنشی هر دو متغیر (ذکر شده در حالت دوم) به بررسی‌ها اضافه می‌شود. تفاوت بین حالت اول و سوم از طریق نمره خی<sup>۲</sup> (Chi square) قابل بررسی است که اگر میزان معناداری را نشان داد می‌توان در نظر داشت که کنش افتراقی یکنواخت یا غیریکنواخت در این آزمون وجود دارد. گام دوم، بررسی تفاوت مدل اول و دوم است که نشانگر وجود کنش افتراقی یکنواخت است، در حالی که مقایسه مدل دوم و سوم بیانگر وجود کنش افتراقی غیر یکنواخت است. در باره حجم نمونه نیز می‌توان از مربع ضریب آر (R<sup>2</sup>) استفاده کرد ([جودوین و گریل، Jodoin & Gierl، ۲۰۰۱](#)).

در این معیار مقادیر کمتر از ۰,۳۵ نوع الف کنش افتراقی بوده و قابل چشم‌پوشی است، بین ۰,۳۵ تا ۰,۷۰ نوع ب کنش افتراقی و میزان متوسط است و بیش از ۰,۷۰ نوع ج وجود کنش افتراقی بوده و مقادیر زیادی را نشان می‌دهد.

در باره سنجش کنش افتراقی با روش متنل-هنزل نیز اگر شاخص سنجش (MH D-DIF) مثبت باشد، حاکی از آن است که آن پرسش برای گروه مرجع دشوارتر از شرکت‌کنندگان در گروه آزمایش بوده است و اگر آن شاخص منفی باشد، نشانگر آن است که این پرسش برای گروه آزمایش دشوارتر بوده است. بر اساس معیارهای سازمان ETS (زیکی، [Zieky، ۱۹۹۳](#)) پرسش‌ها نوع الف، کنش افتراقی را از خود نشان می‌دهند، اگر میزان قطعی شاخص MH D-DIF کمتر از یک باشد و نوع ج کنش افتراقی را نشان می‌دهند، اگر آن شاخص بیشتر از ۱,۵ باشد یا به طور معناداری از یک بیشتر باشد و سایر پرسش‌ها با شاخص‌های متفاوت از آنچه اشاره شد از نوع ب کنش افتراقی‌اند.

#### ۴. نتایج و بحث و بررسی

با بررسی‌های انجام شده و واکاوی داده‌های پژوهش به شیوه‌های رگرسیون لجستیک، متنل-هنزل و مدل رش

تشکیل شده است. همه پرسش‌ها در قالب چند گزینه‌ای ارائه شده‌اند و بخش خواندن و درک مفاهیم شامل متونی است که بی‌درنگ پس از آن‌ها شماری پرسش درک مطلب به شرکت‌کنندگان ارائه می‌شود. شمار پرسش‌های درک مطلب برای هر بخش متفاوت است. در این آزمون یکنمره کل به داوطلبان گزارش می‌شود که به‌سادگی مجموع نمراتی است که در سه خرده آزمون کسب کرده‌اند. در این مطالعه ما به بررسی بخش دستور زبان این آزمون پرداخته‌ایم که شامل پرسش‌های مرتبط به مفاهیم دستور زبان انگلیسی می‌باشد.

#### تحلیل داده‌ها

طی سه دهه گذشته انواع مختلفی از شیوه‌های تخصیص کنش افتراقی پرسش‌ها ارائه شده است که از جمله می‌توان به روش‌های ساده مبتنی بر شاخص‌های دشواری پرسش‌ها (نمودار دلتا) یا شیوه‌های پیچیده مبتنی بر نظریه پرسش و پاسخ (IRT) اشاره نمود. رویکردهای مبتنی بر نظریه سوال-پاسخ به دلیل ظرافت مفهومی خود، از جمله روش‌های تشخیص کنش افتراقی‌اند که به‌طور گسترده استفاده می‌شوند. در مطالعه حاضر، به‌طور خاص از سه روش برای تعزیز و بررسی داده‌ها استفاده شده است: مدل رش، رگرسیون لجستیک و متنل-هنزل.

در رگرسیون لجستیک، پاسخ پرسش‌ها به عنوان متغیر وابسته در نظر گرفته می‌شود که باید از طریق سایر متغیرها پیش‌بینی شود. متغیرهای مورد نظر در این شیوه نمره‌کل، متغیر گروه‌بندی و همکنشی بین این دو هستند. روش متنل-هنزل یک رویکرد تشخیص کنش افتراقی غیرپارامتریک است که بر اساس ایده نسبت شانس (odd ratio) استوار است. نظریه نسبت شانس با "جمع آوری اطلاعات در بین سطوح متفاوت متغیر در حال سنجش (به طور معمول نمره آزمون مشاهده شده) به دست می‌آید تا ارزیابی شود که احتمال موفقیت در یک پرسش خاص برای گروه مرجع در مقایسه با سایر شرکت‌کنندگان در گروه آزمایش چقدر متفاوت است" ([سیرسی و ریوس، Sireci & Rios، ۲۰۱۳](#)، ص ۱۷۵)

همچنین برای ارزیابی کنش افتراقی، سه مدل رگرسیون را می‌توان تعیین کرد که مدل اول فقط نمره کل را شامل می‌شود مدل دوم هم نمره کل و عامل گروه‌بندی، و مدل آخر این دو عامل به‌اضافه مدت تعامل. تفاوت بین مدل اول و

هر پرسش وجود دارد. در باره ستون سوم و چهارم باید عنوان داشت که در بخش اول (MNSQ) باید توجه داشته باشیم که این عدد به صورت کلی باید بین  $+1,3$  و  $-1,3$  باشد، اگر چنین نبود به ستون بعدی آن یعنی ZSTD که اشاره به مربع نمره زی (Z-square) دارد رجوع کرده و در این ستون Z باید عددی بین  $-2$  تا  $+2$  باشد، در غیر این صورت آیتم نرمال نیست و از تناسب در مدل خارج است (Overfit) (لينکر، [Lincare](#)) که در این پژوهش چنین پرسش یافت نشد. قوانین مشابهی برای ستون ۵ و ۶ نیز قابل اعمال است. در ستون شماره ۷ و ۸ به همبستگی‌های اندازه گیری نقطه‌ای (point-measure correlation) پرداخته شده که نشانگر رابطه بین عملکرد افراد در پرسش و توانایی کلی افراد بر مبنای کل آزمون است. در ستون مورد انتظار (Expected) میزان همبستگی مورد انتظار پس از برآش داده‌ها در مدل نشان داده شده در حالی که در ستون مشاهده شده (observed) همبستگی رخ داده در بررسی تبیین شده است.

شاخص‌های کیفی سنجش پرسش‌ها به دست آمد و در جداول زیر ارائه شده است. در این بخش، مقایسه یافته‌های هر سه روش صورت پذیرفته و بررسی شده است. از آن رو که استفاده از مدل‌های بررسی کنش افتراقی بر مبنای نظریه کلاسیک آزمون پیش‌فرض‌های خاصی ندارد، نخست، در این بخش به تبیین تناسب مدل و داده‌ها و مفروضات نظریه رش می‌پردازیم. و سپس نتایج یافته‌های هر سه مدل ارائه می‌شود. در جدول شماره یک در ستون اول (item) شماره پرسش در آزمون نشان داده شده و اینکه ترتیب هر پرسش در آزمون برای شرکت‌کنندگان به‌چه صورت بوده است. در ستون دوم (Measure) به میزان و درجه سختی هر پرسش [Lincare](#) در واحد لوجیت (logit) اشاره شده است (لينکر، [Lincare](#)). به طور مثال در این آزمون سخت‌ترین پرسش شماره ۹ بوده و ساده‌ترین پرسش شماره ۳ بوده است. رابطه مستقیمی بین قدر مطلق عدد نشان داده شده با درجه سختی

جدول شماره ۱: پیش‌فرض‌های مدل رش و تناسب داده‌ها و مدل

Item	Measure	Infit MNSQ	Infit ZSTD	Outfit MNSQ	Outfit ZSTD	PMC Observed	PMC Expected	DIF Contrast
۲۰	۰,۶۳	۱,۲۷	۹,۹	۱,۳۶	۹,۹	۰,۱۴	۰,۴۰	-۰,۳۲
۱۰	۰,۱۷	۱,۱۹	۹,۹	۱,۲۵	۹,۹	۰,۲۱	۰,۳۹	-۰,۷۷
۲۲	۰,۸۳	۱,۱۷	۹,۹	۱,۲۴	۹,۹	۰,۲۳	۰,۴۰	-۰,۳۲
۳۴	۰,۴۲	۱,۱۸	۹,۹	۱,۲۴	۹,۹	۰,۲۲	۰,۴۰	-۰,۲۰
۳۲	۰,۴۸	۱,۱۲	۸,۱	۱,۱۷	۸,۰	۰,۲۸	۰,۴۰	-۰,۳۳
۱	۰,۴۵	۱,۱۱	۷,۷	۱,۱۵	۶,۷	۰,۲۹	۰,۴۰	-۰,۲۸
۲۵	۱,۷۷	۱,۰۱	۰,۰۲	۱,۱۳	۳,۳	۰,۳۵	۰,۳۷	۰,۰۲
۱۱	۰,۷۹	۱,۰۹	۵,۶	۱,۱۲	۵,۳	۰,۳۱	۰,۴۰	-۰,۱۳
۱۸	-۰,۹۷	۱,۰۱	۰,۰۵	۱,۱۱	۲,۶	۰,۳۱	۰,۳۴	-۰,۲۹
۱۴	-۰,۷۹	۱,۰۳	۱,۰۶	۱,۱۰	۲,۶	۰,۳۱	۰,۳۵	-۰,۰۴
۳۱	-۰,۶۵	۱,۰۲	۱,۰۱	۱,۰۹	۲,۷	۰,۳۲	۰,۳۶	۰,۰۰
۱۲	-۰,۸۹	۱,۰۲	۰,۹	۱,۰۸	۲,۰	۰,۳۱	۰,۳۴	۰,۰۵
۳۵	-۰,۴۴	۰,۹۸	-۰,۹	۱,۰۴	۱,۴	۰,۳۸	۰,۳۷	۰,۰۹
۲۷	۱,۳۷	۰,۹۷	-۱,۴	۱,۰۴	۱,۲	۰,۴۰	۰,۳۹	۰,۰۶
۳	-۰,۰۲	۱,۰۲	۱,۴	۱,۰۲	۰,۹	۰,۳۷	۰,۳۹	۰,۰۰
۵	۰,۰۲	۱,۰۱	۰,۶	۱,۰۱	۰,۳	۰,۳۹	۰,۴۰	-۰,۳۴
۳۳	۰,۵۳	۱,۰۰	۰,۳	۱,۰۰	۰,۱	۰,۴۰	۰,۴۰	-۰,۰۶
۱۵	۰,۳۳	۱,۰۰	-۰,۳	۰,۹۹	-۰,۵	۰,۴۰	۰,۳۹	۰,۰۴
۱۳	-۱,۱۵	۰,۹۹	-۰,۲	۰,۹۰	-۱,۱	۰,۳۳	۰,۳۲	۰,۰۰
۶	-۰,۲۰	۰,۹۹	-۰,۵	۰,۹۷	-۱,۰	۰,۳۹	۰,۳۸	۰,۱۰

Item	Measure	Infit MNSQ	Infit ZSTD	Outfit MNSQ	Outfit ZSTD	PMC Observed	PMC Expected	DIF Contrast
۹	۱,۶۶	۰,۸۹	-۰,۲	۰,۹۶	-۱,۱	۰,۴۶	۰,۳۸	۰,۰۰
۲۹	۰,۲۵	۰,۹۵	-۳,۸	۰,۹۱	-۴,۲	۰,۴۵	۰,۳۹	-۰,۰۷
۲۳	۰,۳۲	۰,۹۵	-۳,۹	۰,۹۴	-۲,۹	۰,۴۴	۰,۳۹	۰,۱۶
۲۴	۰,۰۸	۰,۹۵	-۳,۸	۰,۹۳	-۲,۹	۰,۴۴	۰,۳۹	۰,۱۱
۷	۰,۱۹	۰,۹۴	-۴,۲	۰,۹۳	-۳,۱	۰,۴۵	۰,۳۹	۰,۱۷
۳۰	۰,۲۶	۰,۹۴	-۴,۶	۰,۹۱	-۴,۳	۰,۴۵	۰,۳۹	۰,۱۳
۱۹	-۰,۷۷	۰,۹	-۳,۲	۰,۹۲	-۲,۴	۰,۴۲	۰,۳۵	۰,۲۳
۲۱	-۱,۱۱	۰,۹۲	-۳,۲	۰,۸۵	-۳,۴	۰,۴۱	۰,۳۳	۰,۲۲
۱۶	-۰,۳۲	۰,۹۱	-۵,۶	۰,۸۵	-۵,۴	۰,۴۶	۰,۳۷	۰,۳۵
۴	-۰,۲۸	۰,۹۱	-۶,۰	۰,۸۴	-۶,۰	۰,۴۷	۰,۳۸	۰,۲۹
۲۶	-۱,۵۰	۰,۹۸	-۳,۷	۰,۷۶	-۴,۶	۰,۴۲	۰,۳۰	۰,۰۵
۱۷	-۰,۳۵	۰,۸۸	-۷,۶	۰,۸۱	-۷,۱	۰,۴۹	۰,۳۷	۰,۰۰
۲	-۰,۶۱	۰,۸۷	-۶,۸	۰,۷۸	-۷,۰	۰,۴۹	۰,۳۶	۰,۶۷
۸	-۰,۲۲	۰,۸۶	-۹,۰	۰,۸۱	-۷,۷	۰,۵۱	۰,۳۷	۰,۱۸
۲۸	-۰,۸۶	۰,۸۶	-۶,۵	۰,۷۶	-۶,۷	۰,۴۸	۰,۳۴	۰,۲۱

کنش افتراقی هر پرسش (DIF size) سوال‌های شماره ۲ و ۱۰ انتخاب شدند. علی‌رغم وجود میزان احتمال معنا داری برای برخی از پرسش‌ها نظری ۵ و ۱۶ و ۲۶ به دلیل آنکه سایز کنش افتراقی معنی دار نبود آنان به عنوان پرسش‌های با سطح قابل چشم پوشی از کنش افتراقی نشان داده شدند و دسته بنده گردیدند.

می‌توان نتیجه گرفت که علی‌رغم بهره‌بری از شیوه‌های مختلف و روش‌های آماری متفاوت، هر سه شیوه در پژوهش حاضر در مورد نشانه‌گذاری پرسش‌ها با سطح کنش افتراقی متوسط و بالا مشترکات زیادی داشته و عملکرد یکسانی از خود ارائه دادند. در باره انتخاب و تشخیص پرسش‌ها با سطح اندک و قابل چشم پوشی از نظر کنش افتراقی به دلیل آنکه تست‌های آماری متفاوتی در هر شیوه به کار برده شده است یافته اندکی از هم تفاوت دارد و نیازمند مبنای نظری برای تفسیر است. آنچه مبرهن و واضح است این است که قدرت هر سه روش از منظر تشخیص مشابه و تا حد زیادی یکسان بوده است.

سپس در گام بعدی برای مقایسه قدرت تشخیص پرسش‌ها دارای کنش افتراقی از سه شیوه گفته شده در بخش پیشین بهره برده شد و نتایج آن در جدول شماره ۲ نشان داده شده است. در بحث بررسی کنش افتراقی با شیوه متنل-هنزل (Mantel-Haenszel chi square) و متنل-هنزل لور یا الگوریتم طبیعی آلفا (MH) و خطای استاندارد متنل-هنزل (LOR SE) سنجیده شد و بر اساس سطح معنادار بودن نشانگر MH CHI وجود کنش افتراقی بررسی و دسته‌بنده شد. در این بررسی پرسش‌های شماره ۲ و ۱۰ دارای کنش افتراقی با سطح متوسط تشخیص داده شدند.

در بررسی مشابه، با استفاده از رگرسیون لجستیک بر اساس فرضیه یک بعدی بودن داده‌ها (unidimensionality)، شرکت‌کنندگان در دو گروه آزمایش و مرجع با یکدیگر تطبیق داده شدند و پرسش‌ها مشابه یافته‌های روش متنل-هنزل مجدد به عنوان پرسش‌هایی که کنش افتراقی از خود نشان داده‌اند انتخاب شدند. سرانجام، با بررسی داده‌ها از راه مدل رش و بررسی میزان احتمال وجود کنش افتراقی (Welch Prob.) و اندازه

جدول شماره ۲: مقایسه توانش مدل‌های تحلیل در تشخوصی کنش افتراقی پرسش‌ها

سوال	نمانگر	متلب هنوز							رگرسیون لجستیک			مدل رش	
		ETS	Welch Prob.	DIF size	J & G	3 <sup>rd</sup> R <sup>2</sup>	2 <sup>nd</sup> R <sup>2</sup>	1 <sup>st</sup> R <sup>2</sup>	DIF χ <sup>2</sup>	ETS	MH CHI	MH LOR	LOR SE
۱	A	۰,۰۰۴	-۰,۱۳	A	۰,۱۰	۰,۱۰	۰,۱۰	۴,۰۶۶	A	۰,۰۷۸۲	۰,۱۳۴۷	۲,۸۵۱	
۲	B	۰,۰۰۰	۰,۵۶	B	۰,۳۷	۰,۳۷	۰,۳۶	۳۳,۳۳۴	B	۰,۰۹۶۴	۰,۵۶۵۵	۳۴,۴۳۴	
۳	A	۱,۰۰۰	۰,۰۶	A	۰,۱۸	۰,۱۸	۰,۱۸	۳,۶۲۲	A	۰,۰۸۲۲	۰,۰۴۶۱	۰,۲۶۹۷	
۴	A	۰,۰۰۵	۰,۲۰	A	۰,۳۲	۰,۳۲	۰,۳۱	۵,۲۰۸	A	۰,۰۸۹۲	۰,۱۹۲۸	۴,۴۷۹	
۵	A	۰,۰۰۰	-۰,۳۵	A	۰,۲۰	۰,۲۰	۰,۲۰	۲۲,۷۲۵	A	۰,۰۸۳۱	۰,۳۴۰۰	۱۶,۴۱۴	
۶	A	۰,۲۰۴۴	۰,۱۰	A	۰,۲۱	۰,۲۱	۰,۲۱	۲,۷۵۹	A	۰,۰۸۴۵	۰,۱۰۷۶	۱,۵۱۵	
۷	A	۰,۰۳۸۰	۰,۰۹	A	۰,۲۷	۰,۲۷	۰,۲۷	**۳,۱۱۳	A	۰,۰۸۴۴	۰,۰۹۷۹	۱,۲۲۲	
۸	A	۰,۰۲۸۱	-۰,۱۰	A	۰,۳۷	۰,۳۷	۰,۳۷	۰,۴	A	۰,۰۹۲۱	۰,۰۰۰۵	۰,۰۰۱۶	
۹	A	۱,۰۰۰	-۰,۰۱	A	۰,۳۰	۰,۲۸	۰,۲۸	۳۰,۶۴۲	A	۰,۱۰۱۱	۰,۰۱۳۶	۰,۰۰۷۰	
۱۰	B	۰,۰۰۰	-۰,۰۱	B	۰,۰۷	۰,۰۷	۰,۰۵	۴۴,۱۷	B	۰,۰۷۸۹	۰,۰۱۱۰	۴۱,۸۱۸	
۱۱	A	۰,۱۲۰۱	-۰,۰۱	A	۰,۱۲	۰,۱۲	۰,۱۲	۵,۷۹۲	A	۰,۰۸۰۳	۰,۰۰۱۲	۰,۰۰۷	
۱۲	A	۰,۵۷۰۳	۰,۱۰	A	۰,۱۵	۰,۱۵	۰,۱۵	۱,۲۶۸	A	۰,۰۹۲۴	۰,۰۹۹۴	۱,۰۶۶۲	
۱۳	A	۱,۰۰۰	۰,۰۱	A	۰,۱۸	۰,۱۸	۰,۱۸	۰,۳۸	A	۰,۰۹۹۸	۰,۰۰۴۲	۰,۰۰۰۱	
۱۴	A	۰,۸۰۵۹	۰,۰۳	A	۰,۱۴	۰,۱۴	۰,۱۴	۰,۰۶۱	A	۰,۰۹۰۷	۰,۰۱۹۰	۰,۰۲۷۲	
۱۵	A	۰,۰۶۱۴۴	۰,۰۳	A	۰,۲۱	۰,۲۱	۰,۲۱	۰,۰۰۹	A	۰,۰۸۲۲	۰,۰۳۵۱	۰,۱۴۹۰	
۱۶	A	۰,۰۰۰	-۰,۲۸	A	۰,۳۱	۰,۳۱	۰,۳۱	۱۰,۰۶	A	۰,۰۹۰۱	۰,۲۷۳۹	۸,۹۴۷۷	
۱۷	A	۱,۰۰۰	-۰,۱۶	A	۰,۳۵	۰,۳۵	۰,۳۵	۷,۷۴۹	A	۰,۰۹۲۶	۰,۱۶۲۹	۲,۹۴۱۲	
۱۸	A	۰,۰۰۰۱۷	-۰,۲۸	A	۰,۱۵	۰,۱۵	۰,۱۵	۸,۱۰۰	A	۰,۰۹۶۲	۰,۲۸۲۱	۸,۳۲۱۲	
۱۹	A	۰,۰۰۹۲	۰,۱۰	A	۰,۲۶	۰,۲۶	۰,۲۶	۲,۴۹۷	A	۰,۰۹۲۹	۰,۱۰۰۹	۲,۴۸۶۶	
۲۰	A	۰,۰۰۰۱	-۰,۱۰	A	۰,۰۱	۰,۰۱	۰,۰۱	۲,۱۰۹	A	۰,۰۷۶۴	۰,۰۰۱۰	۰,۰۰۰۳	
۲۱	A	۰,۰۲۳۶	۰,۰۹	A	۰,۲۷	۰,۲۷	۰,۲۷	۲,۴۱۰	A	۰,۱۰۳۷	۰,۰۷۴۳	۰,۴۳۹۶	
۲۲	A	۰,۰۰۰۱	-۰,۱۰	A	۰,۰۶	۰,۰۶	۰,۰۶	۴,۲۶۳	A	۰,۰۷۸۳	۰,۱۱۱۰	۱,۹۰۹۰	
۲۳	A	۰,۰۴۳۱	۰,۱۱	A	۰,۲۶	۰,۲۶	۰,۲۶	۳,۷۲۳	A	۰,۰۸۴۰	۰,۱۲۴۹	۲,۰۹۰۷	
۲۴	A	۰,۱۷۶۴	۰,۰۵	A	۰,۲۷	۰,۲۶	۰,۲۶	۳,۱۰۵	A	۰,۰۸۰۰	۰,۰۰۹۶	۰,۴۳۳۱	
۲۵	A	۰,۸۰۶۶	۰,۱۲	A	۰,۱۶	۰,۱۶	۰,۱۶	۱,۹۸۲	A	۰,۰۹۴۷	۰,۱۱۰۰	۱,۲۶۸۳	
۲۶	A	۰,۰۰۰	۰,۳۸	A	۰,۳۳	۰,۳۳	۰,۳۲	۹,۹۲۸	A	۰,۱۱۹۶	۰,۳۶۲۳	۸,۹۰۲۳	
۲۷	A	۰,۵۱۳۰	۰,۱۰	A	۰,۲۱	۰,۲۱	۰,۲۱	*۵,۸۰۶	A	۰,۰۹۰۱	۰,۰۸۶۶	۰,۸۳۷۲	
۲۸	A	۰,۰۲۰۷	-۰,۴۰	A	۰,۳۷	۰,۳۷	۰,۳۷	۰,۱۹۸	A	۰,۱۰۲۶	۰,۰۳۰۱	۰,۰۵۸۶	
۲۹	A	۰,۳۷۳۵	-۰,۱۶	A	۰,۲۷	۰,۲۷	۰,۲۷	۸,۱۹۰	A	۰,۰۸۴۶	۰,۱۰۶۴	۳,۲۷۹۵	
۳۰	A	۰,۱۰۴۳	۰,۰۷	A	۰,۲۸	۰,۲۷	۰,۲۷	۱,۸۷۷	A	۰,۰۸۰۱	۰,۰۸۹۰	۱,۰۰۲۱	
۳۱	A	۱,۰۰۰	۰,۰۲	A	۰,۱۶	۰,۱۶	۰,۱۶	۵,۹۸۶	A	۰,۰۸۹۰	۰,۰۱۷۵	۰,۰۲۲۹	
۳۲	A	۰,۰۰۰	-۰,۲۲	A	۰,۱۱	۰,۱۱	۰,۱۱	۱۱,۳۲۷	A	۰,۰۷۹۱	۰,۲۲۱۵	۷,۷۰۱۸	
۳۳	A	۰,۴۵۱۲	-۰,۰۷	A	۰,۲۰	۰,۲۰	۰,۲۰	۲,۰۵۲	A	۰,۰۸۱۸	۰,۰۵۳۵	۰,۳۷۴۶	
۳۴	A	۰,۰۱۰۹	۰,۰۳	A	۰,۰۵	۰,۰۵	۰,۰۵	۱,۱۸۲	A	۰,۰۷۷۳	-۰,۰۲۶۸	۰,۰۹۴۸	
۳۵	A	۰,۲۶۳۲	۰,۰۷	A	۰,۲۱	۰,۲۱	۰,۲۱	۰,۹۶۱	A	۰,۰۸۷۳	۰,۰۷۶۷	۰,۷۹۷۹	

پرسش ها، کنش افتراقی از خود نشان دادند و سایر پرسش ها عملکرد دقیقی در سنجش سازه مدنظر خود داشتند. در رابطه با تاثیر پیشینه و رشتہ تحصیلی نیز باید در نظر گرفت که در بخش دستور زبان مشارکت و اثرگذاری خاصی مشاهده نشد (هیل، Hale، ۱۹۸۸).

## ۵. نتیجه گیری

هدف از این پژوهش بررسی و تبیین قدرت هر یک از شیوه های سنجش و تشخیص کنش افتراقی پرسش ها آزمون بود که پیشتر معرفی شدند (مدل رش، رگرسیون لجستیک، متنل هنzel). هدف دیگر ارائه شده در این پژوهش، اثرگذاری این شیوه ها بر روایی و صحت تحلیل های آماری انجام شده بود که در خلال آنچه عنوان شده بدان نیز اشاره شد.

بر اساس آنچه در جدول شماره ۲ در بحث مقایسه شیوه های سنجش کنش افتراقی مطرح شد در تشخیص پرسش های با سطح کنش افتراقی متوسط و یا بالا (که البته در اینجا چنین پرسش های وجود نداشت) مدل متنل هنzel کارآمدی مناسب از خود ارائه می دهد.

باید درنظر داشت که به طور کلی مدل رگرسیون لجستیک مدل عمومی تر نسبت به مدل متنل هنzel است و در محاسبات متنل هنzel فقط مدلی از کنش افتراقی یافت می شود که در تمام سطوح یک متغیر وجود داشته باشد و نتایج این پژوهش نشان داد که مدل متنل هنzel و مدل رگرسیون لجستیک هر دو به اندازه مدل رش در تشخیص کنش افتراقی توانمندند. اما این مبحث که آیا در بحث تشخیص کنش افتراقی غیر یکنواحت چگونه عمل می کنند، نیازمند بحث و بررسی بیشتر است. از نظر بررسی های آماری و زمان تحلیل، مدل رگرسیون لجستیک از دو مدل دیگر پیچیده تر و زمان برتر می باشد. باید این نکته را عنوان کرد که برای تست های بسیار ساده و یا بسیار سخت فرض بر آن بوده که شیوه متنل هنzel کارآمد می باشد، اما برای تست هایی با میزان سختی متوسط از کارآمدی آن کاسته می شود. این بخش حداقل با توجه به درجه سختی پرسش ها که محاسبه گردید می توان ذکر نمود که این آزمون از درجه سختی متوسطی برخوردار بوده و ناکارآمدی از بابت اعمال شیوه متنل هنzel در آن مشاهده نگردید.

با توجه با یافته های جدول شماره ۲ تا حدی می توان نتیجه گرفت که برای پرسش های ساده و یا پرسش های بسیار سخت شیوه متنل هنzel به صورت سخت گیرانه تری کنش افتراقی پرسش ها را تبیین نموده است. یکی دیگر از یافته های این پژوهش آن است که درصد پرسش ها داری کنش افتراقی در هر سه شیوه میزان برابری دارد. شباهت و توانش دو مدل متنل هنzel و رگرسیون لجستیک را می توان در تشابه مدل آماری زیربنایی آنان یافت، اما تشخیص پرسش های مشابه، توسط مدل رش حاکم از آن است که دو روش دیگر هم توانایی قابل قبولی برای پیدا کردن پرسش های دارای کنش افتراقی دارند (همیلتون و راجرز، Hambleton & Rogers، ۱۹۸۹).

تأثیر وجود کنش افتراقی بر عدالت آزمون را می توان از منظر تفاوت های کارکردی شرکت کنندگان آزمون سنجید. وجود پرسش هایی با ویژگی کنش افتراقی بدان مفهوم است که در مواردی سازه های غیر مرتبطی به پاسخ افراد تاثیر می گذارد و ملاک توانایی شرکت کنندگان آزمون نیست. در آزمون حاضر، به دلیل وجود تعداد اندک پرسش های با ویژگی کنش افتراقی (۲ سوال) می توان بیان کرد که ۵,۷ درصد پرسش های امکان تأثیر پذیری از عوامل نامرتبطی برای شرکت کنندگان وجود دارد که در برخی از موارد به نفع آنان و در برخی اوقات به ضرر آنان می باشد.

پر واضح است که به صورت کلی بخش بسیار اعظمی از آزمون حاضر هیچ گونه تبعیضی را برای شرکت کنندگان قائل نبوده و عدالت آزمون در بیش از ۹۵ درصد پرسش بخش دستور زبان به خوبی وجود داشته است. لیکن باید در نظر داشت که مفهوم عدالت آزمون مفهومی چند بعدی است و وجود یا فقدان آن متأثر از عوامل مختلفی است که از زمان طراحی آزمون باید آنان را کنترل کرد و سرانجام در زمان تفسیر نتایج، اثرگذاری آن عوامل را در نظر گرفت.

در بحث روایی و وجود کنش افتراقی نیز آنچه به دست آمده، حاکم از آن است که به طور کلی بخش دستور زبان جزء بخش هایی از آزمون های زبان است که کمترین میزان کنش افتراقی را نشان می دهد و از روایی خوبی برخوردار است (مکنمارا، رورو، McNamara & Rover، ۲۰۰۶) در همین پژوهش نیز مشاهده می شود که تعداد بسیار اندکی از

پاسخ فراتر از مدل رشن گام برداریم و مدل های دوپارامتری و سه پارامتری را نیز در تحلیل ها لحاظ نماییم نتایج تغییر می یابد و حساسیت نظریه سوال-پاسخ افزایش خواهد یافت. در بحث نرم افزاری نیز ممکن است که بتوان در راستای افزایش دقت سنجش کوشید و با استفاده از نرم افزار های مختلف برای اعمال این مدل ها اطمینان حاصل نمود که یافته برگرفته از اثربخشی خود شیوه ها و مفروضات مدل های بررسی شده است نه اثر نرم افزار های مختلف استفاده شده در این پژوهش. همچنین وابستگی یافته ها به تعداد پرسش ها هر آزمون نیز نکته دیگری است که پیشنهاد می شود در سطوح مختلف مورد بررسی قرار گیرد.

یافته های این پژوهش برای طراحان آزمون، برگزارکنندگان آزمون و به ویژه تحلیل گران نتایج آزمون ها رویکردهای جدیدی در بهره بری از شیوه های مختلف سنجش کنش افتراقی را پیشنهاد می کند. لیکن در تعیین پذیری این یافته ها باید دقت لازم مد نظر قرار داده شود، زیرا مبحث مورد سنجش در پژوهش حاضر که با نام دستور زبان از آن یاد می شود سازه ای است که در تحلیل های آماری به عنوان سازه تک بعدی (unidimensional) شناخته می شود و در تطبیق مفروضات نظریه سوال-پاسخ می باشد. فلذا اگر سازه ای که مساله ای عدم وابستگی مکانی (local independence) شده سنجیده شود شاید منجر به کسب نتایج متفاوتی گردد. همچنین محتمل است که اگر در بحث مدل های نظریه سوال-

## منابع

Acar, T., & Kelecioglu, H. (2010). Comparison of Differential Item Functioning Determination Techniques: HGLM, LR and IRT-LR. *Educational Sciences: Theory and Practice*, 10(2), 639-649.

Ackerman, P. L. (1992). Predicting individual differences in complex skill acquisition: Dynamics of ability determinants. *Journal of applied psychology*, 77(5), 598.

Bachman, L. F. (1990). *Fundamental considerations in language testing*. Oxford: Oxford University Press.

Cameron, I. M., Scott, N. W., Adler, M., & Reid, I. C. (2014). A comparison of three methods of assessing differential item functioning (DIF) in the Hospital Anxiety Depression Scale: ordinal logistic regression, Rasch analysis and the Mantel chi-square

procedure. *Quality of life research*, 23(10), 2883-2888.

Camilli, G. (2006). Test fairness. In R. Brennan (Ed.), *Educational measurement* (4th ed.) (pp. 221-256). New York: American Council on Education & Praeger series on higher education.

Camilli, G., & Shepard, L. (1994). *Methods for identifying biased test items*. Newbury Park, CA Sage.

Chapelle, C. A. (2020). Validity in language assessment. *The Routledge Handbook of Second Language Acquisition and Language Testing*, 11.

Chen, Z., & Henning, G. (1985). Linguistic and cultural bias in language proficiency tests. *Language Testing*, 2(2), 155–163.

- Chen, M. Y., Liu, Y., & Zumbo, B. D. (2020). A propensity score method for investigating differential item functioning in performance assessment. *Educational and Psychological Measurement, 80*(3), 476-498.
- Clauser, E. B. & Mazor, M. K. (1998) Using Statistical Procedures to Identify Differentially Functioning Test Items. *Educational Measurement: Issues and Practice*.17, 31-44.
- Elena Oliveri, M., Lawless, R., Robin, F., & Bridgeman, B. (2018). An exploratory analysis of differential item functioning and its possible sources in a higher education admissions context. *Applied Measurement in Education, 31*(1), 1-16.
- Fischer, H. F., Wahl, I., Nolte, S., Liegl, G., Brähler, E., Löwe, B., & Rose, M. (2017). Language-related differential item functioning between English and German PROMIS Depression items is negligible. *International Journal of Methods in Psychiatric Research, 26*(4), e1530.
- Hale, G. (1988). Student major field and text content: Interactive effects on reading comprehension in the TOEFL. *Language Testing, 5*(1), 49-61.
- Hambleton, R. K., & Rogers, H. J. (1989). Detecting potentially biased test items: Comparison of IRT area and Mantel-Haenszel methods. *Applied Measurement in Education, 2*(4), 313-334.
- Jodoin, M. G., & Gierl, M. J. (2001). Evaluating type I error and power rates using an effect size measure with the logistic regression procedure for DIF detection. *Applied measurement in education, 14*(4), 329-349.
- Karami, H. (2012) An introduction to Differential Item Functioning. *International Journal of Educational and Psychological Assessment, 11*(2), 59-76.
- Karami, H. (2013) The quest for fairness in language testing. *Educational Research and Evaluation, 19*(2&3), 158-169.
- Linacre, J. M. (2010a). *A User's Guide to WINSTEPS®*. Retrieved May 2, 2010 from <http://www.winsteps.com/>.
- Linacre, J. M. (2010b). *Winsteps®* (Version 3.70.0) [Computer Software]. Beaverton, Oregon: Winsteps.com.
- McNamara, T. & C. Roever (2006). *Language Testing: The Social Dimension*. Malden, MA & Oxford: Blackwell.
- Sireci, S. G., & Rios, J. A. (2013). Decisions that make a difference in detecting differential item functioning. *Educational Research and Evaluation, 19*(2-3), 170-187.
- Steinberg, L., & Thissen, D. (2006). Using effect sizes for research reporting: Examples using item response theory to analyze differential item functioning. *Psychological Methods, 11*(4), 402.

Uyar, S., Kelecioğlu, H., & Doğan, N. (2017).

Comparing differential item functioning based on manifest groups and latent classes.

*Kuram ve Uygulamada Eğitim Bilimleri*,  
17(6), 1977-2000.

Zieky, M. (1993). Practical questions in the use of DIF statistics in test development. In P.W. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 337–348). Hillsdale, NJ: Lawrence Erlbaum Associates.

Zumbo, B. D. (2003). Does item-level DIF manifest itself in scale-level analyses? Implications for translating language tests. *Language testing*, 20(2), 136-147.

Zhu, X., & Aryadoust, V. (2020). An investigation of mother tongue differential item functioning in a high-stakes computerized academic reading test. *Computer Assisted Language Learning*, 1-25.