



دانشگاه علامه طباطبائی
دانشکده علوم تربیتی و روانشناسی
گروه سنجش و اندازه‌گیری

مقایسه‌ی کاربرد اندازه‌ی اثر در سنجش عملکرد افتراقی سوال بین دو روش رگرسیون لجستیک و منتل هنزل

استاد راهنما:

دکتر علی دلاور

استاد مشاور:

دکتر فریبوز درتاج

۱۳۸۸/۹/۱۸

دانشجو:

محمد حسین ضرغامی

آمرت اطلاعات آرک عملی بریز
تسبیر برارک

پایان‌نامه جهت دریافت درجه کارشناسی ارشد در رشته سنجش و اندازه‌گیری

شهریورماه ۱۳۸۸

۱۲۶۷۸۶

۲۰۰۴۴۲۵

تقدیم به روح بلند پدر که برایم بهترین
بود و با عمل و رفتار خود درس‌ها به
من آموخت و تقدیم به تمام کسانی که
نیروهای بیکران وجود آدمی را بیدار
می‌کنند تا آدمی با گیتی پیوند بخورد و
به اوج لذتی ابدی برسد.

سپاس خداوندی که اندیشه را برای پی بردن انسان به جهلش در او قرار داد

بسیار ممنون و متشکرم از استاد بزرگ و مهربانم جناب آقای دکتر علی دلاور که زحمت راهنمایی این پایان نامه را متقبل شدند. همچنین از استاد گرامی جناب دکتر فریبرز درتاج متشکرم که علاوه بر مشاوره‌ی این پایان نامه در طول دوره مشوق بسیار خوبی برایم بوده‌اند. بیشتر سختی‌های این پژوهش بر دوش همسر عزیزم بود، از ایشان نیز نهایت تشکر و سپاس را دارم و امیدوارم زندگی همیشه بر وفق مرادش باشد.

نهایت سپاس و تشکر را از گسترش‌دهندگان شبکه‌های اینترنتی دارم که دسترسی به منابع معتبر را برایمان آسان کرده‌اند و در جهت عدالت اطلاعات گام مهمی را برداشته و معیارهای موفقیت را عوض کرده‌اند. سپاسگذار تمام کسانی هستم که مانند خورشید پشت ابر در عین هدایتگری ناشناخته مانده‌اند.

چکیده‌ی پژوهش:

این پژوهش به مقایسه‌ی دو روش غالب در تعیین عملکرد افتراقی سوال می‌پردازد. این دو روش عبارتند از: روش رگرسیون لجستیک و روش متل هنزل. مبنای این مقایسه، آماره‌ی اندازه‌ی اثر می‌باشد. در این مطالعه تنها به بررسی DIF یکنواخت پرداخته شده است. در ضمن پارامترهای مختلف سوال و بزرگی DIF موجود در داده‌های تجربی دستکاری شده و از طریق اطلاعات بدست آمده داده‌های لازم شبیه سازی شده‌اند. نتایج نشان می‌دهد که کاربرد اندازه‌ی اثر دقت و صحت کشف سوالات دارای DIF را در هر دو روش MH DIF و LR DIF بهبود می‌بخشد. این بهبود برای روش MH DIF در مقایسه با روش LR DIF ناچیز است.

واژگان کلیدی: عملکرد افتراقی سوال، اندازه اثر، رگرسیون لجستیک، متل هنزل.

فهرست مطالب

فصل اول (کلیات پژوهش)..... ۱

مقدمه..... ۲

تعاریف نظری و عملیاتی..... ۷

ضرورت پژوهش و بیان مساله..... ۸

اهداف پژوهش..... ۹

سوالات پژوهش..... ۹

فصل دوم (پیشینه پژوهش)..... ۱۰

مروری بر ادبیات پژوهش..... ۱۱

نگاه تاریخی به سوگیری سوال..... ۱۴

عملکرد افتراقی سوال (DIF)..... ۱۷

تکنیک های ابتدایی سوگیری سوال..... ۱۸

روش های DIF مبتنی بر نظریه IRT..... ۱۹

برازش و ارزیابی گروه اقلیت از طریق برآوردهای گروه کل..... ۲۳

خلاصه ی روش های DIF مبتنی بر IRT..... ۲۴

روش های DIF خنثی دو..... ۲۵

رگرسیون لجستیک به عنوان یک روش در کشف DIF..... ۳۰

محدودیت های روش های DIF..... ۳۱

آزمون های معناداری و اندازه ی اثر..... ۳۴

- اندازه اثر ۴۸
- انواع اندازه اثر ۳۶
- کاربرد اندازه اثر در روش شناسی DIF ۴۲
- نتایج مربوط به پژوهش حاضر ۴۴

فصل سوم (روش پژوهش) ۴۶

- داده های نمونه ۴۷
- جامعه ی داده ها ۴۸
- محاسبه ی فرمول مساحت راجو برای کشف سوالات با عملکرد افتراقی ۵۲
- DIF یکنواخت و DIF غیر یکنواخت ۶۱
- ساخت داده های شبیه سازی شده ۶۲
- تحلیل های آماری ۶۳

فصل چهارم (تحلیل آماری) ۶۶

- سوال اول پژوهش ۶۷
- سوال دوم پژوهش ۷۲
- سوال سوم پژوهش ۷۵

فصل پنجم (بحث و نتیجه گیری) ۷۷

- سوال اول پژوهش ۷۸
- سوال دوم پژوهش ۷۸
- سوال سوم پژوهش ۷۹
- نتیجه گیری کلی ۷۹

محدودیت ها و پیشنهادات ۷۹

منابع ۸۱

پیوست ها ۸۵

فصل اول

کلیات پژوهش

آزمون‌های استاندارد شده^۱ و سایر اندازه‌گیری‌ها در حوزه علوم رفتاری ابتدا به منظور تمایز بین مهارت‌ها و سطوح مختلف توانایی مورد استفاده قرار می‌گرفتند. بعضی از این حوزه‌ها عبارت‌اند از: توانایی‌های دانشگاهی، پژوهش‌های شغلی و ویژگی‌های شخصیتی. عملکرد افتراقی سوالات^۲ نیز ابتدا به عنوان بخشی از روایی^۳ مدنظر بود. وظیفه‌ی عملکرد افتراقی سوالات در پژوهش‌های روایی تعیین درجه‌ی تفاوت توانایی‌های واقعی آزمودنی‌ها در رفتاری غیر سودار است. آزمون سازان و متخصصان روان‌سنجی از عملکرد افتراقی سوال به منظور تعیین امکان وجود سوگیری در یک سوال استفاده می‌کنند، البته محدوده‌ی کاربرد عملکردهای افتراقی سوال از شناسایی سوالات دارای سوگیری بالقوه در امتحانات مدرسه و آزمون‌های استخدامی تا کاربرد تحلیل‌های پیشرفته نظریه‌ی سوال پاسخ به منظور شناسایی پیچیدگی‌های ممکن در زمان تفسیر نتایج پیمایش‌های اجتماعی می‌باشد (راست و گلوبک^۴، ۲۰۰۹). این پژوهش شرایطی را مورد مطالعه قرار می‌دهد که در آن دو روش غالب در تعیین عملکرد افتراقی سوال با یکدیگر مقایسه می‌شوند. مبنای این مقایسه، آماره‌ی "اندازه‌ی اثر"^۵ است. به طور قطع نمرات آزمون از منابع واریانسی غیر از توانایی‌های آزمودنی‌ها، اثر می‌پذیرد. در صورتی که آزمون‌ها به طور کامل آنچه که مدنظر پژوهشگر است را بسنجند، همه نمرات روا و پایا^۶ خواهند بود ولی هرگز نمی‌توان منابع نامربوط واریانس را به طور کامل کنترل کرد، بنابراین باید در استفاده از نتایج آزمون‌ها بسیار دقت شود تا از قضاوت‌های ناعادلانه به نفع یا ضرر گروه یا گروه‌های خاصی اجتناب شود. این قضاوت‌های ناعادلانه زمانی رخ می‌دهند که دو گروه جامعه از لحاظ توانایی‌های مورد سنجش یکسانند ولی منابع نامربوط واریانس برای آنها به طور متفاوتی توزیع شده است (کورویل^۷، ۲۰۰۴).

۱. Standardized tests

۲. Differential Item Functioning (DIF)

۳. validity

۴. Rust & Golombok

۵. Effect size

۶. Valid & Reliable

۷. Courville

روش آماری که در تعیین سوگیری^۱ سوال مورد استفاده قرار می‌گیرد، تحت عنوان عملکرد افتراقی سوال (DIF) نامیده می‌شود. DIF مشاهده‌ی ویژگی‌های آماری متفاوت سوال در دو زیر جامعه است که فرض می‌شود سطح توانایی مورد سنجش در آنها برابر است (هلند و واینر^۲، ۱۹۹۳). DIF یک فرایند دو مرحله‌ای است: گام اول مقایسه‌ی پرونده‌های زیر گروه‌ها در یک سوال و تعیین وجود یا عدم وجود DIF در آن سوال می‌باشد. گام دوم تصمیم‌گیری در مورد این موضوع است که آیا تفاوت بین گروه‌ها (زیر جامعه^۳ ها) به اندازه‌ی کافی بزرگ است که بتوان سوال را حذف کرد و یا تغییر داد. مرحله‌ی دوم در برگیرنده‌ی آزمون‌های معناداری آماری مربوط به DIF است (فن^۴، ۱۹۹۸).

محدودیت اساسی در استفاده از آزمون‌های معناداری آماری^۵ در هر دو مرحله‌ی تعیین DIF و برای سوالی که سوگیری ندارد، حجم بالای نمونه است که ممکن است منجر به خطای نوع اول یا مثبت غلط^۶ می‌شود (تامپسون^۷، ۲۰۰۲). بیشتر روش‌های تعیین DIF، آزمون‌های معناداری آماری را به منظور ارزیابی DIF به کار می‌گیرند. کاربرد ارزش‌های احتمالی (ارزش p)^۸ و آماره χ^2 (خی دو) در برابر حجم‌های متفاوت نمونه مقاوم^۹ نیستند. بنابراین چون حجم نمونه^{۱۰} در مطالعات مختلف یکسان نمی‌باشد، مقایسه‌ی آنها با یکدیگر امکان‌پذیر نیست. ناتوانی در مقایسه‌ی مطالعات مختلف یک ضعف محسوب شده که مانع تعمیم یک مطالعه می‌شود (کیرک^{۱۱}، ۲۰۰۱). استفاده از اندازه‌ی اثر به منظور کمی‌سازی^{۱۲} DIF تنها

^۱ . Bias

^۲ . Holland and wainer

^۳ . Subpopulation

^۴ . Fan

^۵ . Test of statistical significance

^۶ . False positive

^۷ . Thompson

^۸ . P value

^۹ . Robust

^{۱۰} . Sample size

^{۱۱} . Kirk

^{۱۲} . Quantitative

بزرگی^۱ مقدار DIF را در بر نمی‌گیرد بلکه علاوه بر بزرگی، تعمیم‌پذیری^۲ و تکرارپذیری^۳ آن را نیز شامل می‌شود (هیوبرتی^۴، ۲۰۰۲). بنابراین نیاز به استفاده از اندازه‌ی اثر به عنوان یک آماره‌ی مکمل و به منظور کنترل خطای نوع اول به شدت احساس می‌شود (کیرک، ۱۹۹۶) علاوه بر این زمانی که در یک سوال DIF وجود دارد باید مقداری کمی از آن ارائه داد (پتنزا و دورانز^۵، ۱۹۹۴). دو محدودیت دیگر در روش شناسی^۶ DIF که بر کشف DIF اثر می‌گذارند، توزیعات مختلف توانایی و توزیعات متفاوت جامعه‌ی داده‌ها است. ممکن است مفروضه‌ی نرمال بودن توزیع صفت در جامعه که توسط بسیاری از پژوهشگران و متخصصان روان‌سنجی ارائه شده است، برآورده نشود. پیرسون^۷ (۱۸۹۵) مساله بسیار مهمی را بیان کرد. وی به بررسی ۴۴۰ مجموعه داده واقعی پرداخت تا ویژگی توزیع آنها را مشاهده کند، از این مجموع فقط در حدود ۳۲٪ آنها دارای توزیع نرمال و نسبتاً نرمال بودند (میسری^۸، ۱۹۸۹). سوینی^۹ (۱۹۹۶) معتقد است ناهماهنگی در توزیعات توانایی منجر به ناپایی کشف DIF می‌شود. این در حالی است که مفروضه‌ی هماهنگی توزیعات توانایی در گروه‌های مرجع و هدف یکی از مفروضات اصلی DIF می‌باشد.

از دهه‌ی ۱۹۶۰ که روش‌های DIF به عنوان بخشی از پژوهش‌های اندازه‌گیری مطرح شد (آنگوف^{۱۰}، ۱۹۹۳)، پژوهشگران به ضعف‌های موجود در فرایند کشف DIF پی بردند. یکی از خطاهای بارز و معمول در روش شناسی DIF، خطای نوع اول است. این محدودیت مخصوصاً در روش‌های DIF مبتنی بر نظریه سوال-پاسخ، روش‌های مبتنی بر آماره‌ی خنی دو و بالاخص آزمون‌های معناداری معمول، که براساس آنها

^۱ . Magnitude

^۲ . Generlizability

^۳ . Repeatability

^۴ . Huberty

^۵ . Potenza and Dorans

^۶ . Methodology

^۷ . Pearson

^۸ . Micceri

^۹ . Sweeney

^{۱۰} . Angoff

معناداری تفاوت‌ها تعیین می‌شود، خود را بیشتر نشان می‌دهد. به منظور اصلاح این وضعیت باید روش‌های آماری استفاده شود که نسبت به حجم نمونه مقاوم و خدشه ناپذیر باشند.

در ۴۰ سال گذشته روش‌های متعددی برای کشف DIF گسترش یافته ولی تعداد کمی از این روش‌ها مورد مطالعه دقیق قرار گرفته اند. با پیشرفت روش شناسی DIF روش‌های قدیمی یا کنار گذاشته می‌شدند و یا از طریق روش‌های کامل‌تر و جدیدتر استنتاج می‌شدند و نواقص و محدودیت‌های آنها آشکار می‌شد. روش‌های رایج DIF عبارت‌اند از: روش‌های مبتنی بر مدل‌های IRT، روش‌های مبتنی بر آماره‌ی χ^2 (مثل روش استاندارد شده^۱ و روش متل هنزل^۲) و روش رگرسیون لجستیک^۳. روش‌های مبتنی بر مدل‌های IRT، روش متل هنزل و رگرسیون لجستیک در پژوهش‌های مربوط به DIF مناسب‌ترند (سوامیناتان و رگرز^۴، ۱۹۹۰) و از بین این روش‌ها، روش متل هنزل و رگرسیون لجستیک عملی‌ترند، چون اولاً کاربرد آنها ساده‌تر است و در عین حال برای این روش‌ها آماره‌ی "اندازه اثر"^۵ تعریف می‌شود که می‌توان از طریق آن صحت کشف DIF را افزایش داد. مطالعات شبیه سازی شده^۶ توانایی ترکیب "اندازه اثر" را با روش‌های متل هنزل (روسوس و استوت^۷، ۱۹۹۶) و رگرسیون لجستیک (جودوین و گیرل^۸، ۲۰۰۱) ثابت کرده‌اند. سوامیناتان و رگرز (۱۹۹۰) رگرسیون لجستیک را به عنوان روشی در کشف DIF معرفی کرده‌اند و این روش را با روش‌های MH DIF^۹ مقایسه کرده‌اند. مزور و همکاران^{۱۰} (۱۹۹۵) کاربرد تحلیل رگرسیون لجستیک را در کشف DIF به عنوان "یک روش ماندنی و پا بر جا" می‌دانند.

^۱ . Standardized method

^۲ . Mantel Haenszel

^۳ . Logistic regression

^۴ . Swaminathan & Rogers

^۵ . Effect size

^۶ . Simulation studies

^۷ . Roussos & Stout

^۸ . Jodoin & Gierl

^۹ . Mantel Heanseal differential item functioning (MH DIF)

^{۱۰} . Mazor

DIF می‌تواند تحت دو موقعیت مورد ملاحظه قرار گیرد: موقعیت یکنواخت^۱ و موقعیت غیر یکنواخت^۲. DIF یکنواخت زمانی اتفاق می‌افتد که تعاملی بین سطح توانایی و عضویت در گروه وجود نداشته باشد، برعکس DIF غیریکنواخت زمانی اتفاق می‌افتد که تعامل بین سطوح توانایی و عضویت گروهی وجود داشته باشد (سوامیناتان و رگرز، ۱۹۹۰). DIF رگرسیون لجستیک (LR DIF^۳) با DIF روش متل هنزل (MH DIF) مقایسه شده است. مقایسه نشان می‌دهد که رگرسیون لجستیک هم در DIF یکنواخت و هم در DIF غیریکنواخت بهتر عمل می‌کند (سوامیناتان و رگرز، ۱۹۹۰).

“اندازه‌ی اثر” حداقل مجذورات وزنی R^2 (WLS R^۴) توسط زومبو و توماس^۵ (۱۹۹۶) معرفی شده است. این “اندازه اثر” توسط جودین و گیرل^۶ (۲۰۰۱) و همچنین روسوس و استوت^۷ (۱۹۹۶)، به طور تجربی و از طریق داده‌های شبیه سازی شده مورد آزمایش قرار گرفت. تمرکز جودوین و گیرل بر “اندازه‌ی اثر” مربوط به LR DIF بود. آنها به عنوان راهنما، جدول طبقه‌بندی شده‌ای ارائه کردند که از طریق آن “اندازه اثر” بالا، متوسط و ضعیف (قابل اغماض)^۸ مشخص می‌شد. مقیاس گذاری این راهنما شبیه راهنمای بزرگ، متوسط، کوچک کوهن^۹ (۱۹۹۲) مربوط به اندازه می‌باشد. روسوس و استوت (۱۹۹۶) به مقایسه‌ی MH DIF و SIBTEST^{۱۰} پرداختند و از “اندازه اثر” در کنار آزمون‌های آماری استفاده کردند تا به این وسیله میزان خطای نوع اول را کاهش دهند. روش‌های MH DIF و LR DIF می‌توانند از “اندازه‌ی اثر” در کنار

^۱. Uniform

^۲. Non-Uniform

^۳. Logistic regression DIF

^۴. Weighted least squares

^۵. Zumbo & Thomas

^۶. Jodoin & Geirl

^۷. Rousos & stout

^۸. Negligible

^۹. Cohen

^{۱۰}. Simultaneous Item Bias TEST (SIBTEST)

آزمون‌های آماری بهره‌گیرند. این دو روش به صورت تجربی توسط هیدالگو و لویزپینا^۱ (۲۰۰۴) با هم مقایسه شده‌اند.

تعاریف نظری و عملیاتی

عملکرد افتراقی سوال: DIF مشاهده‌ی ویژگی‌های آماری متفاوت سوال در دو زیر جامعه است که فرض می‌شود سطح توانایی مورد سنجش در آنها برابر است (هلند و واینر^۲، ۱۹۹۳). در این پژوهش DIF بر مبنای فرمول مساحت راجو^۳ بدست می‌آید. در صورتی که مساحت بدست آمده بین دو ICC گروه مرجع^۴ و هدف^۵ بیشتر یا مساوی ۴/ باشد آن سوال به عنوان یک سوال دارای DIF شناخته می‌شود.

اندازه اثر: "اندازه اثر" عنوانی است که به مجموعه‌ای از شاخص‌ها که بزرگی اثر آزمایش را می‌سنجد اطلاق می‌شود. "اندازه اثر" در آمار، مقداری است که رابطه‌ی بین دو متغیر را بیان می‌کند. در آزمایشات علمی علاوه بر این که ما باید از معناداری آماری باخبر باشیم؛ باید از اثرات مشاهده شده نیز مقداری کمی داشته باشیم. برای تصمیم‌گیری در موقعیت‌های عملی "اندازه اثر"، شاخص بسیار مناسبی است (گریسون و کیم^۶، ۲۰۰۵). در این پژوهش دو اندازه اثر وجود دارد یکی متعلق به روش رگرسیون لجستیک است و دیگری متعلق به روش متل هنزل. "اندازه اثر" حداقل مجذورات وزنی R^2 (WLS R) که توسط زومبو و توماس^۷ (۱۹۹۶) معرفی شده است، به عنوان اندازه اثر روش رگرسیون لجستیک مد نظر است و لگاریتم نسبت شانس معرفی شده توسط زویک و اریکان (۱۹۸۹) به عنوان اندازه اثر روش متل هنزل قلمداد می‌شود.

گروه مرجع و گروه هدف: در مطالعات عملکرد افتراقی سوال دو گروه مد نظر است گروه مرجع و گروه هدف. گروهی که سوال نسبت به آنها سوگیری دارد گروه هدف است و بر عکس گروهی که سوال به نفع

^۱. Hidalgo & Lopezpina

^۲. Holland and wainer

^۳. Raju's Area Formula

^۴. Reference group

^۵. Focal group

^۶. Grissom & Kim

^۷. Zumbo & Thomas

آنها است گروه مرجع تشخیص داده می‌شود. در این مطالعه گروه زنان به عنوان گروه هدف و گروه مردان به عنوان گروه مرجع می‌باشند.

ارزش P: سطح معناداری که بر مبنای آن معناداری آماری مشخص می‌شود. در این تحقیق این مقادیر از طریق نرم افزار SAS بدست می‌آید.

ضرورت پژوهش و بیان مساله

در حال حاضر اکثر روش‌های کشف DIF برای تعیین وجود DIF در یک سوال بر ارزش‌های P متکی‌اند. هر دو روش MH DIF و LR DIF "اندازه‌ی اثر"ی دارند که می‌تواند به عنوان مکمل در کنار آزمون‌های آماری معمول (مبتنی بر ارزش p) قرار بگیرند (زومبو و توماس^۱، ۱۹۹۶). "اندازه‌ی اثر" می‌تواند صحت کشف DIF را افزایش دهد چون زمانی که با حجم‌های مختلف نمونه مواجه‌ایم "اندازه‌ی اثر" پایتار از آزمون‌های آماری معناداری رایج می‌باشد. آزمون‌های آماری در برابر حجم نمونه مقاوم نیستند، فرض بر این است که کاربرد "اندازه‌ی اثر" منجر به کشف دقیق‌تر و صحیح‌تر DIF موجود در یک سوال می‌شود (فینچ و سایرین^۲، ۲۰۰۱).

یک مطالعه (هیدالگو و لویزینا، ۲۰۰۴) موجود است که روش‌های MH DIF و LR DIF را مقایسه کرده است. نتیجه این مطالعه نشان می‌دهد زمانی که از LR DIF استفاده می‌کنیم سوالات دارای DIF نسبت به زمانی که از MH DIF استفاده می‌کنیم، بیشتر است و در عین حال روش LR DIF به شرایط خاصی در کشف DIF حساسیت ندارد. شرایط خاصی که در روش شناسی DIF و برای انواع مختلف آن (یکنواخت و غیریکنواخت) مدنظر اند عبارت‌اند از: حداقل حجم ۱۰۰۰ برای گروه هدف و مرجع و توزیع نرمال توانایی. البته بعضی از پژوهشگران حجم‌های کمتری را پیشنهاد می‌کنند و مناسب می‌دانند (هیدالگو و لویزینا، ۲۰۰۴).

زومبو و توماس (۱۹۹۶) که حداقل مجذورات وزنی^۳ R^2 را به عنوان "اندازه اثر" روش LR DIF مطرح کردند، به سایر پژوهشگران پیشنهاد کردند مطالعات بیشتری برای آزمون دقت و صحت "اندازه اثر" در تعیین DIF انجام دهند. این دو پژوهشگر پایایی این دو روش را در موقعیت‌های مختلف بررسی نکردند و

^۱ . Zumbo & Thomas

^۲ . Finch

^۳ . Weighted least squares R^2

همچنین مطالعه‌ای راجع به این که آیا "اندازه‌ی اثر" مربوط به هر روش، دقت کشف DIF را افزایش می‌دهد، انجام ندادند. بنابراین لازم است در این زمینه مطالعاتی صورت پذیرد. پژوهش حاضر در این راستا گام برمی‌دارد. در این پژوهش "کاربرد اندازه‌ی اثر در سنجش عملکرد افتراقی سوال بین دو روش رگرسیون لجستیک و متل هنزل" مقایسه می‌شود.

اهداف پژوهش

هدف اصلی این پژوهش یک هدف بنیادی است که منجر به گسترش دانش در حوزه‌ی سنجش عملکرد افتراقی سوال می‌شود. کاربرد و استفاده از ارزش P به عنوان معیاری در تشخیص سوالات دارای DIF مخدوش و بی‌ارزش است (تامپسون و کیفر^۱، ۲۰۰۰) بنابراین باید آماره‌ای جایگزین آن شود که مشکلات آن را نداشته باشد. در ضمن نیاز است تا قدرت اندازه‌ی اثر در کشف سوالات دارای DIF در روش‌های مختلفی که از این آماره استفاده می‌کنند مطالعه شود تا دریابیم در کدام روش در صورت استفاده نکردن این آماره نتایج نامعتبر است. به این ترتیب قضاوت‌هایی که در حوزه‌ی سنجش عملکرد افتراقی سوال بر مبنای دو روش رگرسیون لجستیک و روش متل هنزل می‌شود، اصلاح خواهد شد.

سوالات پژوهش

در این پژوهش سه سوال عمده و اساسی مطرح است:

سوال اول: آیا کشف DIF در یک سوال هنگامی که از "اندازه اثر" مربوط به LR DIF یعنی حداقل مجذورات وزنی R^2 (WLS R^2) به عنوان مکمل آزمون‌های آماری معناداری (P Value) استفاده می‌شود، از دقت بیشتری برخوردار است؟

سوال دوم: آیا کشف DIF در یک سوال هنگامی که از "اندازه اثر" مربوط به MH DIF یعنی لگاریتم نسبت شانس^۲ (Log odds ratio) به عنوان مکمل آزمون‌های آماری معناداری (P Value) استفاده می‌شود، از دقت بیشتری برخوردار است؟

سوال سوم: کاربرد اندازه‌ی اثر در کشف سوالات دارای DIF در روش رگرسیون لجستیک موثرتر است یا در روش متل هنزل؟

^۱ . Thompson & keiffer

^۲ . Log odds ratio

فصل دوم

پیشینه پژوهش

مروری بر ادبیات پژوهش

واژه‌ی عملکرد افتراقی سوال (DIF) در آغاز دهه‌ی ۱۹۹۰ وارد ادبیات سنجش و اندازه‌گیری شد. البته قبل از این دهه، و از دهه ۱۹۶۰ واژه‌ای تحت عنوان "سوگیری سوال"^۱ مطرح بود که هدف آن گسترش کاربرد منصفانه و تخفیف اختلاف در عملکرد آزمون بین گروه‌های جامعه (به عنوان مثال سیاهان و اسپانیایی تبارها)^۲ بود (آنگوف، ۱۹۹۳).

بیشتر این مطالعات برای درک اختلاف نمرات زیر گروه‌های جامعه و مخصوصاً برای اثبات این موضوع هدایت می‌شدند که تفاوت‌های مشاهده شده در نمرات، بیشتر ناشی از سوگیری سوالات می‌باشد و نه اختلاف در توانایی واقعی افراد. به طور کلی تمرکز اصلی مطالعات سوگیری سوال، کشف سوالاتی از آزمون است که ممکن است به نفع یا ضرر یک یا چند گروه جانبداری کرده باشند. سوگیری سوال به صورت زیر تعریف شده است (آنگوف، ۱۹۹۳):

"یک سوال زمانی سوگیری دارد که احتمال پاسخگویی برای افرادی با توانایی‌ها و مهارت‌های یکسان، برابر نباشد."

شپارد و سایرین^۳ (۱۹۸۱) ویژگی سوگیری را نوعی ناروایی^۴ می‌دانند که به یک گروه بیشتر از گروه دیگر ضربه وارد می‌کند. "آنگوف (۱۹۹۳) معتقد است تعاریف مربوط به سوگیری نهایتاً مطلبی را درباره‌ی عملکرد، ارزشیابی، اجرا و یا قضاوت‌های ناعادلانه ارائه می‌دهند.

آیا هر تفاوت موجود در عملکرد افراد را می‌توان به سوگیری سوالات نسبت داد؟ با این سوال مطالعات مربوط به سوگیری سوالات به چالش کشیده شد. ارتباط دادن تفاوت‌های موجود در عملکرد زیرگروه‌ها به

^۱ . Item Bias

^۲ . Hispanic

^۳ . Shepard et al

^۴ . Invalidity

ویژگی‌های سوال (مانند محتوی^۱، واژه بندی^۲ و غیره) و یا به تفاوت در توانایی‌های واقعی آنها ذاتا امری دشوار است (آنگوف، ۱۹۹۳).

در زبان محاوره سوگیری و عادلانه نبودن خیلی به هم شبیه‌اند؛ ولی در زمینه‌ی اندازه گیری و تصمیم‌گیری، سوگیری و انصاف مفاهیم کاملا متفاوتی‌اند. سوگیری یک ویژگی آماری از نمرات آزمون می‌باشد. به همین دلیل گفته می‌شود وقتی سوگیری وجود دارد که خطاهای منظم ایجاد کند. در مقابل منصفانه بودن^۳ به داوری درباره‌ی تصمیم‌گیری یا هر فعالیتی که براساس نتایج نمرات آزمون است، اشاره دارد. سوگیری یک ویژگی اعداد است در حالی که عادلانه بودن یک ویژگی تصمیمات یا ویژگی اشخاص است. اعداد اخلاقا و منطقا خنثی می‌باشند ولی تصمیمات و فعالیت‌های اشخاص می‌تواند عادلانه و یا ناعادلانه باشند. اعداد می‌توانند درست یا نادرست باشند ولی نمی‌توانند عادلانه یا ناعادلانه باشند. عادلانه یا ناعادلانه بودن یک آزمون به هدف و خواسته فرد از آزمون بستگی دارد نه صرفا به نمرات و اعداد که از آزمون بدست می‌آیند (آنگوف، ۱۹۹۳).

به طور کلی تفاوت‌های سوگیری و عادلانه بودن^۴ در جدول مشاهده می‌گردد.

سوگیری	عادلانه بودن
معطوف به نمرات آزمون یا پیش‌بینی‌ها براساس آزمون می‌باشد.	معطوف به تصمیماتی است که گرفته می‌شود.
معطوف به ویژگی‌های آماری نمرات است.	داوری در مورد پیامدها است.
به صورت تجربی تعریف می‌شود.	براساس اصطلاحات فلسفی و سیاسی تعریف می‌شود.
به روش علمی تعریف می‌شود.	نمی‌تواند به روش علمی تعریف و تبیین شود.

^۱. Content

^۲. Wording

^۳. Fairness

^۴. Fairness

بطور کلی دو رویکرد برای اندازه گیری سوگیری وجود دارد (زومبو و توماس، ۱۹۹۶).

الف) قضاوت کارشناسان و ب) روش‌های آماری

باید تأکید کرد که کشف DIF در یک سؤال لزوماً به معنای سوگیر بودن آن سؤال نیست (یعنی جانبداری آن به نفع یک گروه و به ضرر گروه دیگر). از مفهوم جانبداری چنین به ذهن می‌آید که مقایسه‌ای هنجاری و رقابت آمیز با دیگر افراد انجام می‌شود، مثلاً در مواردی که نمره نسبتاً بالا باعث کسب نوعی پاداش و نمره نسبتاً پایین موجب بی‌بهرگی از آن می‌شود. اما همیشه چنین نیست. مثلاً در مواردی که دو گروه فرهنگی یا زبانی با ابزار واحدی اندازه گیری می‌شوند و هدف، صرفاً مقایسه‌ی زیانشناختی آن دو گروه است، حتی در مواردی که در یک سؤال DIF به وضوح مشاهده می‌شود ممکن است آن سؤال کاملاً عادلانه و بدون جانبداری باشد. این مطلب به هدف آن سؤال نیز مربوط می‌شود. برای نمونه اگر سوالی راجع به حرارت لازم برای پختن کیک باشد، چنین سؤالی حتی اگر از لحاظ آماری کاملاً به ضرر مردان سوگیر باشد، در صورتی که در یک آزمون گزینش آشپز و شیرینی‌پزی به کار رود کاملاً عادلانه خواهد بود. مثال دیگر این که دولیتل^۱ و کلیری^۲ مشاهده کردند که دانش آموزان دختر به سؤال‌های هندسه و استدلال ریاضی در مقایسه با دانش آموزان پسر در همان سطح توانایی ریاضی، ضعیف‌تر پاسخ می‌دهند (آنگوف، ۱۹۹۳). پس اگر بخواهیم از نتایج آماری به دست آمده، متابعت کنیم باید چنین سؤال‌هایی را از تست خارج سازیم. اما چنین کاری کاملاً بی‌معنی است چرا که فهم هندسه و توانایی استدلال ریاضی اهداف موجه و لازم آموزش ریاضیات است و از دانش آموز انتظار می‌رود که مسائلی از این دست را حل کند. حذف چنین سوال‌هایی نه تنها تست را ناقص می‌کند بلکه به ضرر دانش آموزانی تمام می‌شود که نمی‌توانند به چنین سؤال‌هایی پاسخ دهند چرا که مشخص نمی‌سازد نقص دانش آموز در کجاست. دیگر این که در تحقیقات مربوط به سوگیری بسیار مشاهده شده که سؤال‌هایی دارای DIF شناخته می‌شوند، اما هیچ‌گونه شواهدی دال بر منشاء وجود این کارکرد افتراقی به دست نمی‌آید. باید خاطر نشان ساخت که صرف وجود تفاوت عملکرد دو گروه در یک سؤال به خودی خود دلیلی کافی بر وجود سوگیری به نفع گروهی که عملکرد قوی‌تری داشته، نیست. زیرا همان‌طور که می‌دانیم بخش اعظم تفاوت افراد و گروه‌ها در تست‌های شناختی، ناشی از ماهیت کیفیت و میزان تعلیم و تربیتی است که چه در مدرسه و چه در خارج از آن به آنها عرضه می‌شود. پس سؤال‌هایی را که کارکرد افتراقی دارند چگونه باید تفسیر و تعبیر کرد؟ مقدار DIF می‌تواند به این معنا باشد که سؤال مزبور در یک گروه علاوه بر سازه مورد نظر تست، سازه

^۱ - Doolittle

^۲ - Cleary