

## معرفی نظریه تعمیم‌پذیری و تشریح فرآیند یک مطالعه اندازه‌گیری جهت سنجش اعتبار

نورعلی فرخی<sup>۱</sup>، لیلا بهرامی<sup>۲</sup>

تاریخ دریافت: ۹۵/۰۳/۲۱

تاریخ پذیرش: ۹۵/۱۰/۱۵

### چکیده

شناسایی و جداسازی منابع چندگانه خطای اندازه‌گیری و برآورد هر یک از آنها، تمایز گذاشتن میان تصمیم‌های نسبی و مطلق، تمایز گذاشتن میان رویه‌های اندازه‌گیری ثابت و تصادفی و همچنین پرداختن به طرح‌های مختلف مطالعه D را می‌توان از جمله نقاط قوت نظریه تعمیم‌پذیری ذکر کرد که نظیر آنها در نظریه کلاسیک آزمون وجود ندارد. با وجود این نقاط قوت، نظریه تعمیم‌پذیری برای محققان کشورمان ناشناخته است و تحقیقات انگشت‌شماری در این زمینه صورت گرفته است. هدف از این مقاله، معرفی نظریه تعمیم‌پذیری و نشان دادن قابلیت کاربرد عملی آن در سنجش اعتبار اندازه‌ها بود. علاوه بر مقایسه میان دو نظریه کلاسیک آزمون و تعمیم‌پذیری، چارچوب مفهومی نظریه تعمیم‌پذیری به سادگی بیان گردید. همچنین، در این مقاله فرآیند طراحی، تحلیل و تفسیر یک مطالعه اندازه‌گیری در قالب یک مثال به همراه محاسبات و معادلات مربوطه، در ۱۵ گام اصلی با جزئیات تشریح شد تا برای محققان و آزمون‌سازان یک راهنمای کاربردی جهت سنجش اعتبار باشد. سودمندی GT نسبت به CTT جهت برآورد اعتبار به خصوص در وضعیت‌های اندازه‌گیری پیچیده نمایان گردید. همچنین، GT محققان را قادر می‌سازد که با اقدامات بهینه‌سازی، سهم خطا را در طرح اندازه‌گیری‌شان کاهش دهند که این کار، افزایش دقت در تعمیم نتایج را به دنبال خواهد داشت.

۱. دانشیار گروه سنجش و اندازه‌گیری، دانشگاه علامه طباطبائی، تهران (نویسنده مسئول)

farrokhinoorali@yahoo.com

۲. کارشناسی ارشد گروه سنجش و اندازه‌گیری، دانشگاه علامه طباطبائی

## واژگان کلیدی: اعتبار، نظریه تعمیم‌پذیری، نظریه کلاسیک اندازه‌گیری

## مقدمه

عوامل مختلفی، نمرات مشاهده‌شده آزمون را تحت تأثیر قرار می‌دهد. در نظریه کلاسیک آزمون<sup>۱</sup> (CTT)، همه این عوامل تحت یک عبارت خطای تصادفی نامتمایز در مدل آماری گنجانیده می‌شود. روان‌سنجان برای لحاظ کردن این عوامل و تفکیک آن‌ها در مدل‌های آماری‌شان، نظریه تعمیم‌پذیری<sup>۲</sup> (GT) را گسترش دادند.

GT، یک نظریه آماری درباره‌ی اتکاپذیری اندازه‌های رفتاری است (شولسون و وب<sup>۳</sup>، ۱۹۹۱). بنا به گفته‌ی برنان<sup>۴</sup> (۱۹۹۲)، باینکه ریشه‌های GT را می‌توان در CTT و تحلیل واریانس (آنوا)<sup>۵</sup> یافت، با وجود این، چارچوب مفهومی GT منحصر به فرد است. GT محقق را قادر می‌سازد که منابع چندگانه‌ی خطا را در یک روش اندازه‌گیری شناسایی، تفکیک و برآورد کند و بدین طریق، CTT را به نحو چشمگیری گسترش داده است (شولسون، وب و رولی<sup>۶</sup>، ۱۹۸۹؛ برنان، ۲۰۰۱). GT با در نظر گرفتن همزمان منابع چندگانه‌ی خطا، اجازه می‌دهد مدل دقیق‌تری از وضعیت اندازه‌گیری ساخته شود که در این صورت، نتایج به دست آمده را با دقت بیشتری می‌توان به سایر موقعیت‌های اندازه‌گیری تعمیم داد (شولسون و وب، ۱۹۹۱؛ ون لیون<sup>۷</sup>، ۱۹۹۷؛ برنان، ۲۰۰۱). این در حالی است که CTT در برآورد اعتبار<sup>۸</sup> وضعیت‌های اندازه‌گیری که شامل منابع چندگانه‌ای از خطاست، کارایی ندارد. باین‌حال، سهولت کاربرد و درک روشن از CTT و همچنین، پیچیدگی مفهومی GT در مقایسه با CTT، باعث شده است که CTT تا به امروز در سنجش اعتبار همچنان به قوت خود باقی بماند. بررسی مطالب مجلات علمی معتبر و پایان‌نامه‌های دانشجویان

1. Classical Test Theory (CTT)
2. Generalizability Theory (GT)
3. Shavelson & Webb
4. Brennan
5. analysis of variance (ANOVA)
6. Rowley
7. VanLeeuwen
8. reliability

تحصیلات تکمیلی در کشورمان، نشان می‌دهد که محققان و دانشجویان حوزه‌های مختلف از جمله آموزش و پرورش، روان‌شناسی، پزشکی و ... برای برآورد اعتبار داده‌هایشان از رویکردهای کلاسیک (مانند آلفای کرونباخ، باز آزمون، فرم‌های موازی، دونیمه کردن و ...) استفاده می‌کنند و بندرت می‌توان به مقاله‌ای اشاره کرد که از طرح‌های اندازه‌گیری<sup>۱</sup> GT جهت سنجش اعتبار استفاده کرده باشد؛ بنابراین، به دلیل اهمیت موضوع، در مقاله‌ی حاضر سعی شده است ضمن معرفی GT به محققان، قابلیت این نظریه در سنجش اعتبار اندازه‌ها نشان داده شود. بدین منظور، ابتدا تاریخچه‌ی GT بیان می‌گردد و بعد از آن، خلاصه‌ی مفیدی از شباهت‌ها و تفاوت‌های موجود میان CTT و GT مطرح می‌شود که به محققان کمک خواهد کرد، درک مفهومی خود را از دو نظریه مذکور گسترش دهند و علاوه بر این، به مزیت‌های GT پی ببرند. سپس، مفاهیم و رویکردهای GT تک متغیری<sup>۲</sup> تشریح می‌شود تا محققان بتوانند با چارچوب مفهومی GT بهتر ارتباط برقرار کنند. در پایان، به منظور آشنایی با چگونگی انجام یک مطالعه اندازه‌گیری جهت سنجش اعتبار، فرآیند طراحی، تحلیل و تفسیر این مطالعه اندازه‌گیری در قالب یک مثال در بافت GT تک متغیری ارائه می‌گردد.

**تاریخچه GT.** سال‌ها قبل از آن که GT به‌طور رسمی توسط کرونباخ<sup>۳</sup> و همکارانش معرفی شود، نویسندگان مختلفی در این زمینه تلاش کرده بودند. به طوری که طبق گفته‌ی کرونباخ (۱۹۹۱، به نقل از برنان، ۱۹۹۷)، GT را می‌توان یک اثر ترکیبی به تألیف حداقل ۲۴ نویسنده پنداشت. با استفاده از مقاله برنان (۱۹۹۷) که در خصوص تاریخچه‌ی GT نگاشته شده است، شرحی از مهم‌ترین کارها و آثار نویسندگان مختلف در جدول زیر ارائه گردیده است.

1. measurement designs
2. univariate
3. Cronbach

جدول ۱. خلاصه‌ای از تاریخچه‌ی GT در خلال سال‌های ۱۹۹۷-۱۹۳۶	
۱۹۳۶	بارت <sup>۱</sup> با مقاله‌ای با عنوان «تحلیل نمرات امتحانی» اولین کسی بود که رویکرد آنوا را برای مسائل اندازه‌گیری بکار گرفت.
۱۹۴۱	هایت <sup>۲</sup> نشان داد از طریق آنوا می‌توان اعتبار آزمون را برآورد کرد.
۱۹۴۷	کرونباخ بیان کرد؛ برای مرتفع کردن ناهمسانی‌ها در برخی از برآوردهای اعتبار انواعی از تحلیل‌های چند رویه‌ای <sup>۳</sup> لازم است.
۱۹۵۰	- پیش از سال ۱۹۵۰، آثار غنی در مورد اعتبار از چشم‌انداز CTT وجود داشت که اکثر آن‌ها به طرز فوق‌العاده‌ای در کتاب گالیکسن <sup>۴</sup> (۱۹۵۰) منتشر شد. - در دهه‌ی ۱۹۵۰ محققان در کارهایشان از این موضوع که آنوا می‌تواند به‌طور همزمان رویه‌های چندگانه را بررسی کند، استفاده کردند.
۱۹۵۱	- احتمالاً اولین بررسی اعتبار برحسب مؤلفه‌های واریانس <sup>۵</sup> ، مطالعه‌ی فینلایسون <sup>۶</sup> درباره‌ی نمرات تخصیص داده‌شده به مقالات بود. - ایبل <sup>۷</sup> مقاله‌ای را در مورد اعتبار درجه‌بندی‌ها منتشر کرد که در آن دو نوع واریانس خطا را در نظر گرفته بود؛ یکی شامل اثرات اصلی ارزیاب بود و دیگری بدون آن. همچنین، طرح‌های متقاطع و آشیانه‌ای <sup>۸</sup> یک‌رویه‌ای را در نظر گرفت.
۱۹۵۲	پلینر <sup>۹</sup> ارتباط نظری میان همبستگی‌های بین طبقه‌ای و آنوا را به دست آورد.
۱۹۵۳	لیندکوئیست <sup>۱۰</sup> گزارش جامعی از نظریه چند رویه‌ای ارائه داد که تمرکز آن روی برآورد مؤلفه‌های واریانس در مطالعات اعتبار بود. لیندکوئیست نشان داد که تحلیل‌های چند رویه‌ای منجر به تعاریف متفاوتی از خطا و ضرایب اعتبار نسبت به قبل شده است. بخش لیندکوئیست به‌وضوح بر قسمت‌های مهمی از GT دلالت می‌کرد.
۱۹۵۵	- ۱۹۵۹- در واقع، اثرات اصلی ارزیاب در مقاله‌ی ایبل نقش اثرات اصلی سؤال را در مقالات لرد <sup>۱۱</sup> (۱۹۵۵، ۱۹۵۷، ۱۹۵۹) که درباره‌ی خطاهای استانداردهای اندازه‌گیری شرطی <sup>۱۲</sup> (SEMs) و اعتبار تحت

1. Burt
2. Hoyt
3. multifacet
4. Gulliksen
5. variance components
6. Finlayson
7. Ebel
8. crossed & nested
9. Pilliner
10. Lindquist
11. Lord
12. conditional standard errors of measurement

<p>مفروضات مدل خطای دوجمله‌ای بود، ایفا می‌کرد. علاوه بر این، لرد در مقالاتش با صراحت از آزمون-های تصادفی موازی<sup>۱</sup> نام برد. موضوعاتی که لرد با آن دست‌وپنجه نرم می‌کرد، تأثیر واضحی بر گسترش GT داشت.</p>	
<p>- کارهای ایبل و لرد سرانجام به تمایز میان خطای نسبی و مطلق<sup>۲</sup> انجامید. - در مقاله‌ی بارت (۱۹۵۵) با عنوان «برآورد اعتبار آزمون به شیوه‌ی تحلیل واریانس» و اثر مدلی، میتزل و دوی<sup>۳</sup> (۱۹۵۶) در مورد مشاهدات کلاسی، رویه‌های چندگانه به‌طور همزمان بررسی گردید. - کورنفلد و توکی<sup>۴</sup> (۱۹۵۶) قواعدشان را در خصوص معادلات میانگین مجدورات مورد انتظار (EMS)<sup>۵</sup> برحسب مؤلفه‌های واریانس منتشر کردند.</p>	
<p>- دستاورد کرونیباخ و همکارانش، یک چارچوب مفهومی و استفاده از یک متدولوژی (تحلیل مؤلفه‌های واریانس) بود که مقالات بسیاری از محققان را، با وجود اینکه برخی از آن‌ها مغایر هم به نظر می‌رسیدند، یکپارچه کرده بود. - خطوط اصلی GT تک‌متغیری با گزارش‌های فنی کرونیباخ، راجاراتنام و گلیسر<sup>۶</sup> (در قالب سه مقاله هرکدام با نویسنده‌ی اول متفاوت) تا حد زیادی کامل شد.</p>	<p>۱۹۶۱- ۱۹۶۰</p>
<p>- بررسی‌های ناندا<sup>۷</sup> درباره‌ی اعتبار درون مجموعه‌ای، انگیزشی برای توسعه‌ی GT چندمتغیری<sup>۸</sup> گردید. - تیم کرونیباخ، توسعه GT چند متغیری را آغاز کردند که در آن هر یک از سطوح یک‌رویه ثابت<sup>۹</sup> یا بیشتر، با یک نمره جهان<sup>۱۰</sup> متمایز مرتبط می‌شد.</p>	<p>۱۹۶۴- ۱۹۶۶</p>
<p>کرونیباخ، گلیسر، ناندا و راجاراتنام، کتابی با عنوان «اتکاپذیری اندازه‌های رفتاری» منتشر کردند که هنوز کامل‌ترین و گسترده‌ترین بررسی از GT است.</p>	<p>۱۹۷۲</p>
<p>کاردینت، تورنر و ال<sup>۱۱</sup>، اصل تقارن<sup>۱۲</sup> را مطرح کردند.</p>	<p>۱۹۷۶</p>
<p>برنان و کانه<sup>۱۳</sup> در مقاله‌های مشترکشان، در پرداختن به موضوعات پیرامون اعتبار نمرات ملاک مرجع، GT را بکار گرفتند.</p>	<p>۱۹۷۷</p>
<p>شولسون و وب، مروری از GT در طی سال‌های ۱۹۸۰-۱۹۷۳ را منتشر کردند.</p>	<p>۱۹۸۱</p>

1. randomly parallel
2. relative & absolute error
3. Medley, Mitzel & Doi
4. Cornfield & Tukey
5. expected mean squares
6. Rajaratnam & Gleser
7. Nanda
8. multivariate
9. fixed facet
10. universe score
11. Cardinet, Tourneur & Allal
12. symmetry
13. Kane

مقاله کانه با عنوان «یک مدل نمونه‌گیری برای روایی»، یکی از مقاله‌های نظری مهم در ادبیات GT در ۲۵ سال گذشته است.	۱۹۸۲
- نوشتن یک مونوگراف درباره‌ی GT با عنوان «کلیات GT» توسط برنان که نسبت به اثر کروناخ و همکارانش (۱۹۷۲) از جامعیت کمتری برخوردار بود.	۱۹۸۳
- برنامه کامپیوتری بنام GENOVA توسط کریک <sup>۱</sup> و برنان جهت اجرای تحلیل GT طراحی شد.	
- علاقه به انجام سنجش در اواخر دهه‌ی ۱۹۸۰ شهرت نسبتاً زیادی را برای GT رقم زد.	
- کاربرد GT در سنجش‌های عملکردی به‌خصوص در مجموعه‌ای از سخنرانی‌ها و مقالات شولسون و همکارانش به‌خوبی نشان داده شد. همچنین، برنان و جانسون <sup>۲</sup> (۱۹۹۵) و برنان (۱۹۹۶) برخی از مسائل نظری و کاربردی را در این خصوص مطرح کردند.	۱۹۹۶-
- گزارش‌های خلاصه و ساده از GT از قبیل؛ اثر فلت <sup>۳</sup> و برنان (۱۹۸۹)، مقاله‌ی شولسون و همکاران (۱۹۸۹)، کتاب درسی مقدماتی شولسون و وب (۱۹۹۱) و مقاله‌ی برنان (۱۹۹۲) با عنوان «GT» منتشر شد.	۱۹۸۹

در سال ۲۰۰۱، کتاب جامعی در مورد GT توسط برنان نگاشته و منتشر شد. همچنین، می‌توان از کتابی که توسط کاردینت، جانسون و پینی<sup>۴</sup> (۲۰۱۰) در این زمینه نگاشته شده است، نام برد.

**نظریه تعمیم‌پذیری در مقابل نظریه کلاسیک آزمون.** طبق گفته‌ی شولسون و همکارانش (۱۹۸۹)، ارتباط بین GT و CTT همانند ارتباط بین آنوای عاملی و ساده<sup>۵</sup> است. محققى که آنوای ساده را بکار می‌گیرد، واریانس را به دو مؤلفه‌ی بین گروهی (واریانس منظم) و درون گروهی (واریانس تصادفی) تقسیم می‌کند. به شیوه‌ای مشابه، CTT نیز، واریانس را به دو مؤلفه، واریانس واقعی (واریانس منظم) و خطا (واریانس تصادفی) تقسیم می‌کند. با به‌کارگیری آنوای عاملی بجای ساده، محقق می‌پذیرد که عوامل چندگانه‌ای در واریانس کل داده‌ها سهیم است. آنوای ساده، سؤالات محدودتری را مطرح می‌کند و نسبت به آنوای عاملی با کارایی کمتری به آنها می‌پردازد. به همین ترتیب، بسط CTT به GT اثرات

1. Crick
2. Johnson
3. Feldt
4. Pini
5. factorial & simple ANOVA

چندگانه در واریانس اندازه گیری را تأیید می کند. GT همانند آنوای عاملی، واریانس را به چندین منبع نظیر؛ واریانس منظم در بین اهداف اندازه گیری<sup>۱</sup>، منابع چندگانه‌ی خطا و به تعاملاتشان تقسیم می کند.

مفروضات زیربنایی GT اساساً مشابه مفروضات CTT است. در هر دو نظریه، خطاها مستقل از نمره‌ی جهان / واقعی و نا همبسته‌اند. همان‌طور که در CTT، SEM در همه‌ی سطوح نمرات (برای همه آزمودنی‌ها) یکسان است در GT نیز، SEM یکسانی برای همه‌ی اهداف اندازه گیری بدون در نظر گرفتن نمره‌ی جهان زیربنایی به کار می رود (استراپ<sup>۲</sup>، ۲۰۰۲، به نقل از الخاروسی<sup>۳</sup>، ۲۰۱۲). از دیگر شباهت‌های دو نظریه‌ی مذکور می توان به موارد زیر اشاره کرد.

از CTT و GT می توان به عنوان نظریه‌های ارزش مورد انتظار<sup>۴</sup> نام برد که نمره واقعی<sup>۵</sup> یا نمره‌ی جهان را به عنوان یک ارزش مورد انتظار از نمرات مشاهده شده تعریف می کنند. نمره مورد نظر در CTT و GT، نمره‌ی مشاهده شده آزمون است و هدف اصلی این دو نظریه، ارزیابی کیفیت نمره مشاهده شده آزمون است که از طریق ضرایب اعتبار و خطاهای استاندارد برآورد می شود، بدون اینکه تلاشی جهت برآورد نمرات در صفت مکنون صورت گیرد. هر دو نظریه، به وضوح خطاهای اندازه گیری تصادفی را شامل می شوند و مفاهیم اعتبار (تعمیم پذیری) در هر دو، به خوبی تعریف شده است. واریانس نمره‌ی جهان در GT مشابه واریانس نمره‌ی واقعی در CTT است. همچنین، واریانس خطای نسبی<sup>۶</sup> ( $\sigma_8^2$ ) و ضریب تعمیم پذیری<sup>۷</sup> ( $E\hat{\rho}^2$ ) در GT به ترتیب مشابه واریانس خطا و ضریب اعتبار در CTT است. از دیگر شباهت‌ها، می توان به خطی بودن مدل‌های هر دو نظریه اشاره کرد و اینکه، در برابر نقض مفروضات مدل‌هایشان مقاوم هستند. به بیانی دیگر، مفروضات نسبتاً ضعیفی دارند.

1. objects of measurement
2. Strube
3. Alkharusi
4. expected value theories
5. true score
6. generalizability coefficient
7. generalizability coefficient

علاوه بر این مستلزم نمونه‌هایی با حجم بزرگ نیستند و واحد تحلیل در آن‌ها، نمرات آزمون است (شولسون و وب، ۱۹۸۱؛ بولس، هِنوفتس و بیلی<sup>۱</sup>، ۱۹۸۲؛ سوئن و لی<sup>۲</sup>، ۲۰۰۷؛ برنان، ۲۰۱۱a). دو نظریه مذکور، به‌رغم شباهت‌ها، تفاوت‌های بسیار مهمی نیز دارند که در زیر به شرح مواردی از آن پرداخته می‌شود.

چارچوب مفهومی: GT نسبت به CTT چارچوب مفهومی قدرتمندتری دارد که منجر به برطرف کردن تعدادی از تناقضات آشکار در چند بحث CTT از اعتبار شده است. دو ویژگی مهم GT که به حل تناقضات کمک می‌کند، عبارت‌اند از: تمایز گذاشتن GT میان رویه‌های<sup>۳</sup> اندازه‌گیری ثابت و تصادفی و همچنین قابلیت این نظریه در پرداختن به طرح‌های مختلف مطالعه‌ی تصمیم<sup>۴</sup> (مطالعه D) (برنان، ۲۰۱۱a).

پیچیدگی مفهومی: چارچوب مفهومی قدرتمند GT و قابلیت آن در جداسازی منابع خطا منجر به پیچیدگی مفهومی آن شده است که این موضوع، شناخت محققان را می‌طلبد (برنان، ۲۰۱۱a). همان‌طور که شولسون و همکاران (۱۹۸۹) مطرح کرده‌اند، پیچیدگی GT می‌تواند روشن کند که چرا CTT روش ارجحی برای برآورد اعتبار باقی مانده است. تعاریف فرم‌های موازی: در CTT تعاریف آزمون‌های موازی و آزمون‌های اساساً تائو معادل، اغلب غیرقابل دفاع هستند. درحالی‌که GT فرض می‌کند که آزمون‌ها تصادفی موازی هستند و محتوای آزمون یک نمونه تصادفی از حیطه یا جهان تعریف شده در نظر گرفته می‌شود (سوئن و لی، ۲۰۰۷). برنان (۲۰۱۱a) بیان می‌کند که GT مفهوم ساده‌ای از موازی بودن را در خود گنجانده است که کاملاً متفاوت با مفهوم کلاسیک فرم‌های موازی در CTT است. همچنین مطرح می‌کند، هر دو نوع موازی بودن ایده‌آل هستند و هیچ‌گاه احتمال اینکه کاملاً واقعیت داشته باشد، نیست. اگرچه یکی یا دیگری ممکن است در زمینه‌ای خاص مناسب‌تر باشد.

1. Bolus, Hinofotis & Bailey
2. Suen & Lei
3. facets
4. decision study



مدل سازی نمرات مشاهده شده: در CTT نمره‌ی مشاهده شده ( $X$ ) یک فرد در آزمون از دو متغیر غیرقابل مشاهده یعنی؛ یک نمره واقعی ( $T$ ) و یک نمره خطا ( $E$ ) تشکیل شده است و به صورت مدل خطی ( $X = T + E$ ) نشان داده می‌شود. در GT نمره مشاهده شده در یک آزمون از رویه‌های مختلف مورداستفاده در آزمون تأثیر می‌پذیرد و با توجه به رویه‌های مورداستفاده در آزمون معرف عملکرد فرد در همان رویه‌هاست (کیامنش، ۱۳۷۴). طبق بیان شولسون و وب (۱۹۹۱)؛ نمره‌ی مشاهده شده‌ی یک فرد به یک نمره جهان و یک یا بیشتر از یک منبع خطا تجزیه می‌شود. در صفحه‌ی ۱۷، مدل سازی نمره مشاهده شده در GT برای طرح  $2 \times 1 \times 1$  نشان داده شده است.

منابع چندگانه‌ی خطا و برآوردهای چندگانه از اعتبار: در وضعیت‌های اندازه‌گیری پیچیده که شامل منابع چندگانه‌ای از خطای اندازه‌گیری (رویه‌ها) است، CTT قادر به برآورد اعتبار نیست؛ زیرا شیوه‌های سنتی اعتبار تنها برای یک منبع خطا طراحی شده‌اند. روش معمول CTT برای برآورد اعتبار این است که از روش‌های مختلفی (همچون باز آزمایی، بین ارزیابان، همسانی درونی،...) استفاده می‌کند. در نتیجه، برآوردهای خطا و برآوردهای اعتبار (نسبت واریانس نمره واقعی به مشاهده شده) مطابق با طرح جمع‌آوری داده‌ها تغییر می‌کند. اینکه روش‌های مختلف، ضرایب اعتبار مختلفی را به دنبال دارند، به نوبه‌ی خود منجر به خطاهای استاندارد اندازه‌گیری متفاوتی می‌شود. سؤالی که اینجا پیش می‌آید این است که در چنین وضعیتی، دقیق‌ترین برآورد ضریب اعتبار کدام است؟ و به منظور ساخت فاصله‌های اطمینان حول نمرات مشاهده، کدام خطای استاندارد اندازه‌گیری را باید به کار برد؟ متأسفانه CTT قادر به پاسخ‌گویی به این سؤالات نیست و تعاریف متغیر نمرات واقعی و خطا را مسکوت می‌گذارد. در حالی که GT اذعان می‌دارد که تعاریف چندگانه‌ای از نمرات جهان و خطا می‌تواند وجود داشته باشد. برخلاف CTT، در GT می‌توان منابع چندگانه خطا را همزمان در ترکیب‌های متفاوتی از تصادفی یا ثابت در نظر گرفت. با تعیین اینکه آیا یک رویه تصادفی یا ثابت باشد، امکان برآورد اعتبار و خطای استاندارد ناشی از منابع معین خطا در GT وجود دارد. به بیانی دیگر، GT سهم هر منبع خطا را در واریانس نمرات آزمون تعیین می‌کند و فرصت محاسبه‌ی برآوردهای متفاوتی از اعتبار

را می‌دهد که بستگی به این دارد که کدام منبع خطا برای هر استفاده‌ی خاص از آزمون مهم در نظر گرفته می‌شود (شولسون و همکاران، ۱۹۸۹؛ ون لیون، ۱۹۹۷؛ سوئن و لی، ۲۰۰۷؛ فن و سان، ۲۰۱۳).

سنجش‌های هنجار مرجع و ملاک مرجع<sup>۱</sup> (تمایز گذاشتن میان تصمیم‌های نسبی و مطلق<sup>۲</sup>): CTT چون نمی‌تواند خطای اندازه‌گیری منظم را در خود جای دهد، تنها برای سنجش اعتبار آزمون‌های هنجار مرجع مناسب است. درحالی‌که GT به دلیل انعطاف‌پذیری که دارد هر دو خطای اندازه‌گیری نسبی و مطلق را در خود جای می‌دهد و میان تصمیم‌های نسبی و مطلق تمایز می‌گذارد؛ بنابراین، هم برای سنجش هنجار مرجع مناسب است و هم نسبت به CTT، برآورد مناسبی از اعتبار را برای آزمون‌های ملاک مرجع فراهم می‌آورد (شولسون و وب، ۱۹۹۱؛ سوئن و لی، ۲۰۰۷؛ کمازاوا<sup>۳</sup>، ۲۰۰۹).

طراحی و بهینه‌سازی پروتکل اندازه‌گیری: در طرح‌های اندازه‌گیری چند رویه‌ای، با به‌کارگیری GT می‌توان هر یک از منابع واریانس تشکیل‌دهنده‌ی خطای کل را شناسایی و با انجام یک مطالعه‌ی تعمیم‌پذیری<sup>۴</sup> (مطالعه‌ی G)، سهم جداگانه‌ی هر یک از آن‌ها را برآورد کرد. به دنبال آن، برای بهبود روند اندازه‌گیری و از طریق انجام مطالعه‌ی D، با دست‌کاری رویه یا رویه‌هایی که سهم بیشتری در خطای کل دارند، می‌توان خطای اندازه‌گیری برآورد شده را کاهش داد. قابلیت GT در پرداختن به طرح‌های مختلف مطالعه D به محقق اجازه می‌دهد که با در نظر گرفتن اندازه اعتبار دلخواه، هزینه، زمان و دیگر محدودیت‌های عملی، پروتکل اندازه‌گیری بهینه‌ای را برای مطالعات واقعی و کاربردهای عملی طراحی و ارزشیابی کند. در مقابل، در طرح‌های چند رویه‌ای، نظریه کلاسیک از طریق روش شناخته‌شده‌ای مثل فرمول پیشگویی اسپیرمن- براون<sup>۵</sup> یک برآورد تک‌بعدی از آنچه خطا را کاهش خواهد داد فراهم می‌کند. به‌بیان‌دیگر، این فرمول برخلاف مطالعات

1. norm-referenced & criterion-referenced
2. relative & absolute decisions
3. Kumazawa
4. generalizability study
5. Spearman-Brown prophecy formula

D، نمی تواند اثر تغییر تعداد سطوح دورویه یا بیشتر را به طور همزمان ارزیابی کند (شولسون و وب، ۱۹۹۱؛ بولس و همکاران، ۱۹۸۲؛ فن و سان، ۲۰۱۳).

نظریه تعمیم پذیری چندمتغیری: GT چندمتغیری، مسائل اعتبار را در راستای جهان های تعمیم چندگانه گسترش داده است که وضعیت متناظر آن در CTT وجود ندارد (برنان، ۲۰۱۱a).

اصل تقارن: در CTT، اشخاص به عنوان هدف اندازه گیری در نظر گرفته می شوند و راجع به مناسب بودن یک اندازه برحسب اینکه چقدر خوب ما را برای تمایز قائل شدن میان افراد قادر می سازد، قضاوت می شود. همه رویه های دیگر (سؤالات آزمون، موقعیت ها، ارزیابان، غیره) به عنوان خطا در نظر گرفته می شوند (شولسون و همکاران، ۱۹۸۹). کاردینت، تورنر و ال (۱۹۷۶) با مطرح کردن اصل تقارن در GT، بیان کردند؛ برخلاف تمرکز سنتی روی افراد، هدف اندازه گیری ممکن است بسته به هدف خاص تصمیم گیرنده تغییر کند و تفاوت های فردی ممکن است به عنوان منبع خطا در نظر گرفته شوند.

**مفاهیم و اصطلاحات در GT.** بنا به گفته ی برنان (۲۰۱۰)، مهم ترین وجه و ویژگی منحصر به فرد GT، چارچوب مفهومی آن است. لذا، برای ارتباط برقرار کردن با این نظریه و درک چارچوب مفهومی آن، مفاهیم و اصطلاحات رایج در آن جداگانه شرح داده می شود.

هدف اندازه گیری: هدف اندازه گیری عاملی است که محقق عمدتاً روی آن تمرکز می کند و تغییر پذیری میان آن ها مطلوب است. در بسیاری از وضعیت های اندازه گیری، افراد هدف اندازه گیری هستند. باین حال، طبق اصل تقارن، هر یک از رویه های موجود در یک طرح را می توان به عنوان هدف اندازه گیری در نظر گرفت. کاردینت و همکاران (۲۰۱۰) مطرح کردند که واریانس حاصل از هدف اندازه گیری مترادف با مفهوم واریانس نمره ی واقعی در CTT است.

رویه‌ها و موقعیت‌ها<sup>۱</sup>: منابع بالقوه‌ی خطا در تعمیم دهی را رویه‌ها و سطوح<sup>۲</sup> رویه‌ها را، موقعیت‌ها یا حالت‌ها می‌نامند. اصطلاحات رویه‌ها و موقعیت‌ها مشابه با عامل‌ها<sup>۳</sup> و سطوح در ادبیات طرح‌های آزمایشی است (وب، رولی و شولسون، ۱۹۸۸؛ شولسون و وب، ۱۹۹۱). می‌توان گفت، یک رویه مجموعه‌ای از سطوح<sup>۴</sup> مشابه‌اندازه‌گیری است (برنان، ۲۰۰۱) که محقق تصمیم دارد عملکرد آزمودنی‌ها را در این سطوح موردسنجش قرار دهد و سطوح یک اصطلاح کلی است که به فرم‌ها یا محرک‌های آزمونی خاص، مشاهده‌کننده‌ها، وضعیت‌ها یا شرایط مشاهده و غیره اشاره می‌کند (کرونباخ، راجاراتنام و گلیسر، ۱۹۶۳). جهان و جامعه<sup>۵</sup>: در GT؛ واژه‌ی جهان به سطوح اندازه‌گیری اختصاص داده شده است، درحالی‌که واژه‌ی جامعه برای اهداف اندازه‌گیری استفاده می‌شود (برنان، ۲۰۰۱).

جهان مشاهدات قابل قبول<sup>۶</sup>: جهان مشاهدات قابل قبول به وسیله همه ترکیبات ممکن از سطوح رویه‌ها تعریف می‌شود. یک اندازه رفتاری (مثل نمره آزمون) به عنوان نمونه‌ای از یک جهان مشاهدات قابل قبول تصور می‌شود که این جهان، شامل همه مشاهدات ممکن است که تصمیم‌گیرندگان جایگزین‌های قابل قبولی از آن برای مشاهده موجود (تحت مطالعه) در نظر می‌گیرند (وب و شولسون، ۲۰۰۵). جهان مشاهدات قابل قبول می‌تواند بزرگ‌تر از جهانی باشد که تصمیم‌گیرنده می‌خواهد به آن تعمیم دهد (شولسون و وب، ۱۹۸۱). طبق گفته‌ی کرونباخ و همکاران (۱۹۶۳)، جهان شامل تعداد محدود یا نامحدودی<sup>۷</sup> از سطوح است که لازم است آن را به وضوح تعریف کرد، به گونه‌ای که سطوحی که در جهان قرار می‌گیرند مشخص باشد؛ اما ضروری نیست که جهان همگن باشد.

جهان تعمیم: عمدتاً محقق علاقه‌مند به تعمیم از یک مشاهده به دیگر مشاهداتی است که عضو جهان همانندی هستند. از آنجاکه یک اندازه معین می‌تواند به طور قابل قبولی به

1. conditions
2. levels
3. factors

۴. در این مقاله، از واژه‌ی سطوح برای conditions استفاده کرده‌ایم.

5. universe & population
6. universe of admissible observations
7. universe of generalization

جهان‌های متفاوت بسیاری تعمیم داده شود، لذا قبل از اینکه محقق مطالعه تعمیم‌پذیری را انجام دهد، باید جهانی که علاقه‌مند به تعمیم دهی به آن است را تعیین کند (کروباخ و همکاران، ۱۹۶۳). جهان تعمیم، یعنی؛ جهانی که محقق قصد دارد نتایج یک روش اندازه‌گیری خاص را به آن تعمیم دهد (برنان، ۲۰۰۱) و می‌تواند با جهان مشاهدات قابل قبول یکسان و یا زیرمجموعه‌ای از رویه‌ها و سطوحشان در جهان مشاهدات قابل قبول باشد. جهان تعمیم را می‌توان از طریق؛ کاهش دادن جهان مشاهدات قابل قبول، یعنی کاهش سطوح یک‌رویه (برای مثال، تبدیل به یک‌رویه ثابت)، انتخاب کردن و در نتیجه کنترل کردن سطحی از رویه‌ای، یا با نادیده گرفتن رویه‌ای تعیین کرد (شولسون و وب، ۱۹۸۱).

نمره جهان: مارکولیدس<sup>۱</sup> (۱۹۹۶) بیان می‌کند، مفهوم نمره‌ی جهان را می‌توان قلب GT در نظر گرفت. متأسفانه به سبب آنکه نمره‌ی جهان را فقط می‌توان برآورد کرد، انتخاب سؤال، فرم آزمون، یا موقعیت سنجش خاصی به‌طور اجتناب‌ناپذیری خطا را وارد روند اندازه‌گیری می‌کند. طبق بیان برنان (۲۰۱۰)، در اصل برای هر فرد (فرد هدف اندازه‌گیری در نظر گرفته شده است)، برای هر تکرار از روش اندازه‌گیری در جهان تعمیم می‌توان نمره‌ی میانگینی تصور کرد. برای چنین فردی، ارزش مورد انتظار از این نمرات میانگین به‌عنوان نمره جهان فرد تعریف می‌شود و واریانس نمرات جهان همه افراد در جامعه، واریانس نمره جهان نامیده می‌شود. در هر اندازه‌گیری سعی می‌شود شرایطی مهیا شود که واریانس نمره جهان زیاد و اندازه‌ی سایر مؤلفه‌های واریانس (واریانس‌های خطا) کم باشد. یک شخص می‌تواند نمرات جهان متفاوتی داشته باشد که بستگی به مشخصات جهان دارد. به بیان دقیق‌تر، نمره‌ی جهان تعریف نمی‌شود تا اینکه هدف اندازه‌گیری در مطالعه‌ی D مشخص شود (وب، شولسون و هرتل<sup>۲</sup>، ۲۰۰۷).

رویه تصادفی<sup>۳</sup>: GT اساساً نظریه‌ی اثرات تصادفی است. یک‌رویه تصادفی از طریق نمونه‌گیری سطوح یک‌رویه به‌طور تصادفی ایجاد می‌شود. حتی اگر سطوح یک‌رویه به‌طور

---

1. Marcoulides  
2. Haertel  
3. random facet

تصادفی نمونه‌گیری نشده باشند، رویه ممکن است تصادفی در نظر گرفته شود اگر سطوح مشاهده نشده در مطالعه  $G$  قابل‌جابه‌جایی با سطوح مشاهده شده باشد (وب و همکاران، ۲۰۰۷).

رویه ثابت: سطوح یک‌رویه‌ی ثابت در مطالعه  $G$  تمام سطوح ممکن موردنظر است. رویه ثابت زمانی رخ می‌دهد که:

- تصمیم‌گیرنده با قصد قبلی سطوح خاصی را انتخاب کند و تمایلی به تعمیم دهی فراتر از آن‌ها را نداشته باشد.
- تعمیم دادن به ورای سطوح مشاهده‌شده غیرمنطقی باشد.
- هنگامی که کل جهان سطوح، کوچک باشد و طرح اندازه‌گیری همه سطوح را در برگیرد (شولسون و همکاران، ۱۹۸۹؛ وب و همکاران، ۲۰۰۷).

مدل‌های تصادفی و آمیخته<sup>۱</sup> با جهان‌های تعمیم نامحدود و محدود: بنا به گفته‌ی شولسون و وب (۱۹۸۱، ۱۹۹۱)، برای معنی‌داری تحلیل‌ها، در هر تحلیل تعمیم‌پذیری حداقل باید یک‌رویه تصادفی باشد. از این رو، در  $GT$  مدل ثابت (مدلی که تمام رویه‌های آن ثابت باشد) وجود ندارد. در مدل تصادفی، همه‌ی رویه‌ها تصادفی هستند و این مدل‌ها با جهان‌های تعمیم نامحدود مرتبط‌اند. در مقابل در مدل آمیخته، ترکیبی از رویه‌های ثابت و تصادفی وجود دارد. تثبیت یک‌رویه، جهان تعمیم محدود را به دنبال خواهد داشت؛ بنابراین جهان تعمیم برای رویه ثابت محدودتر از جهان تعمیم برای رویه تصادفی است. جهان تعمیم محدود از آنجا که نسبت به جهان نامحدود کمتر مستعد خطاست، لذا واریانس نمره جهان و به تبع، ضریب تعمیم‌پذیری بزرگ‌تری دارد. همان‌طور که برنان (۲۰۱۰) متذکر می‌شود، نمی‌توان گفت که جهان محدود رجحان دارد، زیرا محدود کردن یک جهان همچنین گسترده‌ای که یک محقق می‌تواند به آن تعمیم دهد را محدود می‌کند. در اکثر زمینه‌های ارزیابی آموزشی و روان‌شناختی، رویه‌ها را باید به‌عنوان تصادفی در نظر گرفت تا بتوان نتایج

---

1. random & mixed models

را به بیش از یک وضعیت اندازه گیری خاص تعمیم داد (بریچ، سوامیناتان، ولش و چفولز<sup>۱</sup>، ۲۰۱۴).

رویه‌های متقاطع، آشیانه‌ای و درآمیخته<sup>۲</sup>: دو شرط زیر را در نظر بگیرید؛ (الف) دو سطح یا بیشتر (سطوح چندگانه) از A با هر یک از سطوح B مشاهده شود و (ب) سطوح متفاوتی از A با هر یک از سطوح B مربوط باشد. اگر فقط شرط اول برقرار باشد اما سطوح یکسانی از A با هر یک از سطوح B مربوط باشد، گفته می‌شود؛ دورویه متقاطع هستند. اگر هر دو شرط برقرار باشد، گفته می‌شود؛ A درون B آشیانه کرده است. در صورتی که فقط شرط دوم برقرار باشد، یعنی؛ یک سطح متفاوت از A برای هر یک از سطوح B رخ دهد، گفته می‌شود؛ A با B درآمیخته است (شولسون و وب، ۱۹۹۱). به عنوان مثال، هر یک از ارزیابان (r) مجموعه‌ی یکسانی از سؤالات (i) را نمره گذاری می‌کنند. در این صورت، ارتباط بین r و i متقاطع بوده و به صورت  $r \times i$  نشان داده می‌شود. علامت «متقاطع با» خوانده می‌شود. اگر هر یک از ارزیابان (r) مجموعه‌ی متفاوتی از سؤالات (i) را نمره گذاری کنند، در این حالت گفته می‌شود؛ سؤالات در درون ارزیابان آشیانه کرده‌اند و به صورت r: i نمایش داده می‌شود. علامت: «آشیانه کرده درون» خوانده می‌شود. اگر هر یک از ارزیابان (r)، یک سؤال (i) متفاوت را نمره گذاری کنند، گفته می‌شود؛ رویه سؤالات با رویه ارزیابان درآمیخته است.

طرح‌های تعمیم‌پذیری<sup>۳</sup>: عناصر یک طرح (طرح‌های مطالعه G یا D)، عبارت‌اند از: تعداد رویه‌ها، روش نمونه‌گیری که برای انتخاب سطوح رویه‌ها بکار گرفته می‌شود (ثابت در مقابل تصادفی) و ارتباط میان رویه‌ها (متقاطع در مقابل آشیانه‌ای). ترکیب این عناصر می‌تواند طرح‌های مختلف بسیاری را به وجود بیاورد. برای مثال، دورویه‌ای کاملاً متقاطع با رویه‌های تصادفی، سه رویه‌ای آشیانه‌ای با رویه‌های ثابت (براکتر، یودر و مک ویلیام<sup>۴</sup>، ۲۰۰۶). طرح‌ها را می‌توان بر اساس تعداد رویه‌هایشان نام گذاری کرد. اگر فقط یک رویه در

1. Briesch, Swaminathan, Welsh & Chafouleas
2. confounded
3. generalizability designs
4. Bruckner, Yoder & McWilliam

آزمون به کار گرفته شود، طرح آزمون را می‌توان طرح یک‌رویه‌ای<sup>۱</sup> نامید. بنا به گفته‌ی برنان (۲۰۰۳)، برای طرح‌های یک‌رویه‌ای، جهان مشاهدات قابل قبول و جهان تعمیم شامل سطوحی از همان یک‌رویه می‌شود. اگر دورویه به صورت همزمان در آزمون به کار گرفته شود طرح آزمون را می‌توان طرح دورویه‌ای<sup>۲</sup> دانست. همچنان که تعداد رویه‌ها (منابع بالقوه خطای اندازه‌گیری) افزایش می‌یابد، تعداد طرح‌های انتخابی در دسترس برای استفاده از نمره آزمون نیز افزایش می‌یابد. همچنین، طرح‌ها را می‌توان بر اساس روابط میان رویه‌ها نام‌گذاری کرد. طبق بیان شولسون و وب (۱۹۹۱)؛ نظریه تعمیم‌پذیری، طرح‌های کاملاً متقاطع یا متقاطع (همه‌ی رویه‌ها متقاطع هستند)، کاملاً آشیانه‌ای (همه‌ی رویه‌ها آشیانه‌ای هستند)، نسبتاً آشیانه‌ای (هر دورویه‌ی متقاطع و آشیانه‌ای وجود دارد) و طرح‌های درآمیخته را در خود جای داده است. در تقسیم‌بندی دیگر، طرح‌های تعمیم‌پذیری را می‌توان در دو گروه متعادل و نامتعادل<sup>۳</sup> و نیز، تک متغیری و چندمتغیری جا داد.

انواع مطالعات در GT:GT دو نوع مطالعه را شامل می‌شود؛ مطالعه‌ی تعمیم‌پذیری (G) و مطالعه‌ی تصمیم (D). یک مطالعه G به منظور برآورد مؤلفه‌های واریانس زیربنایی یک فرآیند اندازه‌گیری از طریق تعریف کردن جهان مشاهدات قابل قبول طراحی می‌شود. در مطالعه D، تصمیم‌گیرنده از اطلاعات به دست آمده‌ی مطالعه G استفاده می‌کند تا یک اندازه‌گیری که خطا برای هدف خاصی در آن به حداقل برسد را طراحی کند. در طراحی یک مطالعه D، تصمیم‌گیرنده باید: (الف) جهان تعمیم را تعریف کند؛ یعنی، تعداد و گستره‌ی رویه‌هایی را که مایل است به آن تعمیم دهد. (ب) نوع تفسیر در نظر گرفته شده برای اندازه‌گیری را بیان کند؛ زیرا از نمره آزمون یکسانی به شیوه‌های متفاوتی می‌توان استفاده کرد. برای مثال، برخی از تفسیرها ممکن است روی تفاوت‌های فردی تمرکز کنند (تصمیمات نسبی)، برخی ممکن است نمره مشاهده شده را به عنوان برآوردی از نمره جهان شخص بکار گیرند (تصمیمات مطلق). درحالی که سایرین، نمره مشاهده شده را در یک

- 
1. one facet design
  2. two facet design
  3. balanced & unbalanced



برآورد رگرسیون از نمره جهان بکار ببرند. هریک از این تفاسیر مطرح شده با خطای اندازه گیری متفاوتی مرتبط هستند. (ج) از اطلاعات مطالعه G که در مورد اندازه‌ی منابع مختلف خطای اندازه‌گیری به دست آمده است، جهت بهینه‌سازی طرح اندازه‌گیری (کمینه‌سازی خطا و بیشینه‌سازی اعتبار) مطالعه D استفاده کند. این ارزشیابی تا حدی مشابه با فرمول پیشگویی اسپیرمن- بروان در CTT انجام می‌گیرد. به طوری که با افزایش تعداد سطوح یک‌رویه در یک اندازه‌گیری می‌توان سهم خطای متعلق به آن رویه را کاهش داد. البته در مطالعه D، علاوه بر تغییر در تعداد سطوح رویه‌ها، دست‌کاری رویه‌ها شامل تغییر در ماهیت (تصادفی / ثابت) و روابط آن‌ها (مقاطع / آشیانه‌ای) نیز می‌شود (شولسون و وب، ۱۹۸۱؛ شولسون و همکاران، ۱۹۸۹؛ شولسون و وب، ۱۹۹۱).

ملاحظات مطالعات G و D: تمایز قائل شدن میان مطالعات G و D مهم است، زیرا داده‌های تصمیم‌اند که اتکاپذیری‌شان مدنظر است. این دو نوع مطالعه می‌توانند متفاوت از هم باشند و طرح‌های تجربی متفاوتی داشته باشند (کروناخ و همکاران، ۱۹۶۳). در مطالعات G می‌بایست تا جایی که امکان‌پذیر است از طرح‌های مقاطع استفاده کرد؛ زیرا با این طرح‌ها، امکان برآورد جداگانه‌ی همه‌ی مؤلفه‌های واریانس وجود دارد. درحالی که در دیگر طرح‌ها، اثر مستقیم مؤلفه واریانس مربوط به رویه‌ی آشیانه‌ای به صورت جداگانه برآورد نمی‌شود. در مواردی، رویه‌ها ذاتاً آشیانه‌ای هستند و محقق حق انتخاب ندارد؛ مانند؛ سؤالاتی که درون هر یک از فصل‌های یک کتاب درسی آشیانه کردند. با این حال، اغلب در مطالعه D برای سهولت، افزایش اندازه نمونه، یا هر دو، طرح‌های آشیانه‌ای به کار گرفته می‌شود. اگر رویه‌های ثابت به کار روند و سطوح این رویه‌ها قابل تعویض نباشد، باید برای هر سطح، مطالعات G جداگانه‌ای انجام داد. حجم نمونه در مطالعه D ضرورت ندارد همانی باشد که در مطالعه G به کار رفته است. در این مطالعه تصمیم‌ها معمولاً مبتنی بر میانگین همه‌ی مشاهدات چندگانه (مانند سؤالات آزمون) است، درحالی که مطالعه G بر روی مشاهدات منفرد (مانند یک تک سؤالی) تمرکز می‌کند. انتخاب تعداد سطوح هر رویه در مطالعه D و همچنین انتخاب طرح (آشیانه‌ای در مقابل مقاطع، رویه ثابت در مقابل تصادفی)، ملاحظات منطقی و عملی و همچنین مسائل اتکاپذیری را در برمی‌گیرد (شولسون و وب،

۱۹۸۱؛ شولسون و همکاران، ۱۹۸۹؛ وب و شولسون، ۲۰۰۵). در پاسخ به سؤال «چه هنگام مطالعات G و مطالعات D به کار می‌روند؟» براون (۲۰۰۵) بیان می‌کند که ابتدا مطالعه G باید انجام شود؛ سپس و تنها بعد از آن، یک مطالعه D کاربردی را دنبال کرد. انجام هر یک بدون دیگری کمتر معقول است و این دو مطالعه را باید با هم و به‌طور متوالی به کار بست. برآورد مؤلفه‌های واریانس: یکی از اهداف GT، ارزیابی کردن منابع عمده‌ی تغییرپذیری است، به‌طوری‌که تغییرپذیری نامطلوب را بتوان در جمع‌آوری داده‌های آینده کاهش داد. GT، اندازه‌ی تغییرپذیری را برحسب مؤلفه‌های واریانس بیان می‌کند (شولسون و وب، ۱۹۹۱). برنان (۲۰۰۱)، هدف از مطالعه‌ی G را برآورد مؤلفه‌های واریانس برای جامعه و جهان مشاهدات قابل قبول عنوان می‌کند که این مؤلفه‌های واریانس برآورد شده را می‌توان در مطالعات مختلف D به‌منظور طراحی روش‌های اندازه‌گیری کارآمد برای استفاده‌ی عملیاتی و فراهم کردن اطلاعات برای گرفتن تصمیم‌های اساسی در مورد اهداف اندازه‌گیری به کار برد. هنگامی که فرض می‌شود جامعه و جهان مشاهدات قابل قبول به‌طور نامحدودی بزرگ هستند، مؤلفه‌های واریانس؛ مؤلفه‌های واریانس اثرات تصادفی<sup>۱</sup> نامیده می‌شوند. در هر مطالعه تعداد متفاوتی از مؤلفه‌های واریانس به کار گرفته می‌شود که بستگی به طرح مطالعه و نوع تفسیر اندازه‌ها دارد (فن و سان، ۲۰۱۳). طبق بیان مارکولیدس (۱۹۹۶)؛ تمرکز GT روی این مؤلفه‌های واریانس است. اگرچه تحلیل واریانس سنتی رایج‌ترین شیوه‌ی استفاده‌شده برای برآورد مؤلفه‌های واریانس مطالعه‌ی G است، با روش‌های دیگری نیز می‌توان مؤلفه‌های واریانس را برآورد کرد؛ از جمله، روش‌های بیزین، روش‌های واریانس حداقل<sup>۲</sup>، برآورد بیشینه درست‌نمایی محدود<sup>۳</sup> (RMLE) و رویکرد تحلیل ساختار کوواریانس<sup>۴</sup>. موضوع مهمی که باید به آن توجه داشت، منفی بودن مؤلفه (های) واریانس برآورد شده است. این موضوع گرچه از لحاظ نظری ناممکن است، اما در بعضی از وضعیت‌های اندازه‌گیری رخ می‌دهد. هنگامی که اندازه‌ی برآوردهای منفی بزرگ باشد،

1. random effects variance components
2. minimum variance methods
3. restricted maximum likelihood estimation
4. covariance structure analysis approach

ممکن است که آن‌ها، نامعین بودن مدل اندازه‌گیری را منعکس کنند. در این صورت مدل مطالعه‌ی  $G$  باید مجدداً تعیین شود و مؤلفه‌های واریانس دوباره برآورد شوند. هنگامی که اندازه‌ی برآوردهای منفی کوچک (نزدیک به صفر) باشد، احتمالاً به سبب خطای نمونه‌گیری است که در این حالت می‌توان از روش‌های پیشنهادی استفاده کرد. کروناخ و همکارانش، توصیه کردند که صفر جایگزین برآوردهای منفی شود. در صورتی که مؤلفه واریانس منفی در معادلات میانگین مجذورات مورد انتظار (EMS) مربوط به دیگر مؤلفه‌های واریانس ظاهر شد، مقدار صفر به جای برآورد منفی گذاشته شود و برآورد مؤلفه‌های واریانس جدید محاسبه شود. این روش برآوردهای همراه با سوگیری را پدید می‌آورد. برنان، توصیه کرد که صفر جایگزین برآوردهای منفی شود، اما هر جا که مؤلفه واریانس مذکور در معادلات میانگین مجذورات مورد انتظار دیگر مؤلفه‌ها ظاهر شد، برآورد منفی بکار گرفته شود. این فرآیند برآوردهای غیر سوگیری را فراهم می‌کند. با این حال، به لحاظ آماری تغییر برآوردها به صفر رضایت‌بخش نیست. از دیگر توصیه‌ها، استفاده از دیگر رویکردهای برآورد مؤلفه‌های واریانس از جمله، رویکرد بیزین و یا RMLE است (شولسون و وب، ۱۹۸۱؛ ۱۹۹۱).

انواع تصمیم و واریانس‌های خطا: دو نوع تصمیم وجود دارد که واریانس خطا و در نتیجه ضرایب، به‌طور متفاوتی برای هر یک از آن‌ها تعریف می‌شود. اولی، تصمیم نسبی است اگر تصمیم درباره افراد مبتنی بر جایگاهشان در ارتباط با دیگران باشد. این نوع تصمیم بر روی تفسیر هنجار مرجع نمره متمرکز است؛ یعنی نمرات اندازه‌گیری برای متمایز کردن آزمودنی‌ها به کار می‌روند. اعتبار اندازه‌گیری در این حالت مربوط به ثبات جایگاه نسبی افراد است نه در مورد ثبات نمرات واقعی. واریانس خطا برای تصمیم نسبی را با علامت  $\sigma^2_{\theta}$  نشان می‌دهند و آن را واریانس خطای نسبی می‌نامند. این نوع واریانس شامل همه‌ی مؤلفه‌های واریانس تعاملی است که هدف اندازه‌گیری را در برمی‌گیرد. واریانس خطای نسبی به صورت تفاوت میان نمره انحرافی مشاهده‌شده فرد<sup>۱</sup> و نمره انحرافی جهان<sup>۲</sup> او تعریف

- 
1. person's observed deviation score
  2. universe deviation score

می شود. دومی، تصمیم مطلق است اگر تصمیم درباره افراد مبنی بر نمراتشان در ارتباط با یک ملاک باشد. به بیانی دیگر، تصمیم مطلق بر روی سطح عملکرد افراد بدون توجه به رتبه‌ی آن‌ها متمرکز است و در ارتباط با ثابت جایگاه نسبی افراد و هم ثابت نمرات واقعی است. واریانس خطا برای تصمیم مطلق را با علامت  $\sigma_{\Delta}^2$  نشان می دهند و آن را واریانس خطای مطلق<sup>۱</sup> می نامند که شامل همه مؤلفه‌های واریانس مدل به جز هدف اندازه گیری است. این نوع واریانس، بیانگر تفاوت میان نمره مشاهده شده و نمره جهان فرد است. در کل، واریانس خطای نسبی کمتر از واریانس خطای مطلق است، زیرا شامل مؤلفه‌های واریانس کمتری است. این نشان می دهد که تفسیرهای نسبی در مورد نمرات افراد نسبت به تفسیرهای مطلق کمتر مستعد خطا هستند (وب و شولسون، ۲۰۰۵؛ برنان، ۲۰۱۰؛ فن و سان، ۲۰۱۳).

انواع ضرایب: ضریب اعتبار و محاسبه آن بستگی به مفهوم سازی خطای از پیش تعیین شده به عنوان مطلق یا نسبی دارد. به بیان دیگر، GT میان دو نوع ضریب اعتبار تمایز می گذارد: یکی، ضریب تعمیم پذیری است و زمانی به کار می رود که تصمیم ها نسبی هستند. آن را با علائم  $E\rho^2$  یا  $\rho^2$  نمایش می دهند و فرمول آن به صورت زیر است:

$$E\rho^2 = \frac{\sigma_p^2}{\sigma_p^2 + \sigma_{\delta}^2}$$

دیگری، شاخص اتکاپذیری ( $\varphi$ ) است و برای تصمیم های مطلق به کار می رود. این شاخص، بدین صورت فرمول بندی می شود:

$$\varphi = \frac{\sigma_p^2}{\sigma_p^2 + \sigma_{\Delta}^2}$$

تفاوت ضریب تعمیم پذیری و شاخص اتکاپذیری در این است که اولی شامل واریانس خطای نسبی و دومی شامل واریانس خطای مطلق است؛ بنابراین شاخص اتکاپذیری عموماً کمتر از ضریب تعمیم پذیری است (وب و شولسون، ۲۰۰۵؛ برنان، ۲۰۰۱).

1. absolute error variance

فرآیند طراحی، تحلیل و تفسیر یک مطالعه اندازه‌گیری. برای آشنایی با چگونگی انجام یک مطالعه اندازه‌گیری و رعایت مسائل کلیدی موردنیاز آن، در این بخش فرایند طراحی، تحلیل و تفسیر یک مطالعه اندازه‌گیری در قالب یک مثال با داده‌های فرضی گام‌به‌گام تشریح می‌شود. همچنین، برای افرادی که علاقه‌مند به درک نحوه‌ی محاسبات برآورد مؤلفه‌های واریانس، واریانس‌های خطا و همچنین، ضرایب تعمیم‌پذیری هستند؛ فرمول‌ها و چگونگی انجام محاسبات (مربوط به طرح مفروض) ارائه می‌گردد. مثال: فرض می‌شود، از ۱۰ دانش‌آموز خواسته شده است که به یک آزمون تشریحی ۵ سؤالی پاسخ دهند. هر یک از سؤالات توسط ۳ ارزیاب، نمره‌گذاری شده است (حجم نمونه برای مثال کوچک در نظر گرفته شده است). هر دانش‌آموز (S) به تمام سؤالات (i) پاسخ داده و هر ارزیاب (F) نیز، هر سؤال هر دانش‌آموز را نمره‌گذاری کرده است.

۱. طرح سؤال: با توجه به هدف و تصمیم محقق از انجام یک مطالعه‌ی اندازه‌گیری، سؤالات طرح می‌شوند. از این رو یک محقق ممکن است اتکاپذیری نمرات را بسنجد، محقق دیگر، تعمیم‌پذیری نمرات را مدنظر داشته باشد یا هر دو هدف مذکور، ممکن است مطلوب محقق دیگر باشد. با توجه به اصل تقارن، می‌توان سؤالاتی در زمینه‌ی اعتبار هر یک از رویه‌های موجود (در مثال، ارزیابان و سؤالات) مطرح کرد. همچنین، سؤالاتی نیز در ارتباط با بهینه‌سازی طرح مانند؛ تعداد سطوح لازم از رویه‌های موجود برای رسیدن به ضرایب اعتبار دلخواه را می‌توان بیان نمود. در اینجا، سعی می‌شود اتکاپذیری و تعمیم‌پذیری نمرات به‌دست آمده برای دانش‌آموزان و همچنین، تعداد کافی ارزیابان و سؤالات، برای رسیدن به ضرایب دلخواه بررسی گردد. بعد از طرح سؤال، فرایند طراحی، تحلیل و تفسیر یک مطالعه اندازه‌گیری در قالب دو مرحله‌ی کلی؛ الف: مطالعه‌ی G (شامل مراحل ۲ تا ۸) و ب: مطالعه‌ی D (شامل مراحل ۹ تا ۱۵) صورت می‌پذیرد.

**الف: مطالعه‌ی G.** بعد از تعیین جامعه و جهان مشاهدات قابل قبول، مطابق با طرح اندازه‌گیری در نظر گرفته‌شده، داده‌ها جمع‌آوری می‌شوند. سپس با استفاده از نرم‌افزار یا انجام محاسبات دستی، می‌توان هدف مطالعه‌ی G را که برآورد مؤلفه‌های واریانس برای جامعه و جهان

مشاهدات قابل قبول است، تحقق بخشید و در ادامه به تفسیر مؤلفه‌های واریانس برآورد شده پرداخت.

۲. تعیین جامعه.

۲-۱. تعیین هدف اندازه‌گیری: همان‌طور که بیان شد، واژه جامعه برای اهداف اندازه‌گیری به کار می‌رود. هدف اندازه‌گیری در مثال ذکر شده، دانش‌آموزان (s) هستند که منبع خطا نیست، در نتیجه رویه حساب نمی‌شود، لذا تغییرپذیری میان دانش‌آموزان مطلوب است.

۲-۲. وضعیت نمونه‌گیری. در این گام فرعی، به تعداد اهداف اندازه‌گیری و روش نمونه‌گیری از آن‌ها اشاره می‌شود. در مثال، فرض شده است که  $n_s = 10$  (تعداد دانش‌آموزان) نفر یک نمونه‌ی تصادفی از جامعه نامحدود دانش‌آموزان هستند.

۳. تعیین جهان مشاهدات قابل قبول. با انجام گام‌های فرعی زیر می‌توان جهان مشاهدات قابل قبول که شامل همه ترکیبات ممکن از سطوح رویه‌ها است (وب و شولسون، ۲۰۰۵) را به‌طور دقیق تعریف کرد.

۳-۱. شناسایی رویه‌ها. در هر بافت اندازه‌گیری، منابع بالقوه‌ی خطا به‌عنوان رویه در نظر گرفته می‌شوند. بنا به گفته‌ی برنان (۲۰۰۱)، GT می‌تواند هر تعداد از رویه‌ها را لحاظ کند. در مثال مطرح شده، جهان مشاهدات قابل قبول، شامل دورویه‌ی سؤال و ارزیاب است.

۳-۲. وضعیت نمونه‌گیری (تعداد سطوح و تصادفی / ثابت). در این گام، به تعداد سطوح هر یک از رویه‌های موجود در جهان مشاهدات قابل قبول اشاره می‌شود. همچنین، باید تعیین شود که رویه‌ها متعلق به جهان‌های محدود هستند یا نامحدود. به بیانی دیگر، رویه‌ها تصادفی‌اند یا ثابت. طبق گفته‌ی برنان (۲۰۱۱b)، باید از هر یک از رویه‌ها حداقل دو سطح وجود داشته باشد. فرض شده است که  $n_i = 5$  سؤال، یک نمونه تصادفی از جهان سؤالات مشابه ممکن هستند. همچنین فرض شده است که  $n_r = 3$  ارزیاب، یک نمونه تصادفی از جهان ارزیابان مشابه ممکن هستند. یا می‌توان گفت؛ رویه‌های مذکور تصادفی در نظر گرفته می‌شوند، زیرا سطوح مشاهده نشده‌ی آن‌ها قابلیت جابه‌جایی با سطوح مشاهده شده در مطالعه اندازه‌گیری را دارند. مثال ما، فاقد رویه‌ی ثابت است.

۳-۳. تعیین روابط. در این گام، روابط بین رویه‌ها در جهان مشاهدات قابل قبول و همچنین، روابط جامعه و جهان مشاهدات قابل قبول تعیین می‌شود. هر یک از سؤالات در جهان مشاهدات قابل قبول می‌توانند توسط هر یک از ارزیابان در جهان ارزیابی شوند. در این صورت گفته می‌شود که دورویه در جهان مشاهدات قابل قبول، متقاطع ( $i \times r$ ) هستند. چون هر دانش آموز در جامعه می‌تواند به هر سؤال در جهان پاسخ دهد که توسط هر ارزیاب در جهان ارزیابی می‌شود، گفته می‌شود که جامعه و جهان مشاهدات قابل قبول متقاطع هستند که به صورت  $S \times i \times r$  نشان داده می‌شود.

۴. نام طرح اندازه‌گیری. در این گام، مشخصات طرح اندازه‌گیری در مطالعه  $G$  ذکر می‌شود. طرح اندازه‌گیری برای مثالمان، طرح دورویه‌ای متقاطع با رویه‌های تصادفی نام دارد و با نماد  $S \times i \times r$  نشان داده می‌شود.

۵. جمع‌آوری داده‌ها. محقق می‌بایست مطابق با طرح اندازه‌گیری مشخص شده، به جمع‌آوری داده‌ها بپردازد. البته گاهی محقق در وضعیت اندازه‌گیری انجام شده‌ای قرار می‌گیرد که می‌بایست با در نظر گرفتن ملاحظات منطقی و عملی برای داده‌های موجود، طرح اندازه‌گیری مناسبی را تهیه کند.

۶. انتخاب نرم‌افزار. محقق می‌تواند با انتخاب یک نرم‌افزار مناسب، نیاز به محاسبات دستی را مرتفع کرده و به روند تحلیل داده‌ها سرعت و دقت ببخشد. نکته‌ای که باید به آن توجه داشت، این است که نرم‌افزارهای موجود برای انجام تحلیل‌های تعمیم‌پذیری، قابلیت‌های متفاوتی دارند. برای مثال، می‌توان به توانایی آن‌ها در انجام مطالعات  $D$ ، به کارگیری طرح‌های (متعادل-نامتعادل)، طرح‌های (تک متغیری-چند متغیری)، به کارگیری داده‌های گمشده<sup>۱</sup>، نحوه‌ی ورود داده‌ها، سهولت کاربرد و ... اشاره کرد. در این قسمت، شرح مختصری از رایج‌ترین نرم‌افزارها بیان می‌شود. برنان (۲۰۰۱)، از سه برنامه نرم‌افزاری نام می‌برد که هر کدام ویژگی‌های متفاوتی دارند؛ GENOVA، urGENOVA و mGENOVA. دو برنامه‌ی اول، مختص به طرح‌های اندازه‌گیری تک متغیری و سومی

---

1. missing data

مختص طرح‌های اندازه‌گیری چند متغیری است. GENOVA، برای برآورد مؤلفه‌های واریانس اثرات تصادفی طرح‌های متعادل تا حداکثر ۶ اثر (۵ رویه و ۱ هدف اندازه‌گیری) مناسب است. ورود داده‌ها به این نرم‌افزار به دو صورت (داده‌های خام یا میانگین مجذورات) امکان‌پذیر است. urGENOVA برای برآورد مؤلفه‌های واریانس اثرات تصادفی هر دو طرح متعادل و نامتعادل مناسب است. در نرم‌افزار مذکور، ورود داده‌ها تنها به صورت داده‌های خام امکان‌پذیر است. mGENOVA، نیز برای مجموعه‌ای از طرح‌های متعادل یا نامتعادل (با توجه به آشنایی‌ای بودن) بکار می‌رود. ورود داده‌ها در mGENOVA می‌تواند به دو شیوه (داده‌های خام یا ماتریس واریانس و کوواریانس) انجام گیرد. از دیگر نرم‌افزارهای موجود برای انجام تحلیل‌های تعمیم‌پذیری می‌توان از EDUG نام برد که مبتنی بر تحلیل واریانس است و برای طرح‌های متعادل تا حداکثر ۸ اثر (۷ رویه و ۱ هدف اندازه‌گیری) مناسب است. ورود داده‌ها به این نرم‌افزار به دو شیوه (داده‌های خام، مجموع مجذورات) انجام می‌گیرد. برخلاف دیگر نرم‌افزارها، قابلیت اجرای اصل تقارن با این نرم‌افزار وجود دارد و همچنین، هدف اندازه‌گیری را می‌توان به‌عنوان آشنایی شده در رویه‌های دیگر در نظر گرفت. در بین نرم‌افزارهای فوق‌الذکر، GENOVA و mGENOVA و EDUG قابلیت اجرای مطالعه‌ی D را دارند. با وجود این، نرم‌افزارهای نامبرده در صورت کامل نبودن داده‌ها (وجود داده‌های گمشده) قادر به انجام تحلیل نیستند. لذا مشکل مقادیر گمشده، قبل از ورود داده‌ها به این نرم‌افزارها باید برطرف شود. سهولت دسترسی به هر چهار نرم‌افزار از ویژگی‌های مشترک آنهاست (برنان، ۲۰۰۱، a، ۲۰۰۱؛ کاردینت و همکاران، ۲۰۱۰). همچنین، دو نرم‌افزار SPSS VARCOMP و SAS VARCOMP نیز از سوی محققان برای تحلیل‌های تعمیم‌پذیری به کار می‌رود (بریچ و همکاران، ۲۰۱۴).

۷. برآورد مؤلفه‌های واریانس مطالعه‌ی G. در طرح اندازه‌گیری مطالعه‌ی G یعنی؛  $S \times I \times T$ ، نمره کل مشاهده‌شده هر دانش‌آموز در یک سؤال که توسط یک ارزیاب داده شده است، بر حسب مدل خطی زیر بیان می‌شود:

$$X_{sir} = \mu$$

میانگین اصلی  $\mu$



$$\begin{aligned}
 & \text{اثر دانش آموز } \mu - \mu_S + \\
 & \text{اثر سؤال } \mu - \mu_i + \\
 & \text{اثر ارزیاب } \mu - \mu_r + \\
 & \text{اثر دانش آموز } \times \text{ سؤال } \mu - \mu_S - \mu_i + \mu_{Si} + \\
 & \text{اثر دانش آموز } \times \text{ ارزیاب } \mu - \mu_r - \mu_S + \mu_{Sr} + \\
 & \text{اثر سؤال } \times \text{ ارزیاب } \mu - \mu_r - \mu_i + \mu_{ir} + \\
 & \text{اثر باقیمانده } \mu - \mu_S - \mu_r - \mu_i + \mu_{Si} + \mu_{Sr} - \mu_{ir}
 \end{aligned}$$

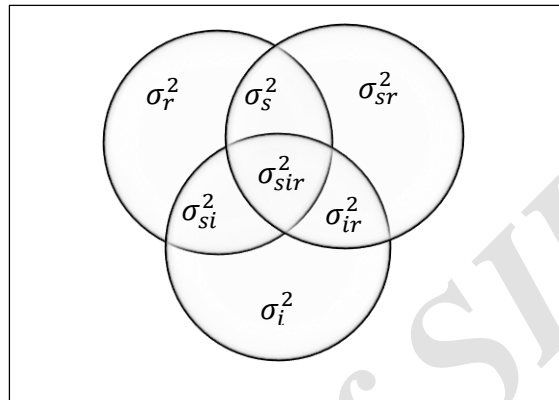
نمره‌ی مشاهده شده، شامل یک اثر برای میانگین اصلی؛ سه اثر برای دانش آموز، سؤال و ارزیاب؛ سه اثر تعاملی دو طرفه (SI، Si، Sr) و اثر باقی مانده (تعامل سه طرفه SIR به علاوه‌ی خطای تصادفی) است. طبق گفته‌ی برنان (۲۰۰۳)، همه‌ی اثرات (به‌غیر از  $\mu$ ) اثرات تصادفی نامیده می‌شوند؛ زیرا آن‌ها مرتبط با روش نمونه‌گیری تصادفی از جامعه و جهان هستند. همچنین، فرض می‌شود که همه‌ی اثرات در مدل، ناهمبسته هستند. توزیع هر مؤلفه یا اثر به‌غیر از میانگین اصلی، یک میانگین صفر و واریانس  $\sigma^2$  (بنام مؤلفه واریانس) دارد. فرض کنید که برای هر دانش آموز در جامعه، نمره‌ی میانگین (نمره مورد انتظار) در تمام سؤالات و همه‌ی ارزیابان در جهان مشاهدات قابل قبول به دست آمده است. واریانس این نمرات میانگین در جامعه‌ی دانش آموزان، واریانس نمره‌ی جهان (واریانس هدف اندازه‌گیری) نام دارد که آن را با نماد  $\sigma_{x_{sir}}^2$  نشان می‌دهند. سایر مؤلفه‌های واریانس برآورد شده، واریانس خطا را تشکیل می‌دهند (شولسون و وب، ۱۹۹۱؛ برنان، ۲۰۰۱؛ وب و شولسون، ۲۰۰۵).

واریانس نمره کل مشاهده شده، در سراسر جامعه‌ی افراد و سطوح در جهان مشاهدات قابل قبول به هفت مؤلفه واریانس مستقل مطابق زیر تجزیه می‌شود:

$$\sigma_{(x_{sir})}^2 = \sigma_S^2 + \sigma_i^2 + \sigma_r^2 + \sigma_{Si}^2 + \sigma_{Sr}^2 + \sigma_{ir}^2 + \sigma_{sir}^2$$

باید در نظر داشت که تعداد مؤلفه‌های واریانس، بستگی به طرح اندازه‌گیری دارد. تجزیه واریانس نمره‌ی کل مشاهده شده برحسب مؤلفه‌های واریانس تشکیل دهنده آن، از

طریق نمودار ون نشان داده شده است. لازم به ذکر است که اندازه‌ی دایره‌ها با نواحی درون آنها، نشان‌دهنده‌ی اندازه‌ی مؤلفه‌های واریانس نیست (شولسون و وب، ۱۹۹۱).



نمودار ۱. مؤلفه‌های واریانس مربوط به طرح اندازه‌گیری  $S \times I \times T$

مؤلفه‌های واریانس را می‌توان با استفاده از ترکیب نظری میانگین مجذورات که میانگین مجذورات مورد انتظار (EMS) نامیده می‌شود، برآورد کرد. بدین صورت که با استفاده از روش تحلیل واریانس (روش‌های مختلف دیگری نیز برای برآورد مؤلفه‌های واریانس وجود دارد که در این مقاله به آنها اشاره شده است)، میانگین مجذورات (MS) برای هر مؤلفه محاسبه می‌شود که با قرار دادن آن در معادلات EMS مربوط به آن مؤلفه، مؤلفه واریانس برآورد شده<sup>۱</sup> (EVC) به دست می‌آید.

در جدول شماره‌ی ۲، فرمول‌های مربوط به برآورد هر مؤلفه واریانس مطالعه G برای طرح اندازه‌گیری  $S \times I \times T$  ارائه شده است. هر طرح اندازه‌گیری، معادلات EMS مخصوص به خود را دارد.

## 2. Estimated variance components

جدول ۲. فرمول‌های مربوط به برآورد مؤلفه‌های واریانس مطالعه G برای طرح اندازه‌گیری  $s \times i \times r$

منبع تغییر	نوع تغییر پذیری	MS	EMS	EVC
s	نمره جهان	$MS_s$	$\sigma_{sir}^2 + n_i \sigma_{sr}^2 + n_r \sigma_{si}^2 + n_i n_r \sigma_s^2$	$\hat{\sigma}_s^2 = \frac{MS_s - MS_{si} - MS_{sr} + MS_{sir}}{n_i n_r}$
i	خطا	$MS_i$	$\sigma_{sir}^2 + n_s \sigma_{ir}^2 + n_r \sigma_{si}^2 + n_s n_r \sigma_i^2$	$\hat{\sigma}_i^2 = \frac{MS_i - MS_{si} - MS_{ir} + MS_{sir}}{n_s n_r}$
r	خطا	$MS_r$	$\sigma_{sir}^2 + n_s \sigma_{ir}^2 + n_i \sigma_{sr}^2 + n_s n_i \sigma_r^2$	$\hat{\sigma}_r^2 = \frac{MS_r - MS_{sr} - MS_{ir} + MS_{sir}}{n_s n_i}$
si	خطا	$MS_{si}$	$\sigma_{sir}^2 + n_i \sigma_{si}^2$	$\hat{\sigma}_{si}^2 = \frac{MS_{si} - MS_{sir}}{n_r}$
sr	خطا	$MS_{sr}$	$\sigma_{sir}^2 + n_s \sigma_{sr}^2$	$\hat{\sigma}_{sr}^2 = \frac{MS_{sr} - MS_{sir}}{n_i}$
ir	خطا	$MS_{ir}$	$\sigma_{sir}^2 + n_s \sigma_{ir}^2$	$\hat{\sigma}_{ir}^2 = \frac{MS_{ir} - MS_{sir}}{n_s}$
sir	خطا	$MS_{sir}$	$\sigma_{sir}^2$	$\hat{\sigma}_{sir}^2 = MS_{sir}$

در جدول شماره‌ی ۳، مؤلفه‌های واریانس برای هدف اندازه‌گیری، هر یک از رویه‌ها و همچنین، برای تعاملات موجود در طرح اندازه‌گیری  $s \times i \times r$ ، برآورد شده است. آنچه در جدول برآوردها باید به آن توجه داشت و راهکار مناسبی برای آن در نظر گرفت، منفی بودن برآوردهاست. همان‌طور که در جدول پیداست، اندازه‌ی مؤلفه‌ی واریانس مربوط به اثر تعاملی دانش‌آموز با ارزیاب، منفی شده است. به سبب کوچک بودن این اندازه (نزدیک به صفر)، صفر را جایگزین آن می‌کنیم. این نکته حائز اهمیت است که این مؤلفه‌ها، برای نمرات هر فرد در یک سؤال و یک ارزیاب است.

جدول ۳. برآوردهای آنوا از مؤلفه‌های واریانس مطالعه‌ی G برای طرح اندازه‌گیری  $s \times i \times r$

منبع تغییر	ss	df	MS	EVC	%	SE
s	۰/۴۶۲	۹	۰/۰۵۱	$\hat{\sigma}_s^2 = \frac{0.051 - 0.012 - 0.005 + 0.0114}{5 \times 3} = 0.0030$	۱۶/۸۵	۰/۰۰۱۵
i	۰/۳۵۵	۴	۰/۰۸۹	$\hat{\sigma}_i^2 = \frac{0.089 - 0.012 - 0.012 + 0.0114}{10 \times 3} = 0.0025$	۱۴/۰۴	۰/۰۰۱۷
r	۰/۰۷۰	۲	۰/۰۳۵	$\hat{\sigma}_r^2 = \frac{0.035 - 0.005 - 0.012 + 0.0114}{10 \times 5} = 0.0006$	۳/۳۷	۰/۰۰۰۵
si	۰/۴۲۷	۳۶	۰/۰۱۲	$\hat{\sigma}_{si}^2 = \frac{0.012 - 0.0114}{3} = 0.0002$	۱/۱۲	۰/۰۰۱۱
sr	۰/۰۹۶	۱۸	۰/۰۰۵	$\hat{\sigma}_{sr}^2 = \frac{0.005 - 0.0114}{5} = -0.0013 = 0$	۰	۰/۰۰۰۵
ir	۰/۰۹۶	۸	۰/۰۱۲	$\hat{\sigma}_{ir}^2 = \frac{0.012 - 0.0114}{10} = 0.0001$	۰/۵۶	۰/۰۰۰۶
sir	۰/۸۲۰	۷۲	۰/۰۱۱۴	$\hat{\sigma}_{sir}^2 = 0.0114$	۶۴/۰۴	۰/۰۰۱۹

۸. تفسیر مؤلفه‌های واریانس. مؤلفه‌های واریانس برآورد شده به مقیاس اندازه‌گیری بکار رفته، وابسته‌اند. از این رو، بر مبنای اندازه‌ی نسبی‌شان تفسیر می‌شوند. برای تفسیر کردن، می‌توان هر یک از مؤلفه‌ها را جداگانه تفسیر کرد و به سهم درصدی هر مؤلفه از واریانس کل (مجموع مؤلفه‌های واریانس) اشاره نمود. همچنین، با یک مقایسه‌ی نسبی مشخص کرد، کدام رویه‌ها و یا تعاملاتشان بیشترین سهم تغییرپذیری را به خود اختصاص داده‌اند. از دیگر شیوه‌های تفسیر، می‌توان به مقایسه‌های دوبه‌دوی مؤلفه‌ها و یا استفاده از انحراف استاندارد اشاره کرد (شولسون و وب، ۱۹۹۱؛ برنان، ۲۰۰۱). در مثال مطرح‌شده، مؤلفه واریانس باقی‌مانده  $\hat{\sigma}_{\text{SIR}}^2$  بزرگ‌ترین منبع تغییرپذیری است به طوری که ۶۴/۰۴ درصد از واریانس کل را به خود اختصاص داده است. این مؤلفه، نشان می‌دهد که سهم بزرگی از واریانس ناشی از؛ تعامل سه راهه بین دانش‌آموزان، سؤالات و ارزیابان و دیگر منابع تغییرپذیری منظم و تصادفی است که در مطالعه، اندازه‌گیری نشده‌اند. در مرتبه‌ی دوم، بیشترین سهم تغییرپذیری به مؤلفه واریانس برای دانش‌آموز یا واریانس نمره جهان  $\hat{\sigma}_{\text{S}}^2$  اختصاص دارد که سهم آن از واریانس کل، ۱۶/۸۵ درصد است. این مؤلفه، تغییرپذیری منظم میان دانش‌آموزان را در توانایی اندازه‌گیری شده نشان می‌دهد که اندازه‌ی بزرگ آن مطلوب یک محقق است. در مرتبه‌ی سوم، سؤالات (با اختصاص دادن ۱۴/۰۴ درصد از واریانس کل) قرار دارند. مؤلفه واریانس سؤالات  $\hat{\sigma}_{\text{I}}^2$  نسبتاً بزرگ است و نشان می‌دهد که سؤالات از بعد دشواری متفاوت هستند. مؤلفه واریانس کوچک برای تعامل دانش‌آموز با سؤال  $\hat{\sigma}_{\text{SI}}^2$  نشان می‌دهد جایگاه نسبی افراد از سؤالی به سؤال دیگر تفاوت چندانی نمی‌کند. مؤلفه‌ی واریانس مربوط به ارزیاب و همه‌ی مؤلفه‌های مرتبط با ارزیاب نسبتاً کوچک هستند. سهم نسبتاً کوچک مؤلفه واریانس ارزیاب  $\hat{\sigma}_{\text{I}}^2$  از واریانس کل (۳/۳۷٪) نشان می‌دهد که ارزیابان در یک مقیاس سخت‌گیری - سهل‌گیری نسبتاً مشابه عمل کرده‌اند. به عبارتی دیگر، میانگین نمره‌گذاری‌هایشان نسبتاً نزدیک به هم بوده است. مطلوب است که سهم رویه ارزیاب از واریانس کوچک باشد. تعامل دانش‌آموز با ارزیاب  $\hat{\sigma}_{\text{SR}}^2$  سهمی در تغییرپذیری نمرات دانش‌آموزان ندارد. بدین معناست که سخت‌گیری - سهل‌گیری یک ارزیاب در بین

دانش آموزان ثبات داشته است. تعامل ارزیابان با سؤالات  $\hat{\sigma}_{IR}^2$  (۵۶/۰٪ از واریانس کل) نشان می‌دهد که ارزیابان در نمره‌گذاری سؤالات مختلف آزمون تقریباً مشابه عمل کرده‌اند. روش دیگر برای تفسیر، مقایسه‌های دوبه‌دوی مؤلفه‌هاست. برای مثال می‌توان گفت؛ واریانس برآورد شده ناشی از سؤالات تقریباً ۴ برابر واریانس ناشی از ارزیابان است که نشان می‌دهد سؤالات نسبت به ارزیابان به‌طور قابل ملاحظه‌ای منبع بزرگ‌تری از تغییرپذیری در نمرات دانش آموزان هستند. به بیان دیگر، تفاوت سؤالات در میانگین دشواری خیلی بیشتر از تفاوت ارزیابان در میانگین سخت‌گیری-سهل‌گیری است.

**ب: مطالعه‌ی D.** همان‌طور که برنان (۲۰۰۱) عنوان کرده است، این مطالعات بر برآورد، استفاده و تفسیر مؤلفه‌های واریانس برای تصمیم‌گیری در مورد اهداف اندازه‌گیری با روش‌های اندازه‌گیری مشخص شده‌ای تأکید می‌کند. بنا به گفته‌ی مارکولیدس (۱۹۹۶) نیز، از مؤلفه‌های واریانس برآورد شده در مطالعه‌ی G می‌توان جهت برآورد واریانس نمره جهان، واریانس‌های خطا و ضرایب اعتبار (تعمیم‌پذیری) برای جهان تعمیم و طرح‌های مطالعه‌ی D استفاده کرد.

۹. تعیین جهان تعمیم. مهم‌ترین مسئله‌ی مطالعه‌ی D، تعیین جهان تعمیم است. در این گام، مشخص می‌کنیم که جهان تعمیم (جهانی که علاقه‌مند به تعمیم دادن نتایج به آن هستیم) با جهان مشاهدات قابل قبول یکسان است و یا زیرمجموعه‌ای از رویه‌ها و سطوح حشان در جهان مشاهدات قابل قبول است. در مورد اخیر، می‌بایست مشخصات جهان تعمیم به‌وضوح بیان شود. در مثال، جهان تعمیم شامل همه ارزیابان و سؤالات در جهان مشاهدات قابل قبول است. از آنجا که هر دورویه نامحدود (تصادفی) فرض شده‌اند، جهان تعمیم نیز نامحدود در نظر گرفته می‌شود.

۱۰. طرح اندازه‌گیری. مطالعه‌ی D می‌تواند طرح اندازه‌گیری مطالعه‌ی G و یا طرح‌های دیگری (در صورت امکان) را در نظر بگیرد. در مورد اخیر، می‌بایست طرح اندازه‌گیری مجدداً تعیین شود. مؤلفه‌های واریانس برآورد شده مطالعه‌ی G، بر روی نمرات واحد دانش آموز-سؤال - ارزیاب تمرکز می‌کند، درحالی که مطالعه‌ی D بر نمرات میانگین تمرکز می‌کند. به همین خاطر در ادبیات برنان، برای رویه‌های بکار رفته در طرح مطالعه‌ی D از حروف بزرگ

استفاده می شود. برای مثالمان، طرح اندازه گیری مطالعه‌ی D همان طرح مطالعه G است که به صورت  $s \times I \times R$  نشان داده می شود. در این طرح نیز، افراد هدف اندازه گیری اند.

۱۱. برآورد مؤلفه‌های واریانس مطالعه‌ی D. همان طور که در جدول شماره ۴ آمده است، جهت محاسبه‌ی مؤلفه‌های واریانس مطالعه‌ی D، از مؤلفه‌های واریانس برآورد شده‌ی مرحله‌ی ۷ (مؤلفه‌های واریانس برآورد شده‌ی مطالعه‌ی G) بدین صورت استفاده می کنیم که: مؤلفه مطالعه G بر  $\hat{n}_1$  تقسیم می شود، اگر مؤلفه شامل i باشد بی آنکه I را شامل شود. همچنین، بر  $\hat{n}_r$  تقسیم می شود، اگر مؤلفه شامل I باشد بی آنکه i را شامل شود. بر  $\hat{n}_1 \hat{n}_r$  تقسیم می شود، اگر شامل هر دو I و i باشد. به عنوان مثال، برای برآورد مؤلفه واریانس ارزیاب  $(\hat{\sigma}_R^2)$ ، مقدار برآورد شده‌ی این مؤلفه در مطالعه‌ی G  $(\hat{\sigma}_r^2 = 0/0006)$  تقسیم بر حجم نمونه ارزیاب در مطالعه D  $(\hat{n}_r = 3)$  می شود. قابل ذکر است که این قاعده‌ی تقسیم بر حجم نمونه، برای مؤلفه هدف اندازه گیری (واریانس نمره جهان) به کار نمی رود. مقدار این مؤلفه در هر دو مطالعه برابر است. همان طور که گفته شد، حجم نمونه برای مطالعه‌ی D  $(\hat{n}_r \text{ و } \hat{n}_1)$  می تواند همان حجم نمونه در مطالعه‌ی G  $(n_r \text{ و } n_1)$  باشد و یا متفاوت از آن باشد. در اینجا، همان حجم نمونه در مطالعه‌ی G در نظر گرفته شده است.

جدول ۴. معادلات و محاسبات مربوط به برآورد مؤلفه‌های واریانس مطالعه‌ی D

EVC در مطالعه‌ی G	EVC در مطالعه‌ی D
$\hat{\sigma}_s^2 = 0/0030$	$\hat{\sigma}_s^2 = 0/0030$
$\hat{\sigma}_i^2 = 0/0025$	$\hat{\sigma}_i^2 = \frac{\hat{\sigma}_i^2}{\hat{n}_1} = 0/0005$
$\hat{\sigma}_r^2 = 0/0006$	$\frac{\hat{\sigma}_r^2}{\hat{n}_r} = 0/0002$
$\hat{\sigma}_{si}^2 = 0/0002$	$\hat{\sigma}_{si}^2 = \frac{\hat{\sigma}_{si}^2}{\hat{n}_1} = 0/0004$
$\hat{\sigma}_{sr}^2 = -0/0013 = 0$	$\hat{\sigma}_{sr}^2 = \frac{\hat{\sigma}_{sr}^2}{\hat{n}_r} = 0$
$\hat{\sigma}_{ir}^2 = 0/0001$	$\hat{\sigma}_{ir}^2 = \frac{\hat{\sigma}_{ir}^2}{\hat{n}_1 \hat{n}_r} = 0/00001$
$\hat{\sigma}_{sir}^2 = 0/0114$	$\hat{\sigma}_{sir}^2 = \frac{\hat{\sigma}_{sir}^2}{\hat{n}_1 \hat{n}_r} = 0/00076$
$n_i = \hat{n}_i = 5 \quad n_r = \hat{n}_r = 3$	

۱۲. نوع تفسیر یا نوع تصمیم. در یک وضعیت اندازه گیری، نوع تفسیر نمره (هنجار مرجع در مقابل ملاک مرجع) تعیین می کند که کدام تصمیم (نسبی یا مطلق) مناسب است (فن و سان، ۲۰۱۳). در این گام، نوع تصمیم و یا نوع تفسیر باید مشخص شود. در تصمیم نسبی،

تصمیم‌گیری درباره‌ی دانش‌آموزان، با توجه به جایگاهشان در ارتباط با عملکرد دیگر دانش‌آموزان انجام می‌گیرد. در تصمیم مطلق، تصمیم‌گیری بر اساس عملکرد هر دانش‌آموز بر اساس ملاک تعیین‌شده صورت می‌گیرد. در مثال، بر اساس گام ۱، هر دو نوع تصمیم در نظر گرفته می‌شوند.

۱۳. محاسبه واریانس‌های خطا و ضرایب مربوطه به همراه تفسیر. هر یک از دو تصمیم نسبی و مطلق، واریانس‌های خطا و ضرایب مخصوص به خود را دارند. در نمودار شماره‌ی ۲، تفاوت میان واریانس خطای نسبی و مطلق در طرح اندازه‌گیری  $s \times I \times R$  نشان داده شده است که قسمت‌های هاشور خورده‌ی آن، سهم واریانس خطا را تحت سطوح مختلف نشان می‌دهد. شکل الف، واریانس خطای نسبی  $\sigma_{\delta}^2$  را نشان می‌دهد که شامل همه‌ی مؤلفه‌های واریانس تعاملی است که هدف اندازه‌گیری را در برمی‌گیرد. شکل ب، مربوط به واریانس خطای مطلق  $\sigma_{\Delta}^2$  است که شامل همه مؤلفه‌های واریانس مدل به‌جز هدف اندازه‌گیری است.



الف: خطای نسبی ب: خطای مطلق

نمودار ۲. سهم واریانس خطای نسبی و مطلق در طرح دورویه‌ای متقاطع

در جدول شماره‌ی ۵، معادلات و محاسبات واریانس‌های خطا و ضرایب مربوطه برای طرح  $s \times I \times R$  آورده شده است.

جدول ۵. معادلات و محاسبات مربوط به واریانس های خطا و ضرایب برای طرح  $s \times I \times R$

$\hat{\sigma}_\delta^2 = \hat{\sigma}_{SI}^2 + \hat{\sigma}_{SR}^2 + \hat{\sigma}_{SIR}^2 = 0.0004 + 0.00076 + 0.0008 = 0.00186$	واریانس خطای نسبی $(\hat{\sigma}_\delta^2)$
$E\hat{\rho}^2 = \frac{\hat{\sigma}_s^2}{\hat{\sigma}_s^2 + \hat{\sigma}_\delta^2} = \frac{0.0030}{0.0030 + 0.00186} = 0.79$	ضریب تعمیم پذیری $(E\hat{\rho}^2)$
$\hat{\sigma}_\Delta^2 = \hat{\sigma}_I^2 + \hat{\sigma}_R^2 + \hat{\sigma}_{SI}^2 + \hat{\sigma}_{SR}^2 + \hat{\sigma}_{IR}^2 + \hat{\sigma}_{SIR}^2 = 0.0005 + 0.0002 + 0.0004 + 0.00076 + 0.0001 + 0.00076 = 0.00276$	واریانس خطای مطلق $(\hat{\sigma}_\Delta^2)$
$\hat{\varphi} = \frac{\hat{\sigma}_s^2}{\hat{\sigma}_s^2 + \hat{\sigma}_\Delta^2} = \frac{0.0030}{0.0030 + 0.00276} = 0.67$	شاخص اتکاپذیری $(\hat{\varphi})$

تفسیر: همان طور که وب و همکاران (۲۰۰۷) مطرح کرده اند، برای تصمیم گیری در مورد افراد مبتنی بر نمرات مشاهده شده شان، ضریب اعتبار  $0.80$  و بالاتر غالباً به قدر کافی معتبر تلقی می شود و در صورتی که تصمیمات، پیامدهای چشمگیری داشته باشند، مقادیر  $0.90$  به بالاتر ترجیح داده می شود. با توجه به اینکه اندازه ی هر دو نوع ضریب، کمتر از حداقل سطح مطلوب ( $0.80$ ) است، می توان نتیجه گرفت؛ اندازه گیری نمرات دانش آموزان از اعتبار کافی برخوردار نیست. به عبارتی دیگر، ضریب تعمیم پذیری ( $0.79$ ) نشان می دهد که امکان تفکیک معتبر دانش آموزان از نظر توانایی اندازه گیری شده وجود ندارد. همچنین، شاخص اتکاپذیری ( $0.67$ ) نشان می دهد که مشخص کردن جایگاه یک دانش آموز در مقیاس توانایی اندازه گیری شده به صورت معتبر، امکان پذیر نیست.

۱۴. خطای استاندارد اندازه گیری و فواصل اطمینان. در استانداردهای سنجش تربیتی و روان-شناختی آمده است که همه ی برآوردهای اعتبار باید همراه خطاهای استاندارد اندازه گیری (SEM) و فواصل اطمینان باشد (بریچ و همکاران، ۲۰۱۴). SEM، کاربر را در مورد اندازه خطایی که نتایج را در بافت اندازه گیری نسبی و مطلق متأثر می سازد، مطلع می کند (کاردینت و همکاران، ۲۰۱۰). اگر فرض کنیم خطاها به طور طبیعی توزیع شده اند، با گرفتن ریشه دوم واریانس های خطای نسبی و مطلق به ترتیب SEM نسبی و مطلق به دست می آید. در مثال مفروض، برای داده هایی که در مقیاس  $0-5$  بودند،  $SEM = 0.282$  نسبی و  $0.387$  مطلق  $SEM =$  مطلق به دست آمد. هر چه مقدار SEM (نسبی یا مطلق) بزرگ تر باشد، دقت اندازه گیری کاهش می یابد و فاصله اطمینان نیز بزرگ تر خواهد بود. فواصل اطمینان رایج،



۹۵٪ و ۹۹٪ هستند که به ترتیب با نمرات  $Z$ ، ۱/۹۶ و ۲/۵۸ مرتبط هستند. برای محاسبه فاصله اطمینان دلخواه، مقدار SEM به دست آمده را در فرمول زیر جایگزین می‌کنیم.

$$\bar{X} \pm (Z \times SEM)$$

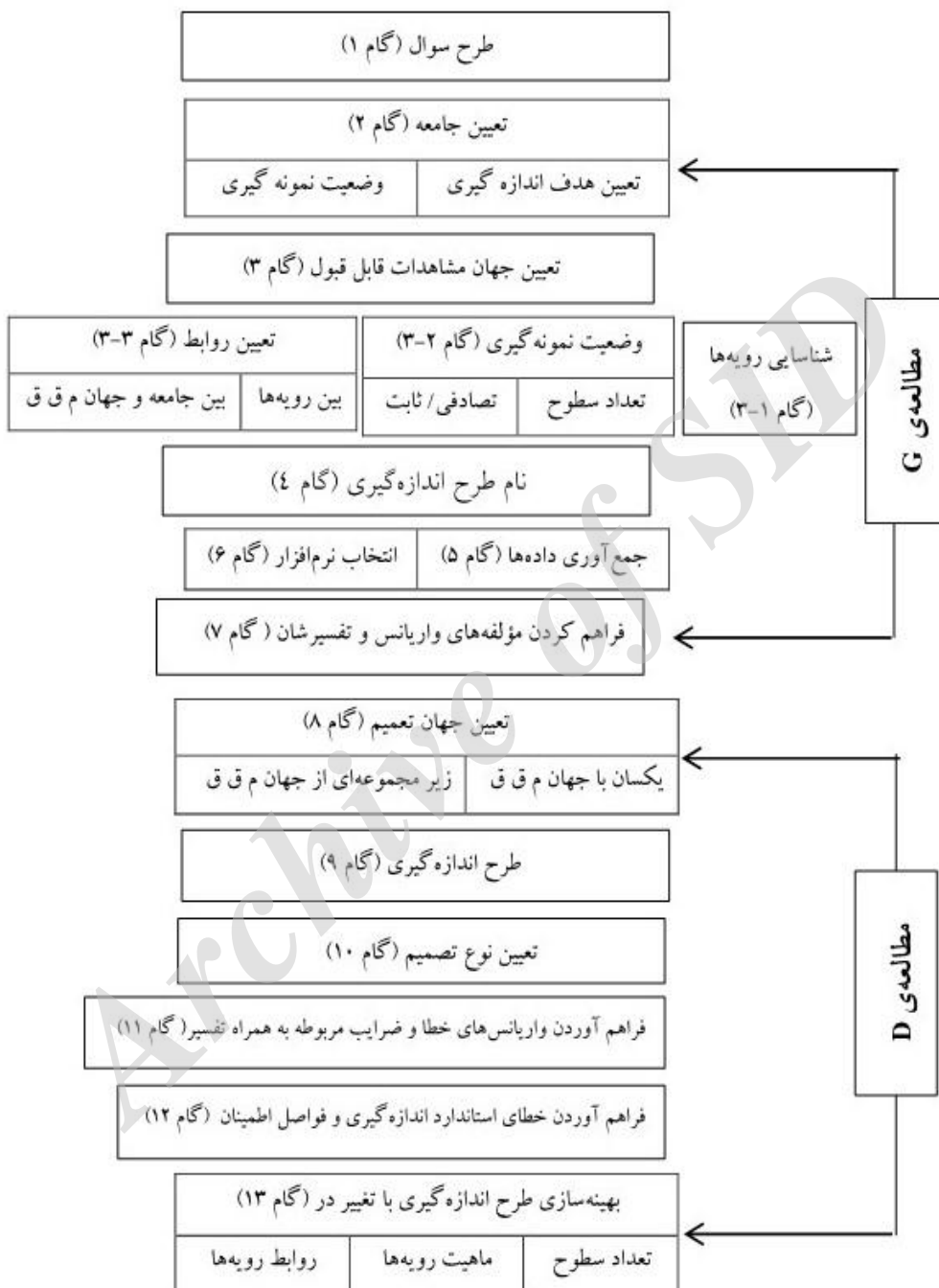
۱۵. بهینه‌سازی طرح اندازه‌گیری. محقق در این گام می‌تواند، با دست‌کاری رویه‌ها (مانند؛ تغییر در تعداد سطوح، ماهیت و روابط رویه‌ها) خطا را در طرح اندازه‌گیری کاهش دهد و اعتبار را به اندازه‌ی دلخواه برساند. در جدول شماره ۶، ضرایب تعمیم‌پذیری و اتکاپذیری برای چند طرح (با تغییر در تعداد سطوح رویه‌ها در طرح اندازه‌گیری  $s \times I \times R$ ) محاسبه شده است. همچنین، در جدول مذکور چگونگی استفاده از اطلاعات به دست آمده در مطالعه‌ی  $G$  برای به کارگیری در طرح‌های مختلف مطالعه‌ی  $D$  نشان داده شده است.

۱۵-۱. تغییر تعداد سطوح رویه‌ها. یکی از سؤالات مطرح شده در گام اول، بررسی تعداد سطوح لازم از ارزیابان و سؤالات برای رسیدن به ضرایب اعتبار دلخواه بود. همان‌طور که در قسمت تفسیر مؤلفه‌های واریانس مطالعه  $G$  گفته شد، سؤالات نسبت به ارزیابان، منبع بزرگ‌تری از تغییرپذیری در نمرات دانش‌آموزان هستند؛ بنابراین، افزودن تعداد سطوح سؤالات، واریانس خطا را بسیار بیشتر از افزودن ارزیابان کاهش می‌دهد که نتیجه‌ی آن، افزایش ضرایب اعتبار خواهد بود. با توجه به جدول شماره‌ی ۶، اگر تعداد ارزیابان را به یک نفر کاهش داده و تعداد سطوح سؤالات را ۲، ۳ و ۴ برابر کنیم، ضرایب تعمیم‌پذیری به دست آمده به ترتیب برابر با ۰/۷۳، ۰/۸۰ و ۰/۸۴ است. باین‌حال، ضرایب اتکاپذیری در هر سه مورد همچنان پایین‌تر از حداقل سطح مطلوب (۰/۸۰) است. اگر تعداد ارزیابان موجود در طرح ۲ نفر باشد و تعداد سطوح سؤالات را ۲، ۳ و ۴ برابر کنیم، ضرایب تعمیم‌پذیری به دست آمده به ترتیب برابر ۰/۸۴، ۰/۸۸ و ۰/۹۱ خواهد شد و ضرایب اتکاپذیری نیز، به ترتیب عبارت‌اند از؛ ۰/۷۲، ۰/۷۸ و ۰/۸۱. اگر تعداد سطوح ارزیابان ثابت نگه داشته شود و تعداد سطوح سؤالات به ۲، ۳ و ۴ برابر افزایش یابد، ضرایب تعمیم‌پذیری از ۰/۷۹ به ترتیب به ۰/۸۸، ۰/۹۲ و ۰/۹۴ و ضرایب اتکاپذیری نیز از ۰/۶۷ به ۰/۷۸، ۰/۸۲ و ۰/۸۵ خواهد رسید. تنها در حالت‌های (۲ ارزیاب و ۲۰ سؤال، ۳ ارزیاب و ۱۵ سؤال، ۳ ارزیاب و ۲۰ سؤال) هر دو ضرایب اعتبار با هم از حداقل سطح مطلوب و بالاتر از آن برخوردارند. از آنجا که طرح

سؤال و روند ارزیابی مستلزم صرف هزینه و زمان است، با انجام این گام می‌توان ترکیب مناسبی از سؤالات و ارزیابان را با توجه به اندازه اعتبار موردنظر و دیگر محدودیت‌های عملی، به دست آورد. بنا به گفته‌ی ون لیون (۱۹۹۷)، به لحاظ عملی، اجرا کردن سؤالات بیشتر که توسط یک ارزیاب نمره‌گذاری می‌شوند به‌صرفه‌تر از داشتن تعداد کمی سؤال با ارزیابان چندگانه است. از این رو، مناسب‌ترین طرح را می‌توان داشتن ۲ ارزیاب و ۲۰ سؤال دانست.

۲-۱۵. تغییر در ماهیت و روابط رویه‌ها. محقق در این گام فرعی می‌تواند با توجه به هدف و سؤالات در نظر گرفته‌شده و همچنین ملاحظات عملی و منطقی، ماهیت رویه‌ها (رویه ارزیاب یا رویه‌ی سؤال را ثابت در نظر بگیرد) یا روابط آن‌ها (مقاطع / آشیانه‌ای) را تغییر دهد و طرح‌های متنوع و متفاوتی را بکار بگیرد.

لازم به ذکر است که استفاده از نرم‌افزار نیاز به محاسبات دستی را مرتفع می‌کند که در این صورت، بعضی از مراحل ذکرشده در قسمت پیشین حذف می‌شوند. در این قسمت، فرآیند طراحی، تحلیل و تفسیر یک مطالعه اندازه‌گیری که به‌طور مفصل تشریح شد، در قالب نمودار شماره‌ی ۳ نشان داده شده است (البته با حذف مراحل مربوط به محاسبات دستی).



نمودار ۳: فرآیند طراحی، تحلیل و تفسیر یک مطالعه اندازه‌گیری

جدول ۶. مؤلفه‌های واریانس و ضرایب تعمیم‌پذیری برآورد شده برای طرح‌های مختلف مطالعه‌ی D

مطالعه‌ی G		مطالعه‌ی D								
$n_r=1$	$n_r$	۱			۲			۳		
$n_i=1$	$n_i$	۱۰	۱۵	۲۰	۱۰	۱۵	۲۰	۱۰	۱۵	۲۰
$\hat{\sigma}_s^2=0.0030$	$\hat{\sigma}_s^2$	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030
$\hat{\sigma}_f^2=0.0025$	$\hat{\sigma}_f^2$	0.0025	0.0017	0.00125	0.0025	0.0017	0.00125	0.0025	0.0017	0.00125
$\hat{\sigma}_R^2=0.0006$	$\hat{\sigma}_R^2$	0.0006	0.0006	0.0006	0.0003	0.0003	0.0003	0.0002	0.0002	0.0002
$\hat{\sigma}_{st}^2=0.0002$	$\hat{\sigma}_{st}^2$	0.0002	0.0001	0.0001	0.0002	0.0001	0.0001	0.0002	0.00013	0.0001
$\hat{\sigma}_{sr}^2=0$	$\hat{\sigma}_{sr}^2$	0	0	0	0	0	0	0	0	0
$\hat{\sigma}_{ir}^2=0.0001$	$\hat{\sigma}_{ir}^2$	0.0001	0.0001	0.00005	0.00005	0.00003	0.000025	0.00003	0.00002	0.00002
$\hat{\sigma}_{sif}^2=0.00114$	$\hat{\sigma}_{sif}^2$	0.0011	0.00076	0.00057	0.00057	0.00038	0.000285	0.00038	0.000253	0.00019
	$\hat{\sigma}_s^2$	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030	0.0030
	$\hat{\sigma}_{(s)}^2$	0.00112	0.0007	0.00058	0.00059	0.00039	0.000295	0.00039	0.000266	0.0002
	$\hat{\sigma}_{(s)}^2$	0.00198	0.00155	0.00131	0.001145	0.000863	0.000722	0.000853	0.000638	0.000527
	$Ep^2$	0.73	0.80	0.84	0.84	0.88	0.91	0.88	0.92	0.94
	$\varphi$	0.60	0.66	0.70	0.72	0.78	0.81	0.78	0.82	0.85

### بحث و نتیجه‌گیری

در حالی که از مدت‌ها پیش اکثر متخصصان اندازه‌گیری برای سنجش اعتبار، استفاده از GT را بجای CTT توصیه کرده‌اند (مشکواش و اکانرا، ۲۰۰۶)، محققان کشورمان صرفاً به CTT متکی هستند. مهم‌ترین دلیل آن را می‌توان، عدم آشنایی محققان و طراحان آزمون با مبانی نظری و قابلیت‌های GT دانست. از این رو، در مقاله حاضر، مبانی و مفاهیم GT تک متغیری (از جمله؛ جهان مشاهدات قابل قبول، جهان تعمیم، نقش مطالعات G و D، برآورد مؤلفه‌های واریانس، انواع طرح‌های تعمیم‌پذیری و انواع ضرایب) به سادگی و به‌دوراز ابهام معرفی و شرح داده شد. علاوه بر این، به برخی از شباهت‌ها و تفاوت‌های این نظریه با CTT اشاره گردید تا نقاط قوت و سودمندی GT آشکار گردد.

همان‌طور که می‌دانیم، یکی از چالش‌های عمده‌ای که هر وضعیت اندازه‌گیری با آن روبه‌روست، شناسایی و کاهش خطاهای اندازه‌گیری است. طبق بیان همبلتون و جونز<sup>۲</sup> (۱۹۹۳)، نظریه یا مدل مناسب به اینکه چگونه در یک وضعیت اندازه‌گیری سهم خطا را می‌توان به حداقل رساند، کمک می‌کند. برخلاف CTT، GT این قابلیت را دارد که منابع

1. Mushquash & O'connor
2. Hambleton & Jones

چندگانه خطای اندازه گیری را همزمان شناسایی، تفکیک و هر یک را برآورد کند. از این طریق می توان منبع (ها) خطایی که مقدار مؤلفه واریانس آن و یا مقدار مؤلفه های واریانس اثرهای تعاملی آن زیاد است را شناسایی کرد و در ادامه با راهکارهای بهینه سازی، مقدار واریانس ناشی از آن را کاهش داد و فرآیند اندازه گیری را بهبود بخشید. همان طور که فن و سان (۲۰۱۳) بیان نموده اند؛ قابلیت این نظریه در انجام طرح های مختلف اندازه گیری مطالعات D، به محققان این امکان را می دهد که اندازه اعتبار را در مطالعه واقعی شان پیش بینی کنند. چارچوب مفهومی محکم GT، تمایز گذاشتن میان تصمیم های نسبی و مطلق و به کارگیری اصل تقارن از دیگر نقاط قوت این نظریه است.

همچنین، در این مقاله تلاش شد که چگونگی انجام یک مطالعه اندازه گیری جهت سنجش اعتبار داده ها در قالب یک مثال نشان داده شود. در GT، طراحی دقیق یک مطالعه اندازه گیری، اجرا، تحلیل و تفسیر نتایج آن، انواع مختلفی از مسائل و ملاحظات را در برمی گیرد که محققان هنگام استفاده از این نظریه می بایست آن ها را در نظر داشته باشند تا بتوانند از سودمندی این نظریه در راستای سنجش اعتبار داده هایشان بهره ی کافی را ببرند. بدین منظور، فرآیند طراحی، تحلیل و تفسیر یک مطالعه اندازه گیری به صورت گام به گام در قالب مطالعات G و D توصیف شد. علاوه بر این، معادلات و محاسبات مربوطه نیز برای علاقه مندان تشریح گردید که به تسهیل درک آن ها از این نظریه کمک می کند.

لازم به ذکر است که برای برآورد اعتبار، از طرح های مختلفی می توان استفاده کرد که در این مقاله تنها به یکی از رایج ترین آن ها (طرح دورویه ای متقاطع با رویه های تصادفی) پرداخته شد. همچنین، با اینکه مثالمان در حوزه ی آموزش و پرورش بود، برای محققان سایر حوزه ها از جمله روان شناختی، پزشکی، بازاریابی، مدیریت و ... نیز مفید خواهد بود. در صورت استفاده از این طرح، محققان می توانند رویه های موجود در وضعیت های اندازه گیری شان را جایگزین رویه های به کاررفته در این مثال کنند و اعتبار داده هایشان را به دست بیاورند. کمبود منابع از مشکلات جدی محققان و علاقه مندان به کار در حیطه ی GT در کشورمان است. با انتشار مقالات، ترجمه و تألیف کتاب در زمینه ی GT از یک طرف و قرار گرفتن GT جزو برنامه آموزشی رشته ی سنجش و اندازه گیری و دیگر رشته های مرتبط

از طرف دیگر، می‌توان انتظار داشت که در آینده قابلیت‌ها و کاربردهای این نظریه از سوی محققان کشورمان به کار گرفته شود.

## منابع

کیامنش، علیرضا (۱۳۷۴). نظریه‌ی تعمیم‌پذیری در اندازه‌گیری آموزشی، نشریه‌ی علوم تربیتی، ۱۶(۱-۲): ۴۸-۲۵.

- Alkharusi, H. (2012). *Generalizability Theory: An Analysis of Variance Approach to Measurement Problems in Educational Assessment*, *Journal of Studies in Education*, 2(1):184-196.
- Bolus, R. E., Hinofotis, F. B., & Bailey, K. M. (1982). An introduction to generalizability theory in second language research. *Language Learning*, 32(2): 245-258.
- Brennan, R. L. (1992). Generalizability theory. *Educational Measurement: Issues and Practice*, 11(4): 27-34.
- Brennan, R.L. (1997). A perspective on the history of generalizability theory. *Educational Measurement: Issues and Practice*, 16: 14-20.
- Brennan, R. L. (2001). *Generalizability theory*. New York, NY: Springer-Verlag.
- Brennan, R. L. (2001a). mGENOVA [Computer software and manual]. Iowa City, IA: University of Iowa.
- Brennan, R. L. (2003). Coefficients and indices in generalizability theory. *Center for Advanced Studies in Measurement and Assessment, CASMA Research Report, 1*: 1-44.
- Brennan, R. L. (2010). *Generalizability Theory*. New York: Springer-Verlag.
- Brennan, R.L. (2011a). Generalizability theory and classical test theory. *Applied Measurement in Education*, 24(1): 1-21.
- Brennan, R. L. (2011b). Using Generalizability Theory to Address Reliability Issues for PARCC Assessments: A White Paper. *Center for Advanced Studies in Measurement and Assessment (CASMA), University of Iowa*.
- Briesch, A. M., Swaminathan, H., Welsh, M., & Chafouleas, S. M. (2014). Generalizability theory: A practical guide to study design, implementation, and interpretation. *Journal of school psychology*, 52(1): 13-35.
- Brown, J. D. (2005). Statistics corner, questions and answers about language testing statistics: Generalizability and decision studies. *Shiken: JALT Testing & Evaluation SIG Newsletter*, 9(1): 12-16.
- Bruckner, C. T., Yoder, P. J., & McWilliam, R. A. (2006). Generalizability and decision studies: An example using conversational language

- samples. *Journal of Early Intervention*, 28(2): 139-153. Cardinet, J., Johnson, S., & Pini, G. (2010). Applying Generalizability theory using Edug. Taylor & Francis Group.
- Cardinet, J., Tourneur, Y., & Allal, L. (1976). The symmetry of generalizability theory: Applications to educational measurement. *Journal of Educational Measurement*, 13: 119-135.
- Cronbach, L. J., Rajaratnam, N., & Gelser, G. C. (1963). Theory of generalizability: A liberalization of reliability theory. *The British Journal of Statistical Psychology*, 16: 137-163.
- Fan, X., & Sun, S. (2013). Generalizability theory as a unifying framework measurement of reliability in adolescent research. *The Journal of Early Adolescence*, 0272431613482044.
- Hambleton, R. K., & Jones, R. W. (1993). An NCME Instructional Module on. *Educational measurement: issues and practice*, 12(3): 38-47.
- Kumazawa, T. (2009). Revision of a criterion-referenced vocabulary test using generalizability theory. *JALT journal*, 31(1): 81-100.
- Marcoulides, G. A. (1996). Estimating variance components in generalizability theory: the covariance structure analysis approach. *Structural equation modeling*, 3(3): 290-299.
- Mushquash, C., & O'Connor, B. P. (2006). SPSS and SAS programs for generalizability theory analyses. *Behavior research methods*, 38(3): 542-547.
- Shavelson, R. J., & Webb, N. M. (1991). Generalizability theory: A primer. Newbury Park, CA: Sage.
- Shavelson, R. J., Webb, N. M., & Rowley, G. L. (1989). Generalizability theory. *American Psychologist*, 44: 922-932.
- Shavelson, R.J., & Webb, N.M. (1981). Generalizability theory: 1973 - 1980. *British Journal of Mathematical and Statistical Psychology*, 34:133-166.
- Suen, H. K., & Lei, P. W. (2007). Classical versus Generalizability theory of measurement. *Educational Measurement*, 4: 1-13.
- VanLeeuwen, D. M. (1997). Assessing reliability of measurements with generalizability theory: An application to inter-rater reliability. *Journal of Agricultural Education*, 38:36-42.
- Webb, N. M., Rowley, G. L., & Shavelson, R. J. (1988). Using generalizability theory in counseling and development. *Measurement and Evaluation in Counseling and Development*, 21: 81-90.
- Webb, N. M., Shavelson, R. J., & Haertel, E. H. (2007). Reliability coefficients and generalizability theory. *Handbook of statistics*, 26(4): 81-124.
- Webb, N.M., & Shavelson, R.J. (2005). Generalizability theory: Overview. *Encyclopedia of statistics in behavioral science*, 2, 717-719.