

## مرواری بر نظریه تعمیم‌پذیری و کاربردهای آن در آموزش پزشکی

\* سارا مرتاض‌هجری<sup>۱</sup>، لیلا جانانی<sup>۲</sup>، محمد جلیلی<sup>۳</sup>

۱- گروه آموزش پزشکی، دانشگاه علوم پزشکی تهران ۲- گروه آمار زیستی، دانشکده بهداشت، دانشگاه علوم پزشکی ایران ۳- گروه طب اورژانس، گروه آموزش پزشکی، دانشگاه علوم پزشکی تهران

\* نویسنده مسؤول: تهران، بلوار کشاورز، خیابان نادری، خیابان حجت دوست، پلاک ۵۷ طبقه سوم، گروه آموزش پزشکی. تلفن: ۰۲۱۸۸۹۵۵۰۱۶۰ نامبر: ۰۲۱۸۸۹۵۵۸۴۶ پست الکترونیک: mjalili@tums.ac.ir

دریافت: ۹۴/۱/۱۷ پذیرش: ۹۴/۳/۳۱

### چکیده

**مقدمه:** تمام اندازه‌گیری‌ها از جمله ارزیابی دانشجویان پزشکی در ذات خود واجد خطا هستند. با توجه به ضرورت حفظ کیفیت ارزیابی در علوم پزشکی، یافتن روش‌هایی برای برآورد میزان خطأ و کاهش آن حائز اهمیت است. طی سال‌های اخیر، پیشرفت‌های چشمگیری در زمینه نظریه‌های اندازه‌گیری صورت گرفته است که در این راه کمک‌کننده است.

**روش کار:** در این مقاله، پس از مرواری بر نظریه کلاسیک آزمون و محدودیت‌های آن، با استفاده از مثال‌هایی در خصوص آزمون عینی ساختارمند بالینی، نظریه تعمیم‌پذیری را بررسی می‌کنیم.

**یافته‌ها:** طبق نظریه کلاسیک آزمون، نمره‌ای که دانشجو در امتحان کسب کرده است، معادل نمره واقعی او نیست و با خطأ همراه است. با تعیین پایایی آزمون می‌توانیم بگوییم میزان تأثیر خطای تصادفی در نمره چقدر بوده است اما در این نظریه نمی‌توان در آن واحد، تمام منابع خطأ را در نظر گرفت. برای رفع این محدودیت از نظریه تعمیم‌پذیری استفاده می‌شود که بر اساس مدل‌های آماری آنالیز واریانس، سهم منابع مختلف خطأ را در اندازه‌گیری مشخص می‌کند و علاوه بر برآورد پایایی آزمون، پیشنهادهایی برای بهبود تعمیم‌پذیری نتایج ارایه می‌دهد.

**نتیجه‌گیری:** در نظریه کلاسیک آزمون، شناسایی میزان اثر منابع مختلف خطاهای احتمالی به تفکیک امکان‌پذیر نیست. در حالی که با استفاده از نظریه تعمیم‌پذیری می‌توان سهم هریک از منابع خطأ را شناسایی نمود و برای آزمونی پایاتر برنامه‌ریزی کرد.

**گل واژگان:** اندازه‌گیری، نظریه کلاسیک آزمون، نظریه تعمیم‌پذیری، پایایی

### مقدمه

الگو و جهت خاصی دارد. در حالی که خطای تصادفی کاملاً شناسی و بدون الگو و جهت مشخص، اندازه‌ها را تحت تأثیر قرار می‌دهد و میانگین آن در تکرار زیاد اندازه‌گیری‌ها نزدیک به صفر است. خطای سیستماتیک را می‌توان با شناسایی الگوهای مشخص آن حذف کرد و خطای تصادفی را می‌توان با تکرار دفعات اندازه‌گیری کاهش داد<sup>(۱)</sup>. آزمون‌های تحصیلی نیز که به نوعی دانش یا مهارت فراگیران را اندازه‌گیری می‌کنند، از این امر مستثنی نیستند و منابع مختلفی ممکن است باعث بروز خطأ

تمام شاخه‌های علم با هدف جمع‌آوری اطلاعات در مورد یک پدیده خاص نیازمند اندازه‌گیری هستند و تمام اندازه‌گیری‌ها در ذات خود واجد خطأ هستند و تحت تأثیر شرایط گوناگون، مقداری مختلفی به دست می‌دهند<sup>(۱)</sup>. خطاهای اندازه‌گیری به صورت کلی به دو نوع خطای سیستماتیک و تصادفی تقسیم می‌شوند. خطای سیستماتیک که نام دیگر آن سوگرایی<sup>(۱)</sup> است، به طور مستمر، منظم و مکرر مقداری را تحت تأثیر قرار می‌دهد و عموماً

<sup>۱</sup> Bias

نظریه کلاسیک آزمون: اساس و پایه نظریه کلاسیک بر این فرض استوار است که نمره مشاهده شده<sup>۶</sup> (X) دانشجو از دو جزء یعنی نمره واقعی<sup>۷</sup> وی (T) و خطای اندازه‌گیری<sup>۸</sup> (E) تشکیل شده است. به عبارت دیگر نمره‌ای که دانشجو در امتحان کسب کرده است، معادل نمره واقعی او نیست و با درجاتی از خطا همراه است (۱ و ۵). این مسأله به این صورت قابل نمایش است:

$$X = T + E$$

بر اساس معادله فوق، خطای اندازه‌گیری حاصل اختلاف بین نمره مشاهده شده و نمره واقعی فرد است. مفهوم نمره واقعی در معادله فوق، یک مفهوم فرضی است که از میانگین نمرات مشاهده شده در بینها تکرار اندازه‌گیری‌ها به دست می‌آید. فرمول زیر در خصوص واریانس این نمرات صدق می‌کند:

$$\sigma_X^2 = \sigma_T^2 + \sigma_E^2$$

در واقع تعریف رایج پایایی، یعنی تعیین میزان ثبات در اندازه‌های به دست آمده، از همین مفهوم گرفته شده است. مفهوم پایایی به این امر دلالت دارد که با تکرار آزمون می‌توان امیدوار بود که به نمره واقعی فرآگیر نزدیک شویم. به عبارت دیگر، هر چند هرگز نمی‌توانیم نمره واقعی دانشجو را دقیقاً به دست آوریم، با تعیین پایایی آزمون می‌توانیم بگوییم میزان تأثیر خطای تصادفی در نمره او چقدر بوده است و سپس با تخمین خطای معیار اندازه‌گیری<sup>۹</sup> می‌توانیم محدوده نمره واقعی دانشجو را برآورد کنیم (۱). ضرب پایایی در نظریه کلاسیک، نسبت واریانس نمره واقعی به واریانس نمره مشاهده شده (به شرط مستقل بودن خطا و نمره مشاهده شده) است و به صورت فرمول زیر نمایش داده می‌شود:

$$r = \frac{\sigma_T^2}{\sigma_X^2} = \frac{\sigma_T^2}{\sigma_T^2 + \sigma_E^2}$$

روش‌های آماری تعیین پایایی متنوع هستند و هر چند اصول کلی آن‌ها در حوزه ارزیابی فرآگیر همانند سایر حیطه‌های اندازه‌گیری است اما به علت شرایط خاص حاکم بر آزمون‌ها برخی از روش‌ها کمتر قابل استفاده هستند. به صورت مشخص‌تر دو روش بازآزمایی<sup>۱۰</sup> و آزمون‌های همارز<sup>۱۱</sup> علی‌رغم استفاده زیاد در تعیین پایایی پرسشنامه‌ها، کاربرد کمتری در

در آن‌ها شود. یکی از روش‌های ارزیابی فرآگیران در علوم پزشکی، آزمون عینی ساختارمند بالینی<sup>۲</sup> است. در این آزمون، دانشجو به ترتیب وارد ایستگاه‌های مختلف می‌شود و در هر یک از آن‌ها در زمان مشخصی یک وظیفه بالینی را در مواجهه با بیمارانما یا مولاژ انجام می‌دهد و آزمون گر مستقر در ایستگاه با مشاهده مستقیم عملکرد دانشجو بر اساس چک‌لیست‌های از پیش تدوین شده به ارزیابی وی می‌پردازد (۳). تأثیر هر عاملی به غیر از سطح علمی دانشجویان روی نمرات OSCE به عنوان خطای اندازه‌گیری محسوب می‌شود. به عنوان مثالی از خطای سیستماتیک در این آزمون، آزمونگری است که سخت‌گیر است و احتمالاً برای همه دانشجویان سخت‌گیر است و نمره تمام آن‌ها را کمتر از حد واقعی برآورد می‌کند. در اینجا آموزش دادن و توجیه وی می‌تواند کمک کننده باشد. در حالی که برای کاهش خطای تصادفی OSCE، افزایش تعداد آزمونگران یا افزایش تعداد ایستگاه‌های آزمون کمک کننده است.

با توجه به ضرورت حفظ کیفیت روش‌های ارزیابی در علوم پزشکی و اطمینان از عملکرد قابل قبول دانشجویان و فارغ‌التحصیلان، یافتن روش‌هایی برای برآورد میزان خطا و کاهش آن در اندازه‌گیری حائز اهمیت است (۴). به موازات تغییرات صورت گرفته در سایر علوم، پیشرفت‌های چشمگیری در زمینه نظریه‌های اندازه‌گیری<sup>۳</sup> صورت گرفته است که در ساخت و تحلیل آزمون‌ها و پرسشنامه‌ها کاربرد زیادی دارند. اولین نظریه منسجم در این زمینه تحت عنوان نظریه کلاسیک آزمون<sup>۴</sup> در اوایل قرن بیستم شکل گرفت (۵) و توسعه ابزارهای ارزیابی طبق اصول آن به سرعت گسترش یافت. در عین حال، برخی مطالعات، محدودیت‌هایی را در این نظریه نشان دادند و زمینه‌ساز ظهور نظریه‌های جدید از جمله نظریه تعمیم‌پذیری<sup>۵</sup> شدند که امکان شناسایی میزان خطای را در ابعاد مختلف اندازه‌گیری فراهم می‌سازد (۱ و ۲).

در این مقاله قصد داریم به صورت خلاصه مروری بر این دو نظریه اندازه‌گیری انجام دهیم. با توجه به منابع متعدد خطای OSCE، در این مقاله با ذکر نمونه‌هایی از این آزمون، ابتدا به بیان ویژگی‌ها و محدودیت‌های نظریه کلاسیک می‌پردازیم و سپس نظریه تعمیم‌پذیری را مورد بررسی قرار می‌دهیم.

<sup>6</sup> Observed score

<sup>7</sup> True score

<sup>8</sup> Measurement error

<sup>9</sup> Standard Error of Measurement (SEM)

<sup>10</sup> Test-retest

<sup>11</sup> Parallel tests

<sup>2</sup> Objective Structured clinical examination (OSCE)

<sup>3</sup> Measurement Theories

<sup>4</sup> Classical Test Theory (CTT)

<sup>5</sup> Generalizability Theory

آن‌ها می‌توان به EduG، GENOVA و G string اشاره کرد (7 و 9). در ابتدا به مرور چند مفهوم پایه‌ای این نظریه می‌پردازیم و سپس کاربردهای عملی آن را با ذکر چند مثال تبیین می‌کنیم:

رویه و انواع آن: نمره دانشجو تحت تأثیر خصوصیاتی از آزمون است که در نظریه تعمیم‌پذیری رویه<sup>19</sup> نامیده می‌شوند و هر یک دارای سطوح مشخصی هستند. به عنوان مثال، در آزمون OSCE، ایستگاه‌ها، آیتم‌های چک‌لیست‌ها، عملکرد بیمارنماها، عملکرد ارزیابان و... هر یک می‌توانند با ایجاد خطای اندازه‌گیری موجب شوند دانشجو نمره‌ای غیر از نمره واقعی خود کسب کند. رویه‌ای که مربوط به خود مورد اندازه‌گیری<sup>20</sup> است، مانند دانشجو در مثال آزمون OSCE، رویه تمیز<sup>21</sup> نامیده می‌شود زیرا تفاوتی که بین نمرات دانشجویان مختلف دیده می‌شود، ناشی از خطای نیست و منعکس‌کننده تفاوت واقعی در سطح آن‌ها است. به سایر رویه‌ها، مانند ایستگاه یا آزمون‌گر، رویه تعمیم<sup>22</sup> گفته می‌شود (7 و 8). کار نظریه تعمیم‌پذیری این است که مشخص کند نمره دانشجو تا چه حد به سایر ایستگاه‌ها یا آزمون‌گرانی که به صورت بالقوه می‌توانستند وجود داشته باشند، قابل تعمیم است. هر چه واریانس ناشی از دانشجویان بیشتر باشد، پایایی آزمون بیشتر خواهد بود و هر چه سهم رویه‌های تعمیم در پراکندگی نمرات بیشتر باشد، پایایی آزمون کمتر است (4).

جهان و نمره جهانی: در نظریه تعمیم‌پذیری به جای نمره واقعی از لفظ «نمره جهانی»<sup>23</sup>، برای مجموعه شرایط اندازه‌گیری از اصطلاح «جهان»<sup>24</sup> و برای کل مجموع افراد مورد اندازه‌گیری از اصطلاح «جمعیت»<sup>25</sup> استفاده می‌شود. پایه و اساس ادراکی<sup>26</sup> نظریه تعمیم‌پذیری را «جهان مشاهدات قابل قبول»<sup>27</sup> تشکیل می‌دهد که طبق تعریف شامل کل مشاهدات ممکنی است که برای اندازه‌گیری در شرایط خاص از طرف محقق قابل قبول است (8). به عنوان مثال، به صورت بالقوه برای انجام OSCE، 50 ایستگاه قابل طراحی است که در حال حاضر تنها 10 ایستگاه از بین آن‌ها انتخاب شده است. اگر نمونه‌های دیگری انتخاب شده بودند، احتمالاً نمرات دانشجویان

حوزه آموزش دارند. در مقابل، روش‌هایی که مبتنی بر یک بار اجرای آزمون هستند و تحت عنوان همسانی درونی<sup>12</sup> شناخته می‌شوند، مانند روش‌های دو نیمه کردن<sup>13</sup>، کودر-ریچاردسون<sup>14</sup> و آلفای کرونباخ برای سنجش پایایی ابزارهای ارزیابی فرآیند بیشتر استفاده می‌شوند (4). اگرچه نظریه کلاسیک به دلیل سادگی مفاهیم ابتدایی، قابلیت اجرای بالا و داشتن نرم‌افزارهای آماری متنوع، متدائل ترین نظریه اندازه‌گیری مورد استفاده است، محدودیت‌هایی نیز برای آن ذکر شده است. همان‌طور که قبل از ذکر شد، منابع متنوع خطا مانند خود سوالات، شرایط امتحان و آزمون‌گر با دور کردن نمرات مشاهده شده از نمرات واقعی دانشجویان، بر پایایی آزمون تأثیر می‌گذانند. نکته اینجاست که برای محاسبه پایایی آزمون در نظریه کلاسیک نمی‌توان در آن واحد تمام منابع خطا را در نظر گرفت و در هر لحظه تنها اثر یکی از این منابع قابل بررسی و برآورد است. در واقع نظریه تعمیم‌پذیری برای رفع این محدودیت ارایه شد و برتری اصلی آن نسبت به نظریه کلاسیک این است که سهم منابع مختلف خطا را در اندازه‌گیری مشخص می‌کند (1).

**نظریه تعمیم‌پذیری:** با توجه به محدودیت نظریه کلاسیک، نظریه تعمیم‌پذیری، توسط کرونباخ<sup>15</sup> معرفی شد و بعدتر توسط برنان<sup>16</sup> گسترش یافت (1 و 6). این نظریه که عموماً به عنوان جی-تئوری<sup>17</sup> شناخته می‌شود، در واقع بسط نظریه کلاسیک است و در شناخت منابع خطا بسیار کارایی دارد (7). در نظریه کلاسیک می‌توان تنها اثر یک منبع خطا را مشخص نمود اما در نظریه تعمیم‌پذیری، می‌توان منابع مختلف خطا را شناسایی نمود و مقدار اثر هر یک از آن‌ها را برآورد کرد تا تصویر واضح‌تری از خطای اندازه‌گیری به دست آید و تفسیر دقیق‌تری از نمرات قابل ارایه باشد. به این ترتیب متولیان آزمون قادر خواهند بود منابع اصلی خطا را شناسایی کنند و برای برگزاری آزمونی پایاتر برنامه‌ریزی کنند (1 و 8).

پایه محاسبات در نظریه تعمیم‌پذیری بر اساس مدل‌های آماری آنالیز واریانس<sup>18</sup> است که هرچند از طریق نرم‌افزارهایی مانند SPSS قابل اجرا است اما با توجه به دشواری‌های آن، مخصوصاً هنگامی که تعداد منابع خطا زیاد است، برنامه‌های اختصاصی نظریه تعمیم‌پذیری نیز تدوین شده است که از جمله

<sup>19</sup> Facet

<sup>20</sup> Object of Measurement

<sup>21</sup> Facet of differentiation

<sup>22</sup> Facet of generalization

<sup>23</sup> Universe Score

<sup>24</sup> Universe

<sup>25</sup> Population

<sup>26</sup> Conceptual

<sup>27</sup> Universe of admissible observations

<sup>12</sup> Internal consistency

<sup>13</sup> Split-halves

<sup>14</sup> Kuder-Richardson

<sup>15</sup> Cronbach

<sup>16</sup> Brennan

<sup>17</sup> G theory

<sup>18</sup> ANOVA

یک معیار مشخص می‌شود)، فرمول ضریب مطلق استفاده می‌شود اما اگر قرار است از نتایج آزمون به صورت هنجارمحور (مانند آزمون پذیرش که ردی و قبولی دانشجو بر اساس رتبه وی تعیین می‌گردد) استفاده شود، از فرمول ضریب نسبی استفاده می‌شود.

محاسبه ضریب تعمیم‌پذیری در حالت تکرویه: در عالم واقع آن چه در یک آزمون به عنوان منابع خطا دخیل است، محدود به یک رویه نیست. اما می‌توان در نظر گرفت که در یک آزمون فرضی تنها نقش یک رویه برای ما مهم است و می‌خواهیم تأثیر آن را بستجیم. به عنوان مثال، یک ایستگاه با سه آزمونگر برای 10 دانشجو برگزار می‌شود. آزمونگران با مشاهده عملکرد فرآگیران نمره‌ای بین 1 تا 10 به آن‌ها اختصاص می‌دهند (جدول 1).

جدول 1- توزیع نمرات ده دانشجو در یک ایستگاه OSCE با سه آزمونگر

شماره دانشجو	آزمونگر سه	آزمونگر دو	آزمونگر یک
8	7	6	1
6	5	4	2
2	2	2	3
5	4	3	4
6	4	5	5
10	9	8	6
9	7	5	7
8	7	6	8
8	6	4	9
8	9	7	10

برای برآورد میزان خطای این آزمون تنها یک رویه یعنی آزمونگر را در نظر می‌گیریم. اگر برای داده‌های فوق ANOVA را به صورت معمول انجام دهیم، دو جدول برای منابع تغییرات بین گروه<sup>35</sup> (دستیاران) و درون گروه<sup>36</sup> (آزمونگران و خطا) خواهد داد که خلاصه داده‌های آن‌ها در جدول 2 آمده است. باید از جدول‌ها مقادیر میانگین مربوطات را انتخاب کنیم و سپس به واریانس تبدیل نماییم یا مستقیماً خروجی اجزای واریانس را دریافت کنیم، با جایگزین کردن مقادیر در فرمول ضریب تعمیم‌پذیری، مقدار آن به صورت زیر به دست می‌آید:

$$G = \frac{\sigma_{student}^2}{\sigma_{student}^2 + \sigma_{error}^2}$$

$$G = \frac{4.037}{4.037 + 0.556} = 0.88$$

همان گونه که مشخص است ضریب تعمیم‌پذیری مقدار خوبی دارد و سهم آزمونگران در ایجاد خطا ۱۷٪ واریانس کل بوده است که زیاد نیست. یعنی آزمونگران توانسته‌اند ارزیابی خود را با

<sup>35</sup> Between-Subjects

<sup>36</sup> Within-Subjects

تابستان ۹۴، دوره هجدهم، شماره دوم، پایاپی 69

متفاوت بود. بنابراین، خطای در آزمون وجود دارد که ناشی از ایستگاه است و مقدار آن باید برآورده شود و مشخص شود نمره‌ای که دانشجو در این آزمون ۱۰ ایستگاهی کسب کرده است، تا چه حد قابل تعمیم به تمام ۵۰ ایستگاه فرضی است. آن حالت بالقوه، «جهان تعمیم»<sup>۲۸</sup> نام دارد و به حالت عملی شده، جهان مشاهدات قابل قبول گفته می‌شود. میانگین نمره دانشجو از تمام حالات جهان تعمیم همان نمره جهانی وی است. در برخی از موارد این دو جهان با هم یکسان هستند. مثلاً هنگامی که در یک آزمون از پنج آزمونگر استفاده می‌شود و کل تعداد استادان گروه مربوطه هم پنج نفر است. در اینجا آن چه ممکن و عملی شده بر آن چه تعمیم می‌دهیم، منطبق است. تعیین این موضوع در استفاده از آنالیز مناسب اهمیت دارد. به این ترتیب که اگر تعداد رویه مورد نظر در جهان ممکن و تعمیم منطبق بر هم باشند (مثلاً همان پنج آزمونگر)، از مدل «اثرات ثابت»<sup>۲۹</sup> استفاده می‌شود و اگر مانند مثال ایستگاه، تعداد رویه در جهان تعمیم بیشتر باشد، از مدل «اثرات تصادفی»<sup>۳۰</sup> استفاده می‌گردد (8).

ضریب تعمیم‌پذیری: ضریب تعمیم‌پذیری<sup>۳۱</sup> مقداری بین صفر تا یک دارد و معادل ضریب پایابی آزمون در نظر گرفته می‌شود. این ضریب مشخص می‌کند که نمره دانشجو را تا چه حد می‌توان به تمام رویه‌ها تعمیم داد. به عبارت دیگر مشخص می‌کند که نمره دانشجو تا چه حد به نمره واقعی او نزدیک است. نحوه محاسبه ضریب تعمیم‌پذیری به این ترتیب است که ابتدا باید رویه‌های موردنظر را تعیین کرد. هر رویه واحد مقداری به نام «جزء واریانس»<sup>۳۲</sup> است که از طریق ANOVA محاسبه می‌شود. از اجزاء واریانس رویه‌های مختلف در نهایت طبق فرمول مشخصی ضریب تعمیم‌پذیری به دست می‌آید. فرمول‌های متنوعی برای محاسبه ضریب تعمیم‌پذیری یک آزمون وجود دارد که بسته به این که کدام منابع خطا در نظر می‌گیریم و چگونه نمرات را تفسیر می‌کنیم، انتخاب می‌شوند. به عنوان مثال، از لحاظ مفهومی می‌توان دو نوع ضریب تعمیم‌پذیری مطلق<sup>۳۳</sup> و نسبی<sup>۳۴</sup> برای آزمون در نظر گرفت که به هدف برگزاری آزمون بر می‌گردد (7 و 8). اگر قرار است از نمرات برای تصمیم‌گیری به صورت معیارمحور استفاده شود (مانند آزمون پایان ترم که ردی و قبولی هر دانشجو نسبت به

<sup>28</sup> Universe of generalization

<sup>29</sup> Fixed effects

<sup>30</sup> Random effects

<sup>31</sup> G coefficient

<sup>32</sup> Variance component

<sup>33</sup> Absolute

<sup>34</sup> Relative

## مرواری بر نظریه تعمیم‌پذیری و...

مطلق از مقدار عددی ضریب تعمیم‌پذیری نسبی کمتر باشد. واریانس رویه‌ها در حالت هنجاری در جدول ۴ آمده‌اند. ضریب تعمیم‌پذیری نسبی که با  $\Phi$  (فی) هم نشان داده می‌شود، با فرمول زیر به دست می‌آید:

$$G_{\text{relative}} = \Phi = \frac{\sigma_s^2}{\sigma_s^2 + \sigma_t^2 + \sigma_{sxr}^2 + \sigma_{\text{residual}}^2}$$

$$G_{\text{relative}} = \Phi = \frac{60}{60 + 15 + 5 + 20} = \frac{60}{100} = 0.60$$

مطالعه تعمیم‌پذیری و مطالعه تصمیم‌گیری: تخمین سهم منابع مختلف خطا در آزمون و تعیین میزان تعمیم‌پذیری نتایج آزمون که قبلاً شرح داده شد، به «مطالعه تعمیم‌پذیری»<sup>37</sup> معروف است. در عین حال، از دیگر قابلیت‌های این نظریه، امکان پیش‌بینی پایایی آزمون در شرایط مختلف فرضی است که «مطالعه تصمیم‌گیری»<sup>38</sup> نامیده می‌شود (8). در واقع تلاش می‌شود با استفاده از یافته‌های مطالعه تعمیم‌پذیری، طرح‌های دیگری برای آزمون پیشنهاد شود (مثالاً افزایش تعداد ایستگاه‌ها یا آزمونگران) که در آن‌ها مقدار خطای آزمون به حداقل رسیده و تعمیم‌پذیری آزمون حداکثری باشد.

مزایا و محدودیت‌های نظریه تعمیم‌پذیری: همان‌طور که قبلاً اشاره شد، مزایای نظریه تعمیم‌پذیری در طراحی و تحلیل آزمون‌های مختلف شامل مشخص کردن سهم منابع مختلف خطا، برآورد پایایی و ارایه پیشنهاد برای بهبود تعمیم‌پذیری نتایج است. از محدودیت‌های نظریه تعمیم‌پذیری می‌توان به اجرایی نبودن استفاده از نظریه تعمیم‌پذیری در شرایط واقعی اشاره کرد. به دلیل محدودیت‌های عملی، معمولاً از این نظریه در مطالعات پایلوت (تعداد محدود ایستگاه‌ها و ارزیابان) استفاده می‌شود و سپس نتایج حاصل به شرایط واقعی ارزیابی تعمیم داده می‌شوند (10) و (11). همچنین پیچیدگی و کمبود نرم‌افزارهای تحلیل از محدودیت‌های دیگر نظریه تعمیم‌پذیری است (12).

خطای اندکی انجام دهنده نمرات قابل اطمینان است. سهم بزرگی از واریانسی که در نمرات وجود دارد، یعنی بیش از 70 آن، مربوط به عملکرد فراگیران و قابل قبول است.

محاسبه ضریب تعمیم‌پذیری در حالت دو رویه: اکنون یک OSCE را تصور کنید که با 3 ایستگاه و 2 آزمونگر در هر ایستگاه برگزار شده است. برای سهولت محاسبات، مدل اثرات ثابت را در نظر می‌گیریم. به عبارت دیگر فرض بر این است که جهان تعمیم و جهان مشاهدات ممکن یکی هستند. در واقع همین سه ایستگاه و همین دو آزمونگر موجود بودند. یا این که هدف ما این نیست که نتایج آزمون را به setting دیگری تعمیم دهیم. واریانس رویه‌ها در این حالت در جدول 3 آمده است. طبق فرمول زیر، ضریب تعمیم‌پذیری آزمون دست می‌آید (ایستگاه با  $t$  نشان داده شده است) یعنی واریانس تمام منابع به همراه واریانس خطای مخرج کسر قرار می‌گیرند:

$$G = \frac{\sigma_s^2}{\sigma_s^2 + \sigma_t^2 + \sigma_{sxr}^2 + \sigma_{\text{residual}}^2}$$

$$G = \frac{60}{60 + 20 + 15 + 15 + 5 + 15 + 20} = \frac{60}{150} = 0.40$$

قبلاً اشاره شد که در نظر گرفتن هدف آزمون، یعنی هنجاری یا معیاری بودن آن، در محاسبه ضریب تعمیم‌پذیری مهم است. فرمول فوق مربوط به ضریب مطلق است. از آنجا که در حالت معیارمحور، باید نمره دقیق تک تک دانشجویان تعیین شود، در تعیین میزان تعمیم‌پذیری نمرات نیز تمام خطای هر دانشجو نسبت شوند. در ضریب نسبی مشخص کردن جایگاه هر دانشجو نسبت به بقیه کفایت می‌کند و تعیین نمره وی مدنظر نیست. در واقع، تفاوتی که آزمونگران ایجاد می‌کنند، تأثیری در رتبه دانشجویان نسبت به یکدیگر ندارد. به همین دلیل منابع خطای کمتری وارد فرمول می‌شوند یعنی واریانس آزمونگر، واریانس تعامل آزمونگر با دانشجو و واریانس تعامل آزمونگر با ایستگاه حذف می‌شوند. این موضوع باعث می‌شود مقدار عددی ضریب تعمیم‌پذیری

جدول 2- نتیجه آزمون ANOVA درون گروه و بین گروه و اجزای واریانس مربوطه

	درصد واریانس از کل	میانگین مربعات	مقدار واریانس	درصد واریانس از کل	مجموع مربعات	درجه آزادی	مجموع خطای
دانشجو	73	4/037	12/67	9	114		
آزمونگر	17	0/944	10	2	20		
خطای باقی‌مانده (دستیار × آزمونگر)	10	0/556	0/56	18	10		

جدول 3- واریانس رویه‌های یک آزمون OSCE در حالت معیاری (مطلق)

	درجه آزادی	مجموع مربعات	میانگین مربعات	مقدار واریانس	درصد واریانس از کل	مجموع خطای
دانشجو (10 نفر)	9	3915	435	60	40	
آزمونگر (2 نفر)	1	815	815	20	13	
ایستگاه (3 ایستگاه)	2	960	480	15	10	
دانشجو × آزمونگر	9	585	65	15	10	
دانشجو × ایستگاه	18	540	30	5	3	
ایستگاه × آزمونگر	2	340	170	15	10	
خطای باقی‌مانده (دانشجو × آزمونگر × ایستگاه)	18	360	20	20	13	

جدول 4- واریانس رویه‌های یک آزمون OSCE در حالت هنجاری (نسبی)

					منع خطا
					دانشجو (10 نفر) ایستگاه (3 ایستگاه) آزمونگر : ایستگاه دانشجو × ایستگاه خطا یا باقی مانده (دانشجو × آزمونگر : ایستگاه)
52	60	435	9	3915	
13	15	480	2	960	
17	20	815	3	1155	
4	5	30	18	540	
13	15	170	27	945	

### نتیجه‌گیری

تفکیک امکان‌پذیر نیست در حالی که با استفاده از نظریه تعمیم‌پذیری می‌توان سهم هر یک از منابع خطای شناسایی نمود و به این ترتیب ضمن برآورد ضریب پایایی امتحان برای برگزاری آزمونی پایاتر برنامه‌ریزی کرد.

<sup>37</sup> Generalizability study (G study)

<sup>38</sup> Decision study (D study)

### References

- Brennan R. *Generalizability theory*. New York: Springer Verlag; 2001
- Raykov T, Marcoulides GA. *Introduction to Psychometric Theory*. 1st ed. New York: Routledge; 2010
- Khan KZ, Ramachandran S, Gaunt K, Pushkar P. The Objective Structured Clinical Examination (OSCE): AMEE Guide No. 81. Part I: an historical and theoretical perspective. *Medical Teacher* 2013;35(9):e1437-46.
- Tavakol M, Dennick R. Post-examination interpretation of objective test data: monitoring and improving the quality of high-stakes examinations: AMEE Guide No. 66. *Medical Teacher* 2012;34(3):e161-75.
- Courville TG. An empirical comparison of item response theory and classical test theory item/person statistics. Texas, US: Texas A&M University; 2004.
- Cronbach L, Gleser GC, Harinder N, Nageswari R. The dependability of behavioral measurement: Theory of generalizability for scores and profiles. New York: Wiley; 1972.
- Bloch R, Norman G. Generalizability theory for the perplexed: a practical introduction and guide: AMEE Guide No. 68. *Medical Teacher* 2012;34(11):960-92.
- Streiner DL, Norman GR. *Health Measurement Scales: A practical guide to their development and use*. 4th ed. Oxford University Press: USA; 2008
- Cardinet J, Johnson S, Pini G. Applying Generalizability Theory using EduG (Quantitative Methodology Series). 1st ed. Routledge; 2011
- Lawson DM. Applying generalizability theory to high-stakes objective structured clinical examinations in a naturalistic environment. *J Manipulative Physiological Therapy* 2006;29(6):463-7.
- Clauser BE, Harik P, Margolis MJ, Mee J, Swygert K, Rebbecki T. The generalizability of documentation scores from the USMLE Step 2 Clinical Skills examination. *Acad Med* 2008;83(10 Suppl):S41-4.
- Webb NM, Shavelson RL. Generalizability Theory: Overview. *Encyclopedia of Statistics in Behavioral Science* 2005; 2:717-719.

## An Overview of the Generalizability Theory and its Implications in Medical Education

Mortaz Hejri S<sup>1</sup> (MD, MSc, PhD candidate), Janani L<sup>2</sup> (PhD), Jalili M<sup>3\*</sup> (MD)

<sup>1</sup> Department of Medical Education, Tehran University of Medical Sciences, Tehran, Iran

<sup>2</sup> Department of Biostatistics, School of Public Health, Iran University of Medical Sciences, Tehran, Iran

<sup>3</sup> Department of Medical Education, Department of Emergency Medicine,  
Tehran University of Medical Sciences, Tehran, Iran

Original Article

Received: 6 Apr 2015, Accepted: 21 Jun 2015

### Abstract

**Introduction:** All measurements, including medical students' assessment, are potentially prone to error from various sources. Considering the importance of maintaining the quality of assessment methods in medical education, finding methods to estimate the error and to decrease it, is essential. Many developments have been achieved in measurement theories during recent years, which can be used in this way.

**Methods:** In this study, following reviewing the classical test theory (CTT) and its limitations, we discuss the generalizability theory (GT) using some examples of the Objective Structured Clinical Examination.

**Results:** According to the CTT, students' scores in exams are not their true scores and are accompanied with some levels of error. Determining reliability of tests helps us to assess how much the random error has been affected the score. In the CTT, it is impossible to consider all error sources simultaneously. To cope with this limitation, the GT specifies share of various sources of error based on the ANOVA models. This theory estimates reliability of tests and offers some recommendations for improving their generalizability.

**Conclusion:** In CTT, specification of effect of various sources of errors is not possible, whereas GT can be used to identify share of each error source. Hence, the test reliability coefficient can be estimated and also a more reliable test may be planned.

**Key words:** measurement, classical test theory, generalizability theory, reliability

### Please cite this article as follows:

Mortaz Hejri S, Janani L, Jalili M. An Overview of the Generalizability Theory and its Implications in Medical Education. Hakim Health Sys Res 2015; 18(2): 146- 152.

\*Corresponding Author: No. 57, Hojatdust St., Keshavarz Blvd., Tehran, Iran. Tel: +98- 21- 88955846, Fax: +98- 21- 889550160. E-mail: [mjalili@tums.ac.ir](mailto:mjalili@tums.ac.ir)