

## استخدام نظرية إمكانية التعميم فى تقدير ثبات بطاقة ملاحظة لتقييم أداء الطالب المعلم

جوزيف سمير وهيب كامل

مدرس مساعد بقسم علم النفس التربوى

كلية التربية - جامعة السويس

تاريخ استلام البحث : ١ / ١٢ / ٢٠٢٢ م

تاريخ قبول البحث : ١٩ / ١٢ / ٢٠٢٢ م

البريد الالكتروني للباحث : [joseph.samir@gmail.com](mailto:joseph.samir@gmail.com)

DOI: JFTP-2303-1266

## المخلص

إمكانية التعميم، ومعرفة مكونات التباين الأكثر تأثيراً على ثبات درجات تقييم الأداء بواسطة تلك البطاقة، ومعرفة أفضل الطرق لتحسين إجراءات القياس المستقبلية وأكثرها فعالية في تصميمات القياس للوصول لدرجات قابلة للتعميم، وقد طبقت الدراسة الحالية على (٢٧١) طالب وطالبة من طلاب الفرقة الرابعة بكلية التربية- جامعة السويس حيث طبقت عليهم بطاقة ملاحظة، وكشفت النتائج وصول بطاقة الملاحظة لدرجة جيدة إلى مقبولة من معاملات إمكانية التعميم النسبية والمطلقة مما يدل على ثبات درجات بطاقة الملاحظة، وإن أكبر مكونات التباين تأثيراً على معامل إمكانية التعميم النسبية والمطلقة تمثلت في تباين الأبعاد، وتفاعل طالب- بعد، وتفاعل طالب- بعد- فترة، وتفاعل طالب- بعد- مصحح- فترة الممزوج بالأخطاء العشوائية الغير مقدرة في التصميم، وجاءت باقى مكونات التباين الأخرى ضعيفة ولا تؤثر على معاملات إمكانية التعميم، كما أوضحت النتائج أن زيادة عدد الأبعاد وزيادة عدد مرات التطبيق تؤدي إلى زيادة معاملات إمكانية التعميم النسبية والمطلقة أفضل من زيادة عدد المصححين.

## الكلمات المفتاحية:

نظرية إمكانية التعميم- تقييم أداء - بطاقة الملاحظة.

## Using generalizability theory to estimate the stability of an observation card to Assesse the Student Teacher's performance

### ABSTRACT

The current study aims to estimate the stability of scores of an observation card to sssesse the student teacher's performance by using generalizability theory, and to know the components of the variance that most affect the stability of scores of performance assessment by observation card, and to know the best ways to improve future measurement procedures and the most effective in measurement designs to reach generalizable scores, The current study was applied to (271) male and female students from the fourth year at the Faculty of Education - Suez University, where an observation card was applied to them. Components of variance influencing the coefficient of relative and absolute generalizability were represented in the dimensional variance, the student-dimension interaction, the student-dimension-period interaction, and the student-dimension-corrector-period interaction mixed with unestimated random errors in the design, and the rest of the other variance components were weak and did not affect on the generalizability coefficients, and the results showed that increasing the number of dimensions and increasing the number of times of application to increase in the coefficients of relative and absolute generalizability better than increase the number of correctors.

### KEYWORDS:

generalizability theory - performance assessment - observation card.

**المقدمة:**

نتيجة للانتقادات المتعددة التي أبرزتها البحوث والدراسات التربوية في السنوات الأخيرة لاستخدامات التقييمات التقليدية للاختبار؛ فقد تزايد الاهتمام بالتقييمات البديلة القائمة على الأداء والتي جاءت كرد فعل للتقييمات التقليدية التي تعتمد على اختيار الإجابة، والتدريس من أجل الاختبار، وقياس معارف ومهارات منزلة، إضافة إلى تركيزها على قياس العمليات المعرفية الدنيا، واختلاف النواتج التعليمية عن المرغوبة من قبل المعلم، والتقييم المباشر المبني على المعارف دون الأداء؛ لذا أنتقلت بعدها الممارسات التربوية نحو استخدام تقييمات الأداء في مجالات التكوين والعمل مما يسمح بتعزيز المعارف والمهارات والاتجاهات ذات القيمة في التربية، وإحداث علاقة أفضل بين التعلم والتقييم (صلاح الدين علام، ٢٠٠٠؛ ٢٠٠٤).

وحيث أن المعلم العامل الرئيس في أي نظام تعليمي والقوة الفاعلة في المنظومة التعليمية، ويعد عنصراً فعالاً في تحقيق أهداف التربية الأمر الذي جعل الاهتمام به مدخلاً من المداخل الأساسية لإصلاح التعليم، لذا فإن تقييم أداء المعلم عملية قد تتجدد معاييرها التي يجب أن تستند عليها حتى تكون عملية وموضوعية وتفيد المعلم نفسه، وتفيد الأطراف التربوية الأخرى ذات الصلة بتدريبه وتأهيله؛ لذا يجب أن تواكب ما يستحدث من تطبيقات تقنية حتى يستطيع المعلم أن يؤدي أداءً تدريسياً فعالاً ومبدعاً تظهر نتائجه في مخرجات التعلم (Kerry, 2015).

ولقد استخدمت النظرية الكلاسيكية للقياس لفترة طويلة من الزمن في عملية بناء الاختبارات وتطبيقها وتحليل وتفسير نتائجها، إلا أنها لم تتمكن من معالجة بعض القضايا في القياس والتقييم التربوي، وقد وجهت لهذة النظرية الكثير من الانتقادات مثل معالجتها لمفهومي الصدق والثبات نظراً لأنها لا تميز بين أخطاء القياس المتعددة، وإنما تعطي قيمة تقديرية كلية لمصادر أخطاء القياس المتعددة المتعلقة بإحدى صور الاختبار المستخدمة بواسطة فاحص معين، وفي ظروف معينة. لذلك تتعدد أساليب تقدير الثبات، ونتيجة لذلك أوضح (فؤاد ابو حطب، ٣٢١، ١٩٩٢) إن علماء القياس النفسى والتربوى قاموا بتطوير وتعديل نظريات القياس حتى تتلائم مع مختلف أدوات القياس التي تُستخدم في الحكم والتقييم؛ للوصول إلى دقة تقدير الصدق والثبات لتلك الأدوات، وتحقيق موضوعية القياس.

وتعد نظرية إمكانية التعميم من النظريات المعاصرة المستخدمة لتقدير صدق وثبات تقييمات الاداء نتيجة لما فرضته تقييمات الأداء من تعدد أبعاد القياس والواقعية وإمكانية التعميم ووجود العديد من المصادر المحتملة المختلفة لأخطاء القياس لذا وجب الرجوع الي طرق إحصائية أكثر مرونة وتلائماً ودقة لمعالجة تلك الأدوات بدلاً من النظرية الكلاسيكية التي أصبحت غير كافية ومجدية. (Brennan, 2010).

**مشكلة البحث:**

بالرغم من حدود النظرية الكلاسيكية للاختبارات في معالجتها لأخطاء القياس إلا إنها ساهمت بشكل مباشر في ظهور نظرية إمكانية التعميم، وقدمت إطاراً فكرياً ساعدت على تطوير نظرية إمكانية التعميم بواسطة استخدام طرق تحليل التباين في تقدير ثبات القياسات السلوكية. لذلك لا يمكن اعتبار

نظرية إمكانية التعميم كبديل للنظرية الكلاسيكية ولكن يمكن اعتبارها كإمتداد لها لتلافي عيوبها، فهي عبارة عن مفهوم يجمع بين طرق تحليل التباين ونظرية القياس الكلاسيكية بأسم إمكانية التعميم. حيث أعتبر برينان (Brennan, 2001) أن النظرية الكلاسيكية وتحليل التباين كأبوين لنظرية إمكانية التعميم حيث قدموا حلول متعددة عن كيفية تجزئة مصادر التباين المتعددة في القياس وتقدير مكونات تباين الدرجات وتفاعلاتها المتعددة في نفس الوقت وبالتالي التحرر من فكرة خطأ أحادي البعد غير مميز.

وتؤكد نظرية إمكانية التعميم على أن هناك مصادر متعددة للخطأ في موقف القياس منها مثلاً خصائص المفحوصين، وخصائص الفاحصين، ووقت تطبيق الاختبار، ونوع فقرات الاختبار و... غيرها، وإن هذه الأخطاء قد تحدث معاً بصورة متزامنة أثناء القياس أو منفردة. فنظرية إمكانية التعميم ترى أن الدرجة الملاحظة يمكن تقسيمها الي درجة شاملة و واحد أو أكثر من مكونات الخطأ (Shavelson & Webb, 1991) لذا فلن يكون من المناسب أن يعتمد الباحث علي درجة واحدة أخذت من موقف واحد ووفق ظروف دون ظروف أخرى، وأعتمدت على مصدر (بُعد) واحد للخطأ دون بقية الابعاد مثلما الحال في النظرية الكلاسيكية (Eason, 1989) في: صبرى محمود، (٢٠١٦).

وعلى الرغم من محاولة نظرية الاستجابة للمفردة من معالجة أوجه القصور في النظرية الكلاسيكية ومحاولة كلاً من النظريتين أن تعالج تباين أخطاء القياس معالجة مختلفة، إلا أن كليهما ينظران إلى مصدر هذة الاخطاء نظرة ضيقة ومحدودة (صلاح الدين علام، ٢٠٠٠، ٦٩٩).

ومن الناحية التطبيقية لا يوجد أى نموذج أو نظرية متكاملة حتى الآن غير نظرية إمكانية التعميم يمكنها أن تعالج في آن واحد مختلف مصادر الخطأ في موقف القياس، وتوفير معلومات دقيقة عن كل مصدر من مصادر التباين الاخرى وتفاعلاتها مع بعضها البعض، إضافة إلى توفير طرق مختلفة لتعميم وتحسين إجراءات القياس لم تكن متوفرة في كل من النظرية التقليدية والاستجابة للمفردة (طباع فاروق، ٢٠٢٠).

كما أشار كاردينيت وآخرون (Cardinet et al., 2010) إن نظرية إمكانية التعميم لا تقدم وسائل لتقدير ثبات قياسات أنجزت من قبل فحسب، ولكن أيضاً تقدم معلومات حول مساهمة الاخطاء لإستخدامها في تحسين إجراءات القياس في المستقبل. لذلك تقوم دراسات القرار على إستخدام المعلومات المتحصل عليها من دراسات إمكانية التعميم لتصميم أفضل للقياس السلوكي.

كما أن التركيز في الماضى على المهارات الاساسية أدى إلى افتقار الطلبة إلى المهارات العقلية العليا، وإن الاختبارات التي تقتصر على الورقة والقلم والاختيار من متعدد فقط كان لها آثار سلبية كثيرة، حيث ركزت توجه المعلمين والطلبة نحو التدريب على هذة الاختبارات والاعداد لها، وبذلك عززت الممارسات التربوية الآلية وجعلت التعلم سطحي وسلبى مما أسهم في إعاقه الإصلاح الفعلى

للتعليم حيث تم التركيز على درجات الاختبارات بدلاً من تحسين تعلم الطلبة؛ مما أدى إلى إصلاحات تستند إلى تقويم بديل يركز على التقييم القائم على الأداء نظراً لتعدد وتنوع مستويات النواتج التعليمية (صلاح الدين علام، ٢٠٠٤).

فإستخدام مقاييس الأداء يسهم في عمليات ضبط وتوثيق نتائج القياس يزيد من الدقة في الاداء، وتحديد الفروق بشكل واضح، ويوفر الوقت والجهد والمال في عمليات التقويم، ويزيد من الموضوعية ويقلل من التحيز، لذا إكتسب أهمية في شتى العلوم فجميع العلوم تسعى لتطوير أساليب موضوعية دقيقة لقياس الظواهر المتعلقة بها من أجل فهم هذه الظواهر وتفسيرها، والتنبؤ بالعلاقات القائمة بين متغيراتها، ومحاولة ضبطها والتحكم بها (Cary, 1994).

ونظراً لشيوع إستخدام التقييمات التقليدية مثل الاختيار من متعدد والصواب والخطأ و... غيرها فإن تقدير الخصائص السيكومترية لها نالت اهتمام الباحثين لكثير من الوقت، وبنفس الطريقة إنتقل الاهتمام أيضاً نحو تقييمات الأداء خلال العقدين الماضيين لقياس جودتها الفنية خاصة ما يتعلق بصدق وثبات تلك التقييمات نظراً لاختلافها عن التقييمات التقليدية وتعدد أبعادها وإعطائها حرية أكبر للطلاب في الإجابة عن المهام التي تعكس إدماج معارف ومهارات متعددة، وواقعية سياقها، وإستدعائها لمهارات تفكير عليا؛ مما يتطلب إعادة النظر في مفاهيم الصدق والثبات الكلاسيكية (يسرى زكى، ٢٠١٦).

حيث أن استخدام النظرية الكلاسيكية لتقدير صدق وثبات تقييمات الأداء أصبحت غير كافية لوجود العديد من المصادر المحتملة لأخطاء القياس لذا وجب الرجوع إلي طرق إحصائية أكثر مرونة وتلائماً ودقة لمعالجة ذلك بأستخدام مفاهيم ومبادئ نظرية إمكانية التعميم ( Brennan, 2010; 2001; Suen & lei, 2007), حيث تبحث نظرية إمكانية التعميم في العلاقات بين أخطاء القياس المختلفة مع بعضها البعض بشكل متزامن وليست بشكل منفصل كما تفعل النظرية الكلاسيكية للقياس؛ فهي لاتستطيع أن تبحث أخطاء القياس الناتجة من المصححين والفترات وصور المهام و... غيرها في وقت واحد معاً كما تفعل نظرية إمكانية التعميم ( Shavelson& webb, 1991; Atilla&Ezel, 2010).

لذا عند فحص الدراسات في الأدب النفسى والسيكومتري، يلاحظ أن نظرية التعميم تُستخدم غالباً في تقدير ثبات تقييمات الأداء ( Lee& Kantor, 2007; Yilmaz&Gelbal, 2011; Yilmaz&Başbaşa, 2015), وقد قام هوانج (Huang, 2009) بجمع ودراسة نتائج أكثر من ٥٠ دراسة أشتملت علي ١٣٠ مجموعة بيانات مستقلة خلال الاعوام من ١٩٨٠ إلى ٢٠٠٦ تناولت تقييمات الاداء من خلال تطبيقات نظرية إمكانية التعميم في مجال تقييمات الاداء، وقام بالتوصل إلى نتائج كمية في عدة أبعاد من حيث طريقة التقييم (إختبار ورقة وقلم، وأنواع أخرى)، ومجالات

الموضوع (رياضيات، والعلوم، والاصغاء، والقراءة، والكتابة، واللغة الأجنبية، والاداء العسكري، والاداء الطبى).

وبالرغم من اعتماد أغلب تقييمات الاداء بمختلف أنواعها وتخصصتها علي نظرية إمكانية التعميم، فقد أعتدت بعض الدراسات العربية والاجنبية الأخرى مثل دراسة كلاً من ( Demet & Erkut, 2018)، (سامية مداح، ٢٠١٧)، (ياسر عبدالحافظ وعبدالله سالم، ٢٠١٣)، ( Maria & Maria, 2013)، (Mohamed, 2010) بالرغم من حداثهم في تقييم أداء المعلم - في مختلف تخصصاتهم- على استخدام نظرية القياس التقليدية في القياس بالرغم من عدم فاعليتها في مقاييس تقييم الأداء دون اللجوء الي نظرية إمكانية التعميم رغم قدرتها وأفضليتها في معالجة تقييم الأداء، كما أن أغلب تلك الدراسات تقوم على إستبيانات تقليدية لمجموعة من المهام والكفايات للمعلم دون الاعتماد على تحليل عمل مهنة المعلم.

وأستخدم كاري (Carrie, 2013) نظرية إمكانية التعميم لتحديد مصادر التباين في أداة الملاحظة التجريبية المصممة لتقييم فعالية معلم التربية الخاصة، وقد خلصت النتائج إلى ضرورة وجود محكمين مختلفين مثل (مديري المدارس وأعضاء هيئة التدريس بالجامعة ... الخ)، وضرورة استخدام أحجام أكبر لعينات الواجهه(محكمين/مواقف) لزيادة دقة القياس للأداة المستخدمة.

كما أوضح كلاً من (Bell et al., 2012; Hill et al., 2012) أن النموذج الكامل للحكم علي جودة المعلم يتضمن مجموعة من الأبعاد/المكونات وهي: كل مستويات المفردات، المقيمين، مناسبات القياس عبر الايام (الوقت)، أيام التقييم، الاقسام، كما أوضح (Bakker et al., 2008) أن زيادة عدد المقيمين وزيادة مهام التقييم وتنوعها تؤدي إلى معاملات أكثر تعميم عند تقييم أداء المعلم.

كما أوضحت العديد من الدراسات التي تناولت تقييم الأداء بأستخدام نظرية إمكانية التعميم مثل دراسة كلاً من ( Shavelson et al., 1993; Lane et al., 1996; Webb et al., 2000; Gao& Brennan, 2001; Nie et al., 2007;Guler& Gelbal, 2010; Taylor& Pastor, 2013) وداسة كلاً من (صبرى محمود، ٢٠١٦)، ودراسة ( طباع فاروق، ٢٠٢٠) أن أكبر مصدر للتباين يعود إلى تفاعل الطالب مع الاختبار مقارنة بباقي مكونات التباين الأخرى التي جاءت ضعيفة، وهذا يتفق مع التحليل الذى أجراه هوانج (Huang, 2009) حول الدراسات التي أهتمت بتقييمات الاداء فى تخصصات مختلفة و أوضحت تلك الدراسات أن نسبة كبيرة من مصدر الخطأ يرجع إلى تفاعل الطالب مع المهمة أو الاختبار.

فى حين اظهرت دراسات أخرى ( Chen et al., 2007; Lee& Kantor, 2007; Gebril, 2009) بأن أكبر مصدر للتباين يرجع إلى تفاعل الطالب مع الاختبار مع المصحح الممزوج بالاختبار العشوائية غير المقدرة، بالإضافة إلى دراسة كازانوف وديماس (Casanova& Demeuse, 2011) التي كشفت أن اكبر مصدر للتباين راجع إلى المصححين، وتفاعل المصححين مع الاختبارات.

وفيما يتعلق بمستوى ثبات الاختبارات المستخدمة في تقييم أداء الطلاب، فقد تبين نتائج الدراسات حيث أظهرت بعض الدراسات ( Lane et al., 1996; Lee& Kantor, 2007; Guler& Gelbal, 2010) بلوغ درجات الاختبارات مستويات مقبولة من الثبات، إلا أن البعض الآخر من الدراسات مثل ( Shavelson et al., 1993; Nie et al., 2007; Taylor& Pastor, 2013) كشفت عن ضعف في ثبات الاختبارات.

ومن ناحية أخرى أثبتت معظم الدراسات ( Nie et al., 2007; Lee& Kantor, 2007; Guler& Gelbal, 2010; Taylor& Pastor, 2013)، وداسة كلا من (صبرى محمود، ٢٠١٦)، و( طباع فاروق، ٢٠٢٠) أن زيادة عدد الاختبارات تساهم أكثر في زيادة معاملات إمكانية التعميم مقارنة بزيادة عدد المصححين، وهو على عكس ما توصلت إليه دراسة كازانوف وديماس (Casanova& Demeuse, 2011) التي كشفت أن زيادة عدد المقدرين أفضل من زيادة عدد المهمات في زيادة معاملات إمكانية التعميم، كما أوضح هوانج (Huang, 2009) أنه يمكن استخدام طرق أخرى لخفض مكونات التباين تمثلت في إدماج أكبر عدد من الأبعاد في موقف القياس، وإستخدام تصميمات متقاطعة بدلاً من تصميمات متداخلة، وإدماج الفترة كبعد من أبعاد القياس.

إن الاختلاف في نتائج الدراسات السابقة حول أختلاف مصادر تباين الخطأ المؤثرة على ثبات درجات تقييم الطلاب بالإضافة إلى مدى تحقيق الاختبارات لمستويات مقبولة من الثبات، وقلة الدراسات في البيئة العربية والمصرية التي تتناولت نظرية إمكانية التعميم؛ يكون مبرراً لإجراء البحث الحالي؛ لذلك سوف تهتم الدراسة الحالية بالاجابة على مجموعة من التساؤلات المرتبطة والمنقسمة إلى تساؤلات ترتبط ارتباطاً وثيقاً بتصميمات نظرية إمكانية التعميم التي تهتم بتفسير وتقدير مكونات التباين المؤثرة على دقة القياسات، وتساؤلات ترتبط ارتباطاً وثيقاً مع دراسات القرار التي تهتم بتحسين إجراءات القياس للوصول إلى مستويات مقبولة من إمكانية التعميم في المستقبل. وبالتالي يمكن صياغة مشكلة البحث في الاسئلة التالية:

- ١- ماهو حجم التباين الذي تساهم به كل مصدر من مصادر التباين (الأبعاد) في التباين الكلي لمتوسط درجات أداء الطالب/المعلم على بطاقة الملاحظة ؟
- ٢- ما تأثير زيادة كل مصدر من مصادر التباين (الأبعاد) المؤثرة على ثبات درجات أداء الطالب/المعلم على بطاقة الملاحظة ؟

### أهداف البحث:

يهدف البحث الحالي إلى :

- ١- معرفة مكونات التباين الاكثر تأثيراً على ثبات درجات بطاقة الملاحظة لتقييم أداء الطالب المعلم.

٢- معرفة أفضل طرق تحسين إجراءات القياس وأكثرها فعالية في تصميمات القياس، عن طريق زيادة عدد الأبعاد لمصادر التباين للحصول على درجات قابلة للتعميم .

### أهمية البحث:

يكتسب البحث الحالي أهميته من العديد من المنطلقات منها :

١- المساهمة في توضيح وفهم مفاهيم ومبادئ نظرية إمكانية التعميم، وتقييمات الاداء، ودلالة وأهمية كل منهما من الناحية النظرية والتطبيقية في تقدير دقة القياسات وإجراءات تحسين عمليات القياس للكثير من الباحثين لعدم الالمام بها حتى الآن.

٢- المساهمة في تقييم الاداء المهني للمعلم باستخدام أداة مبنية علي تحليل عمل مهنة المعلم ونظرية إمكانية التعميم تكون بمثابة اختبار قبول للطلاب بكليات التربية لضمان الحد الأدنى من الأداء المطلوب للقبول بكليات التربية، ويمكن أن تستخدم كرخصة لمزاولة مهنة التدريس.

### مصطلحات البحث:

نظرية إمكانية التعميم (Generalizability Theory) (GT)

نظرية إحصائية لتقييم صدق وثبات القياسات السلوكية، وتركز النظرية على مصادر أخطاء القياس المتنوعة حيث تقوم بفصلها وتقديرها وترتيبها في نفس الوقت (Salkined, 2008, 436)، تعرفها (Urbina, 2004, 138) بأنها تقدير بديل للثبات، وامتداد للنظرية الكلاسيكية للقياس تستخدم طرق تحليل التباين في تقدير التأثيرات الكلية للمصادر المتعددة لتباين الخطأ في درجات الاختبار في آن واحد.

تقييم الاداء Performance Assessment

هو أن يظهر المتعلم بوضوح، أو يبرهن، أو يقدم أمثلة، أو تجارب أو نتائج أو... غير ذلك، تتخذ دليلاً على تحقيقه مستوى تربوياً أو هدفاً تعليمياً معيناً. وتتطلب مهام الاداء إجراء العمليات والتوصل إلى نتائج، وأهم ما يميز هذه المهام أنها مباشرة ووظيفية وواقعية أي تمثل مواقف حياتية فعلية خارج نطاق الصف المدرسي كما أنها لا تتطلب بالضرورة الورقة والقلم (صلاح الدين علام، ٢٠٠٤، ١٠٥).

الإطار النظري والدراسات المرتبطة:

أولاً: نظرية إمكانية التعميم

يري كرونباخ Cronbach وراجارتنام Rajaratnam وجليزر Glaser أن درجات الاختبار أو أي درجات ملاحظة تُعد عينة من نطاق شامل للدرجات الممكنة وتتباين وتختلف ظروف جمع تلك الملاحظات وفقاً لأبعاد متعددة، ومثال على ذلك تطبيق اختبار معين مرة واحدة، أو تطبيقه وإعادة تطبيقه واشتماله على عينة ممثلة من المفردات دون بقية العينات المستمدة من النطاق الشامل، وتطبيقه تحت ظروف معينة دون باقى الظروف، وهذا التباين في الأبعاد المتعلقة بأداة القياس وما تتضمنه من مفردات وظروف تطبيقها وخصائص مجموعة الأفراد التي يطبق عليها الاختبار تؤدي إلى

مشكلات في عمليات القياس وإمكانية تعميم نتائجها، وقد اقترح كرونباخ وزملائه عام ١٩٧٢م نظاماً لتقدير تأثير كل من هذه الأبعاد في تباين الدرجات الملاحظة وتقدير دقة الدرجة الحقيقية من الدرجات الملاحظة التي تتأثر بجميع أبعاد عملية القياس وإمكانية تعميم تلك الدرجات الملاحظة في ضوء هذه الأبعاد المتعددة عند القياس، وقد عرف ذلك النظام باسم نظرية إمكانية التعميم **Generalizability Theory** أو نظرية الثبات متعددة الأوجه **Multifaceted Reliability Theory** أو نظرية عينة المجال **Domain Sampling Theory** فهي نظرية تبرز أهميتها واسهاماتها في توسيع نطاق استخدامات كل من مفهومي الصدق والثبات وتكاملهما (صلاح الدين علام، ٢٠٠٠، ١٧١).

### أسباب ظهور نظرية إمكانية التعميم:

١- يري كرونباخ ومعاونيه أن ضعف النظرية الكلاسيكية تكمن في التبسيط الشديد لمفهوم الصدق والثبات مما جعل الباحثين يستخدمون هذين المفهومين استخداماً إلثاً، مما أدى إلى غموض معناهما الحقيقي واستخدامهما الفعلي؛ لذا فقد اهتم كرونباخ ومعاونيه بمراجعة وتعميق مفهومي الصدق والثبات ومعالجتهما معالجة أكثر شمولاً من منظور متعدد الأبعاد (Nicholas, 2016).

٢- لا تسمح الطرق المختلفة لتقدير الثبات بالنظرية الكلاسيكية مثل إعادة التطبيق، والصور المتكافئة، والاتساق الداخلي لصانع القرار من التحكم في أخطاء القياس الأخرى التي قد تؤثر على الدرجة الملاحظة أثناء استخدام إحدى هذه الطرق مما يؤدي إلى غموض في تفسير صدق وثبات درجات المقاييس في معناهما واستخداماتهما الفعلية؛ فتعدد الطرق المستخدمة في تقدير الثبات يعبر عن عدم قدرة النظرية الكلاسيكية على التحكم في الأخطاء المتعددة التي يمكن أن تؤثر في عملية القياس، كما لا تتميز طرق قياس الثبات بالمرونة الكافية لتناسب مع مشكلات الثبات التي تظهر في جميع الاختبارات العقلية (Crocker & Algina, 2006) كما تبين أيضاً أن تطوير النظرية الكلاسيكية لطرق متعددة في تقدير الثبات تُحدد الثبات بشكل مختلف لا يتماثل مع التعريف النظري له (Webb et al., 2006).

٣- على الرغم من أن النظرية الكلاسيكية ونظرية الاستجابة للمفردة تعالجان تباين أخطاء القياس معالجة مختلفة؛ إلا أن كليهما ينظر إلى مصدر هذه الأخطاء نظرة ضيقة ومحدودة ولا يمكنهما التمييز بين أخطاء القياس؛ فالمعلومات أو البيانات التي تستمد سواء من العلوم الطبيعية أو العلوم السلوكية تكون مشوبة بدرجة معينة من الخطأ، وقد تكون بعض أنواع أخطاء القياس غير متضمنة في عملية القياس وذلك يضع قيداً على نطاق ظروف القياس التي يمكن التعميم عليها؛ فنظرية إمكانية التعميم تعتمد على تحديد ظروف القياس التي تسهم في الخطأ لكي يمكن تعميم القياس عليها (صلاح الدين محمود، ٢٠٠٠، ٦٩٩).

### أهمية نظرية إمكانية التعميم

نظرية إمكانية التعميم تُعرف بنظرية متعددة الأبعاد فهي تزودنا بمعلومات مفصلة عن مختلف مصادر الخطأ في درجات الاختبار عن طريق جمع ملاحظات متعددة لنفس مجموعة الأفراد في كل

المتغيرات المستقلة التي تساهم في تباين الخطأ في القياس مثل المهام، والمصححين، والفترات الزمنية، وصيغ مهام الاختبار، وأيضاً التفاعلات بين الأبعاد وبعضها البعض (طباع فاروق، ٢٠٢٠). كما أن نظرية إمكانية التعميم لم تهتم فقط بكيفية بناء أدوات قياس تكون نتائجها متكافئة في ظروف ومواقف اختبارية مختلفة، وإنما اهتمت أيضاً باستخدامات أدوات القياس في صنع قرارات معينة، فنظرية إمكانية التعميم تميز تحليلات الثبات من منظورها إلى مرحلتين مختلفتين، المرحلة الأولى تسمى بدراسة إمكانية التعميم تهدف إلى تحديد المصادر أو الأبعاد المتعددة لأخطاء القياس وتقييم مقدار التباين في الدرجات المحصلة من كل مصدر أو بعد من أبعاد القياس، وتسمى المرحلة الثانية بدراسة القرار وتهدف إلى تحديد الثقة في درجات القياس وإمكانية تصميم إجراءات قياس أكثر فعالية (Meyer, 2010).

وتفترض نظرية إمكانية التعميم أن الدرجة الشاملة هي بديل للدرجة الحقيقية في نظرية القياس الكلاسيكية، وتقسّم خطأ القياس إلى مصادر خطأ متعددة حيث يعتبر التمييز بين مصادر الأخطاء عند القياس من المساهمات الأساسية لنظرية إمكانية التعميم (Bertrand & Blais, 2004). يرى شافلسون، وويب (Shavelson & Webb, 1991) أن قوة نظرية إمكانية التعميم تكمن في إمكانية تقدير مختلف مصادر القياس بشكل مستقل في تحليل واحد، إضافة إلى تمكين صانع القرار من تحديد المصادر الكامنة التي تسهم في تباين الدرجات للحصول على درجة عالية من ثبات القياس. المفاهيم الأساسية لنظرية إمكانية التعميم

إن ما تتطلبه تلك النظرية من تصورات فكرية وأساليب إحصائية مختلفة عن باقي نظريات القياس الأخرى تتطلب بعض المفاهيم الأساسية لفهماها، ومن المفاهيم الأساسية للنظرية:

#### (١) النطاق الشامل: Universe Domain

هو نطاق يتضمن شروط جمع الملاحظات أو الحصول على القياسات حتى نتمكن بعد ذلك من إمكانية تعميم هذه القياسات من موقف إلى آخر ومن عينة الأفراد إلى عينة أخرى؛ فهو الهيكل أو الإطار الذي يحدد خصائص كل بعد متضمن في عملية القياس (Meyer, 2010).

#### (٢) الدرجة الشاملة: Universe Score

هي القيمة المتوقعة لدرجات الفرد الملاحظة والتي يحصل عليها في مختلف المواقف التي تنتمي إلى النطاق الشامل المحدد من قبل الباحث. ولل فرد أكثر من درجة شاملة على حسب الأبعاد التي يحتوي عليها النطاق الشامل. كما أن الدرجة الشاملة هي درجة تقابل الدرجة الحقيقية في النظرية الكلاسيكية (صلاح الدين علام، ٢٠٠٠، ٧٠٢).

**٣) البعد/ الوجه/ مصادر الخطأ: Facet**

هو عبارة عن المصادر التي قد تؤثر على موقف القياس مثل الاختبارات والمصححين والاهداف والفترات وصيغ الاختبار و... غيرها التي يمكن أن تختلف من ملاحظة إلى أخرى وبذلك يمكن أن تؤثر على درجة الفرد الملاحظة (Brennan, 2010).

**٤) المستويات/ الشروط: Levels/ Conditions**

كل بعد من الأبعاد يُشكل من مجموعة من المستويات التي تعبر عن مختلف شروط أو ظروف القياس في هذا البعد، مثل عدد الأفراد الذين يطبق عليهم الاختبار، وعدد المقدرين الذين يقدر أداء الطلاب، ومختلف فترات أداء الطالب للمهام و... هكذا (Cardinet et al., 2010, 11)

**٥) تحليل التباين: Analysis of Variance**

هو أسلوب يستخدم لعزل تباين متغير واحد من تباين المتغيرات الأخرى بواسطة تجزئة التباين الكلي لمجموعة الملاحظات التي ترجع إلى عوامل معينة مثل الجنس، ومجموعات المعالجة الأخرى (Everitt & Skrondal, 2010).

**٥) دراسات إمكانية التعميم: Generalizability Study (G-Study)**

هي دراسات تهتم بتحديد درجة تكافؤ نتائج الاختبارات التي نحصل عليها في ظروف اختبارية مختلفة وتحت شروط متباينة؛ أي أنها دراسات تستخدم في جمع معلومات وبيانات تفيد في تقدير مكونات تباين القياسات التي نحصل عليها بطريقة ما من الطرق، وهذا يفيد في تطوير هذه القياسات من خلال محاولة تقليل أثر تلك التباينات على أخطاء القياس (صلاح الدين علام، ٢٠٠٠، ٧٠٤).

**٦) معامل إمكانية التعميم: Coefficient Of Generalizability**

هو النسبة بين تباين الدرجة الشاملة إلى تباين الدرجة الملاحظة، وبالتالي يمكن اعتبار معامل إمكانية التعميم بديلاً عن معامل الثبات الذي نحصل عليه باستخدام النظرية الكلاسيكية للقياس، وكل اختبار أو مقياس يكون له أكثر من معامل إمكانية تعميم على حسب الأوجه التي يتم أخذها بعين الاعتبار في دراسة إمكانية التعميم (صلاح الدين علام، ٢٠٠٠، ٧٠٤).

**٧) دراسات القرار: Decision Study (D-Study)**

هي دراسات تهدف الوصول إلى أفضل إجراء قياس يتمتع بالثبات في موقف معين اعتماداً على المعلومات المتوفرة من دراسات إمكانية التعميم، وبالتالي فإنها تجيب على التساؤل كيف يمكن الوصول بثبات الاختبار إلى أفضل مستوى، وكيفية تحسين إجراءات القياس بالاعتماد على تقارير دراسة إمكانية التعميم وتقديرات خطأ التعميم الخاص بالتصميم المستخدم في الدراسة (طباع فاروق، ٢٠١٩).

**الخصائص السيكومترية للاختبارات في إطار نظرية إمكانية التعميم**

يري كلا من كرونباخ Cronbach وجليسر Gleser أن الصدق والثبات مفهومان مترابطان، ويمكن أن يندرجا تحت أسم مقياس إمكانية التعميم، ويكمن الفرق الرئيسي بينهما في الأبعاد التي نود التعميم عليها (صلاح الدين علام، ٢٠٠٤) لذا تعد نظرية إمكانية التعميم من نظريات القياس المعاصرة التي ساهمت في إبراز التكامل بين مفهومي الصدق والثبات، ويمكن استغلال نتائج ثبات القياسات السلوكية في التأكد من صدق أنظمة التقييم خاصة إذا تعلق الأمر بصدق التكوين الفرضي، ويتم ذلك من خلال تحليل مكونات التباين التي تفسر تباين الأداء في الاختبارات مما يساعد على تحديد النطاق الشامل للشروط المقبولة.

كما أشار بيركر وآل (Berker & Al, 1993)، و موس (Moss, 1992) أن التمييز بين الصدق والثبات في نظرية إمكانية التعميم أصبح ضئيلاً، حيث يمكن الأخذ بالاعتبار في تقييم الأداء العلاقات بين المهام كأدلة لكلاً من الصدق والثبات، كما أوضح شيفلسون و ويب (Shavelson & Webb, 2009) أنه يمكن اعتبار بعض أبعاد موقف القياس أبعاداً للصدق في ضوء نظرية إمكانية التعميم، وركزا على مجالات المحتوي المتعددة وأنواع المهام (اختيار من متعدد، مفتوحة، يدوية،...) كأبعاد للصدق تكون ضمنية في تصميم دراسة إمكانية التعميم.

وقد أوضح أيضاً شافيلسون وآخرون (Shavelson et al. 1993) التكامل بين الصدق والثبات في إطار نظرية إمكانية التعميم في تقييم الأداء من خلال فحص تغير معاينة مختلف طرق تقديم المهام دليل على الصدق التقاربي للاختبار. لذا يمكن اعتبار مختلف إجراءات القياس التي تعني بقياس نفس السمة وتعطي نفس الصورة دليلاً على الصدق التقاربي، ويمكن اعتبار تعميم عينة مهام أو بنود الاختبار على مجال واسع دليلاً على صدق المحتوى، ويمكن اعتبار مقدار التغير داخل القياس إذا كانت الفترة تمتد إلى أسابيع أو أشهر أو أبعد من ذلك إلى سنة دليلاً على الصدق التنبؤي (Shavelson & Webb, 2009) فالصدق والثبات مفهومان مترابطان يختلفان فقط في الأبعاد التي يود صانع القرار التعميم عليها.

**ثانياً: تقييم الأداء**

يطلق الباحثين على هذا النوع من التقويم عدة أسماء، وهي جميعها تمثل مرادفات لمعنى واحد، ومن أمثلة ذلك التقويم الشامل Comprehensive Assessment، والتقويم الأصيل أو الواقعي أو الحقيقي Authentic Assessment، والتقويم المعتمد على الأداء Performance Assessment.

ويعرفه (صلاح الدين علام، ٢٠٠٤، ١٠٥) بأنها إظهار المتعلم بوضوح، أو ببرهن، أو يقدم أمثلة، أو تجارب، أو نتائج أو غير ذلك لتحقيق مستوى تربوياً أو هدفاً تعليمياً معيناً، ويمكن تقييم أداء الطالب داخل الصف الدراسي من خلال: المحادثة الشفوية، والتعبير التحريري، وإجراء التجارب العملية، وتكوين المجسمات، ورسم الخرائط، وغير ذلك. ويمكن في تقييم الأداء تقييم العمليات

المتضمنة في الأداء أثناء تنفيذه، ويمكن أيضاً تقييم النتائج النهائية وتقدير درجة جودتها استناداً إلى موازين تقدير، وأهم ما يميز هذه المهام أنها مباشرة، ووظيفية، وواقعية أي تماثل مواقف حياتية فعلية خارج نطاق الصف المدرسي.

ويعتبر أحد الأساليب المتنوعة أو الأدوات المستخدمة في التقويم التربوي البديل، ويعرف بأنه أحد المقاييس التي تتطلب من الطالب أداء شيء ما مثل حل مسألة رياضية أو كتابة مقال صحفي أو تصميم ملصق أو... غير، وتقدير أدائه أو ملاحظته بواسطة أدوات قياس معينة في ظروف عمل حقيقية أو ظروف عمل محاكية لها، ومن ثم الحكم على هذا الأداء وفق معايير محددة سلفاً (يسرى زكى، ٢٠١٦، ١٨١٥).

ويعرفه (عبد العزيز محمد، ٢٠١٩، ٥٣٦) على أنه توظيف للطالب معلوماته ومهاراته في مواقف حياتية حقيقية أو مواقف تحاكي المواقف الحقيقية، أو إجراء للبحوث المتعلقة بموضوعات المقرر الدراسي، أو أدائه لبعض الاختبارات العملية، أو قيادة بأنشطة أو عروض أو أداءات عملية يظهر من خلالها مدى إتقانه لما أكتسبه من مهارات.

### أسباب ظهور تقييم الأداء

سيطرت الاختبارات الموضوعية ذات الاختيار من متعدد لمدة زمنية كبيرة على الاختبارات وذلك بسبب سهولة إعدادها، وتطبيقها، وتغطيتها للمحتوي الدراسي، وسهولة تصحيحها، وتميزها بدرجة عالية من الموضوعية والصدق والثبات، إلا أنها غير قادرة على قياس مهارات التفكير العليا كالقدرة على حل المشكلات، والتحليل، والاستدلال، ومهارات التفكير؛ فهي تهتم بقياس المعارف الفعلية والاتجاهات السلوكية للطلاب في مجال معين (Haladyna, 2004) إلا إنها غير قادرة على إنتاج إجابات مبنية من الطلاب فهو يختار إجابة من بين مجموعة من الاجابات.

لذا فقد أشار لين وآخرون (Linn et al., 1991) إلى مجموعة من الانتقادات للاختبارات الموضوعية تتمثل في: التأكيد على المعارف الفعلية والمهارات المنفصلة، والسرعة في الاجابة، والاجابات القصيرة المصطنعة، وعدم القدرة على تغطية مختلف محتويات أو مجالات المعرفة، والتدريس من أجل الاختبار.

لذا فرضت تلك الانتقادات على المهتمين بالتربية بإيجاد بدائل جديدة أكثر فاعلية لتلافي عيوب التقييم التقليدي، ولتحقيق ذلك اتجهت الجهود إلى الاهتمام بالتقييمات البديلة المعتمدة على الأداء لأثارة عمليات التفكير العليا. كما لم تتوقف الجهود إلى إنجاز مهام مجردة من السياق وبعيدة عن الواقع، ففي أواخر سنوات الثمانينات زاد الطلب بكثرة لاستخدام التقييم الواقعي للاداء ( Keller et al., 2010).

كما يركز تقييم الأداء على قدرة تنفيذ مهام محددة بنجاح بدلاً من اختزال بعض المعلومات والمعارف المتفرقة، لذا يمكن تفسير التحصيل من خلاله بأنه أداء أو قدرة على إنجاز شيء ما،

ويتضمن الأداء توسيع نطاق المعرفة والفهم وتنمية مهارات متعددة ومتكاملة يصعب قياسها وفق الأساليب الكلاسيكية (صلاح الدين علام، ٢٠٠٠).

### مكونات وأقسام تقييم الأداء

يري (صلاح الدين علام، ٢٠٠٤)، و(عبدالله السعداوي، ٢٠١٠) أن تقييم الأداء يشتمل على مكونين رئيسيين هما مهام الأداء، وقواعد تقدير الأداء، وفيما يلي عرض تلك المكونات:

١) مهام الأداء: وتتطلب إجراء عمليات أو سلسلة من الأنشطة، أو أداء عمل معين بطريقة مناسبة، أو إبتكار وتكوين نتائج مركبة تحقق مستويات معينة من الجودة، وعلى الرغم من إتساع نطاق مهام الأداء وتنوع صيغها وأنماطها إلا أنه يمكن تصنيفها إلى:

أ) مهام محدودة: تتطلب أن يقوم الطالب بأداء محدد في زمن قصير، وتستخدم هذه المهام المحددة عادة لأغراض المساءلة، والتحقق من كفايات الطالب في مهارات معينة ذات أهمية، مثال على ذلك: كتابة مقال من خلال استدعاء عمليات التحليل والمقارنة والاستنتاج والتقييم أو غيرها من المهام.

ب) مهام موسعة: تتطلب أن يختار الطالب موضوعاً معيناً ويعطى وقتاً كافياً للتفكير فيه، وبحثه من جوانب متعددة لكي يظهر تمكنه من الموضوع وعمق فهمه له، مثال على ذلك: إجراء تجربة علمية تتعلق بالتمثيل الضوئي وتفسير نتائجها، أو مراجعة قصيدة شعر والتعليق عليها، أو غير ذلك.

ج) عرض الاعمال: هو توظيف الطالب لمهاراته المتنوعة لإبراز كفاءته في مجال معين فهو بمثابة تقويم ختامي لانجازات الطالب.

د) ملفات الاعمال (البورتفوليو): تشمل على عينة من الأعمال المتجمعة المتنوعة التي قام بها الطالب خلال مدة دراسية معينة، ويمكن أن تتضمن هذه الأعمال نتائج ابتكارها الطالب، أو أفضل الأعمال التي أنجزها خلال مدة معينة، أو بعض الأعمال التي لا تزال في مرحلة الإعداد لكي يبين تحسنة خلال المدة المحددة، وتقويمه لجوانب القوة والضعف في بعض هذه الأعمال.

هـ) مهام أخرى: تشمل على أي أنشطة أخرى قام بها ودونها الطالب، أو لاحظها المعلم وحكم على نوعيتها وجودتها.

كما يرى (عبد العزيز محمد، ٢٠١٩، ٥٥٠) أن الأقسام الرئيسية الخمسة لمهام الأداء أشكالاً متعددة تتطلب جميعها إنشاء أو تكوين إجابات، أو ابتكار نتائج تظهر قدرة الطالب على استخدام وتوظيف معرفة ومهاراته في مجال معين، ومن تلك الأشكال:

- الإجابات الحرة أو المستفيضة أو الكتابة: تتمثل في أسئلة المقال التي تتطلب أن يكون الطالب إجابته بنفسه بدلاً من أن يختار الإجابة الصحيحة من بين بدائل معطاه، وتسمح تلك الإجابات

بتوظيف مهارات عقلية عليا كالفهم المتعمق أو الاستدلال، والتحليل، والتفكير الناقد، وحل المشكلات.

- الاجابات المحدودة: تتمثل في أسئلة الاختيار من متعدد، وأسئلة الصواب والخطأ، وأسئلة المزوجة، وغيرها من الأسئلة بشرط أن تعكس مهام واقعية من الحياة العملية وليست أجزاء مقتضبة من المعلومات أو الحقائق والمطلوب من الطالب تذكرها.

- التعبير الشفوي: يعد من الأنشطة الأدائية التي ينبغي أن يكتسبها الطلاب، ومن المهارات التي تتعلق بالتواصل اللغوي ويتم تقويمها شفويًا: القراءة، والتعبير، وإلقاء الكلمات، والمناظرات، والمناقشات، وتظهر قدرة الفرد على التعبير والتلخيص وربط الأفكار.

- عروض الأعمال: تعد من المهام الأدائية المتسعة النطاق، وتتطلب إظهار الطالب كفاياته المتعددة في تصميم أعمال معينة، وتتعدد صيغ عروض الأعمال مثل المقال، أو الورقة البحثية، أو المقابلة، أو العروض الدرامية، أو غيرها.

- التجارب: تعد التجارب من المهام الأدائية الأساسية في تعلم الرياضيات، والعلوم بخاصة، وتتضمن هذه التجارب ممارسة الطالب الفعلية لعمليات التخطيط والتصميم وإجراء التجارب العملية في مواقع العمل بدلاً من أختبارات الورقة والقلم.

(٢) قواعد تقدير الأداء: هي عبارة عن محكات الحكم على جودة مهام الأداء أو نوعيتها، وترتبط ارتباطاً وثيقاً بمهام الأداء حيث أن كلا منهما يعتمد على الآخر ويكمله، وتساعد تلك القواعد على تصميم مهام يمكن تقدير مكوناتها أو نتائجها، ويمكن أن تنقسم إلى:

(أ) قواعد تقدير محدودة: هي عبارة عن قواعد توضح خصائص الأداء الجيد للمهام المحدودة وما يشتمل عليها من مهارات وكفايات، وهذه القواعد تكون نوعية بمعنى أقتصرها على مهام محددة.

(ب) قواعد تقدير عامة: هي عبارة عن قواعد توضح خصائص الأداء الجيد للمهام الموسعة وما يشتمل عليها من مهارات وكفايات، ويمكن أن تنطبق هذه القواعد على مهام متنوعة ومختلفة.

(ج) ملاحظات المعلم المنظمة: من خلال ملاحظة أداء الطلاب وسلوكهم، ويمكن أن تكون هذه الملاحظات موضوعية من خلال ربطها بمستويات المحتوى، وتحليل مكونات الأداء المراد ملاحظة في مواقف طبيعية، وإعداد قواعد أو محكات لتقدير هذه المكونات والحكم على نوعيتها.

### خصائص مهام الأداء

إن مهام الأداء تتميز بعدة خصائص تميزها عن التقويم التربوي التقليدي، ومن تلك الخصائص:

١- الارتباط بالمستويات أو النواتج التعليمية: فالتقويم التربوي البديل ومهام الأداء تهدف إلى تحقيق مستويات تربوية متميزة أو نواتج تعليمية محددة تتطلب إبراز مهارات الطالب المتنوعة التي تتميز بالواقعية، في إطار مناهج دراسية قائمة على التفكير وعمليات تدريس تؤكد على الانشطة

التعليمية وشمولها وإيجابية الطالب ومشاركة بفاعلية وتوجيهه ومساعدته على بناء معرفته من خلال تفاعله مع بيئته ومرورة بخبرات متنوعة بدلاً من عمليات التدريس القائمة على التلقين من جانب المعلم (محمد بن راشد، ٢٠٢١، ١٤٧).

٢- الوضوح: بمعنى وضوح المطلوب أداة أو تكوينية، وشروط الأداء من حيث حدود الوقت، والمصادر المتاحة، والاستعانة بالآخرين، ومحكات تقدير الأداء والحكم عليه، وصياغة كل تلك التعليمات في عبارات بسيطة لضمان فهم جميع الطلاب لها (عادل سرايا، ٢٠٠٥).

٣- واقعية المهام: يقوم بتقييم الأداء على حل مشكلات واقعية أو كتابة أسلوب الحل الذي يراه مناسباً أو عرض الحل شفويًا أو بأي طريقة مماثلة تعتمد على عمليات عقلية عليا (تهانى المزينى، ٢٠١٥).

٤- التقييم المباشر للسلوك أو الأداء: يسعى التقييم البديل إلى تقويم أداء الطالب وسلوكه بطريقة مباشرة لذا فيجب أن تكون المهام المستخدمة في العمليات الاختبارية تمثل عينة ممثلة للمهارات، ويقوم المعلم بالحكم على سلوك أو أداء أو نتائج الطالب بناءً على محكات تقدير محددة وواضحة وتقدير درجاتها يتم بواسطة مجموعة من المصححين المؤهلين ويجب أن تكون تلك المحكات المستخدمة في التقدير واضحة أيضاً للطلاب من أجل توجيههم أنشطتهم وتعلمهم (الفريق الوطنى للتقويم، ٢٠٠٤، ١١-١٢).

٥- الاستناد إلى عينات من الأداء عبر الزمن: مهام الأداء المختلفة تجرى على الطلاب عبر الزمن، وذلك لأنها تهتم بفحص أنماط أعمال الطالب المختلفة وخططة ومخططاته ونموه فيها لتعطى صورة نهائية كاملة عن أداء الطالب النهائى، لذا تتسم مهام الأداء بالاستمرارية (حسن زيتون، ٢٠٠٨).

٦- شمولية وتكامل التقويم: فتقييم الأداء شمولي لانه يعتمد على استخدام أدوات ووسائل متنوعة لتقييم أداء الطالب في مواقف التعلم المختلفة من مختلف جوانبه ليست المعرفة فقط بل يمتد ليشمل المهارات والمعارف واتجاهات الطلاب وعاداتهم وسماتهم (فوزية الدوسرى، ٢٠٠٥).

### تقييم أداء المعلم

تقييم الأداء يمثل مدخلاً لتطوير أداء المعلم، حيث أن توظيف التقويم الحقيقى يتطلب من المعلم الانتقال من ثقافة الصف الدراسي القائم على عمليات التذكر والاسترجاع والمهام التنافسية والاختبارات ذات الفرصة الواحدة والمقارنات بين الطلاب إلى الصف القائم على المهام التعاونية وربط المادة العلمية بالمواقف الحياتية والمشكلات ذات الصلة ببيئة الطالب، وتنوع المهام، وتنوع استراتيجيات وأدوات التقويم، ومقارنة الطالب بذاته لدراسة مدى التحسن فى أدائه (يحيى عبدالخالق، ٢٠١٨، ٢٩٥).

ومما يؤكد ذلك ما أشار لة (على عبد العظيم، ٢٠١٥، ٩٤) بأن التقييم البديل يُعد الحل المناسب للحكم على أداء المعلم وتطوير برامج إعدادة، وهذا ما يؤكد أيضاً دراسة (محمد موسى، ٢٠١٨) من وجود ارتباط موجب دال احصائياً بين درجة استخدام أدوات التقييم الواقعي وبين درجة التطور المهني الذاتي للمعلمين.

كما أن استخدام التقييم التربوي البديل أو الواقعي في تقييم أداء المعلم يهدف إلى تقويم الجوانب المختلفة في شخصية المتعلم المعرفة والمهارية والوجدانية وليس الجانب المعرفي فقط (حجاج غانم، ٢٠٠٧، ٢).

وتؤكد دراسة (محمد يوسف، ٢٠١٧) أن التقييم البديل هو أساس عملية إعداد الطالب المعلم؛ وذلك لانه يتضمن العوامل التي تمكن الطالب من مواكبه متطلبات العصر واكتساب المهارات اللازمة التي تؤهله لأداء مُرضي، لذا فقد أوصت هذه الدراسة بضرورة أستحداث أساليب ووسائل الإعداد والتدريب في كليات التربية المختلفة، وإستخدام التقييم البديل كأحد حلول عملية الاعداد بما يسهم في إعداد الطالب المعلم وتنمية أفكاره والخروج من قيود الاعداد التقليدي الذي يبدأ بالتدريس الالقائي وينتهي بالاختبارات ومنح الدرجات والتقديرية إلى فكر ووسائل وأدوات التقييم التربوي البديل.

ويري (على عبد العظيم، ٢٠١٥، ٩٤) إن التدريس عملية تتطلب مهارات وقدرات وسمات شخصية ومعرفة بمجال التخصص تتضامن جميعها لرسم صورة كاملة عن أداء المعلم؛ وقد يؤدي قياس كل جانب منها منعزلاً عن الجوانب الأخرى خارج نطاق الموقف التعليمي إلى نتائج مضللة عن الأداء، ويعد التقييم البديل الحل المناسب للحكم على أداء المعلم وتطوير برامج إعدادة.

لذا يرى (صفوت توفيق، ٢٠٢١) بضرورة إعتداد استراتيجيات وأدوات ومفاهيم التقييم التربوي البديل ضمن مقررات برامج إعداد المعلمين في كليات التربية المختلفة.

### منهج البحث:

أغلب الابحاث السيكومترية تعتمد في الاصل على الطبيعة الوصفية الكشفية، لذا فإن المنهج المناسب لتلك الدراسة هو المنهج الوصفي الذي يهتم بإستكشاف مكونات التباين الأكثر تأثيراً على ثبات وصدق درجات التقييم، ومعاملات إمكانية التعميم النسبية والمطلقة، وإستكشاف الأبعاد ذات الفعالية في زيادة إمكانية التعميم وثبات درجات أداء الطلاب للمهام المختلفة .

### إجراءات البحث:

#### ١- مجتمع البحث:

تكون مجتمع البحث من طلاب الفرقة الرابعة بكلية التربية - جامعة السويس من الشعب العلمية والأدبية للتعليم العام (إعدادي وثانوي) والتعليم الاساسي (إبتدائي) وشعب (تربية فنية ورياض أطفال وتكنولوجيا التعليم والتربية الخاصة) في العام الجامعي ٢٠٢١ / ٢٠٢٢م، حيث بلغ عدد مجتمع الدراسة (٦٤٢) طالب وطالبة بمتوسط حسابي (٢١.٣) لاعمارهم وإنحراف معياري (٠.٦٤٣).

## ٢- عينة البحث الاستطلاعية

تكونت العينة الاستطلاعية من (٢٠٠) طالب وطالبة من طلاب الفرقة الرابعة بكلية التربية - جامعة السويس من الشعب العلمية والادبية للتعليم العام (إعدادي وثانوي)، بمتوسط عمر زمني = ٢١.١، وانحراف معياري = ٠.٦٨٧، وتم استخدام العينة للاستطلاعية للتحقق من بعض التحليلات الأولية لاختبار المواقف باستخدام النظرية التقليدية للتحقق من الخصائص السيكومترية لها قبل إجراء تحليلات إمكانية التعميم.

## ٣- عينة البحث الاساسية:

تكونت عينة الدراسة النهائية من (٢٧١) طالب وطالبة من طلاب الفرقة الرابعة بكلية التربية - جامعة السويس من الشعب العلمية والادبية للتعليم العام (إعدادي وثانوي) و التعليم الاساسي (ابتدائي).

ثانياً: أدوات البحث:

بطاقة الملاحظة للطالب المعلم (إعداد الباحث، ٢٠٢٢)

تهدف بطاقة الملاحظة إلى جمع بيانات عن الأداء الفعلي للطالب المعلم خلال أدائهم للمهام التدريسية المباشرة مع الطلاب داخل حجرة الدراسة، وذلك من خلال الابعاد التالية:

أ- القدرات المتمثلة في: سرعة استدعاء الافكار، وطلاقة التعبير، وتذكر الاسماء والاشخاص، والانتباه لاشياء كثيرة.

ب- السمات الشخصية المتمثلة في: العرض، والتواد.

ج- المهارات التدريسية المتمثلة في: التخطيط والتنظيم، والتنفيذ، والتقويم، وإدارة الصف، وأستئارة دافعية التلاميذ، واستخدام التكنولوجيا الحديثة، والتواصل.

وفي ضوء قوائم مهارات التدريس التي تضمنتها الدراسات السابقة والاطار النظرى، وقوائم كفايات المعلم، وإستمارات تقييم الطالب المعلم المختلفة من وحدة التربية العملية بكليات التربية بعدة جامعات، وفي ضوء ملاحظة أداء عدد من المعلمين أثناء تدريسهم بالاضافة إلى الاستعانة بأراء أساتذة التربية وعلم النفس والموجهين والمعلمين؛ أمكن جمع بنود بطاقة الملاحظة لتقييم أداء الطالب المعلم داخل حجرة الدراسة.

وتم تصنيف تلك البنود والمهارات والاداءات فى بطاقة الملاحظة للمحاور الاتية:

## جدول (١)

## يوضح تصنيف بنود بطاقة الملاحظة لتقييم الطالب المعلم

بطاقة الملاحظة								
المهارات التدريسية			السمات الشخصية			القدرات		
الاداءات	البعد	م	الاداءات	البعد	م	الاداءات	البعد	م
٤ اداءات	التخطيط والتنظيم	١	٤ اداءات	العرض	١	٣ اداءات	سرعة إستدعاء الافكار	١
٣ اداءات	التنفيذ	٢	٤ اداءات	التواد	٢	٣ اداءات	طلاقة التعبير	٢
٤ اداءات	التقويم	٣				٣ اداءات	تذكر الاسماء والأشخاص	٣
٤ اداءات	إدارة الصف	٤						
٤ اداءات	أستثارة دافعية التلاميذ	٥				٤ اداءات	الانتباه لاشياء كثيرة	٤
٣ اداءات	استخدام التكنولوجيا الحديثة	٦						
٤ اداءات	التواصل	٧						

وبذلك اصبحت بطاقة الملاحظة تتكون من (٤٧) أداءً موزعاً على ثلاثة عشر بعداً تغطي كلاً من القدرات، والسمات الشخصية، والمهارات التدريسية اللازمة لنجاح الطالب/المعلم في عمله. كما تم تقدير مستوى الأداء لكل بند من بنود بطاقة الملاحظة بأستخدام مقياس ثلاثي التدرج (بدرجة كبيرة- بدرجة متوسطة- بدرجة منخفضة).

عرضت الصورة الاولى لبطاقة الملاحظة لتقييم أداء الطالب/المعلم بالاضافة إلى موازين تقدير المهام المختلفة (نموذج الاجابة لكل مهمة) على عينة من المحكمين والخبراء من أساتذة التربية وعلم النفس التربوي والقياس والتقويم لاستطلاع آرائهم في مدى مناسبة محتويات بطاقة ملاحظة للهدف الذي أعدت من أجلها، وقد أبدى المحكمين ملاحظتهم على البطاقة، وميزان التقدير، وأقتصر الباحث على إختيار الاختبارات والبنود بإجماع آراء الخبراء والتي حصلت على نسبة اتفاق ٨٠% فأكثر، وحصلت على قيم ومعايير مقبولة لمعادلة لوشي، ومعادلة Kappa كوهين، ودالة لمعامل كا تربيع. لتصبح بطاقة الملاحظة لتقييم الأداء في صورتها النهائية المعدة للتطبيق على الطالب المعلم لتقييم أداءة تتكون من (١٣) بعداً مختلفاً.

- صدق الاختبار قام الباحث بإخضاع أبعاد بطاقة الملاحظة لتحليل العاملى الاستكشافى المحدد العوامل للتعرف على تشبعات الأبعاد التى تتكون منها بطاقة الملاحظة، حيث تم حساب الصدق العاملى بطريقة المكونات الأساسية ثم طريقة بروماكس وليس التدوير المتعامد لفاريماكس بسبب عدم إستقلالية الأبعاد المستخدمة على العينة الاستطلاعية التى قوامها ٢٠٠ طالب وطالبة، وتم تحديد عدد العوامل لعامل واحد للتحقق من بنية بطاقة الملاحظة، وأوضحت النتائج تشبعات الابعاد على العامل المحدد الذى بلغ نسبة التباين الكلية المفسرة للدرجات لة (٣٤.٤٥٦%) من التباين الكلى، حيث تراوحت قيمة التشبعات من (٠.١٢٥ إلى ٠.٥٩٢) وبالتالي فالتحليل العاملى يتفق مع الأبعاد التى أفترضها الباحث مما يدل على أن المقياس يتصف بصدق بنائى، ويدل على ترابط بين الأبعاد المكونة للعامل وتجانسها، وبالتالي تم التحقق من صدق التكوين الفرضى لأختبار المواقف.

- ثبات الاختبار كما قام الباحث بعد الانتهاء من التحليل العاملي الاستكشافي بحساب ثبات بطاقة الملاحظة بأستخدام طريقة إعادة تطبيق الاختبار بفواصل زمني قدرة (١٥ يوم) علي العينة الاستطلاعية التي بلغت (٢٠٠) طالب وطالبة، كما تمّ أستخدام طريقة ماكدونالد أوميجا بدلاً من ألفا - كرونباخ لان حساب معامل ثبات ألفا كرونباخ يفترض تحقق بعض الافتراضات مثل: أفترض أن تقيس جميع المفردات نفس البنية بنفس درجة الدقة، وإفترض التوزيعات الطبيعية المستمرة للمفردات، وإفترض التوزيع الطبيعي للنتائج الاجمالية، وإفترض عدم ارتباط أخطاء المفردات، وإفترض أحادية البعد، وإفترض البيانات المتصلة (Yang & Green, 2011) وهو ما يصعب تحقيقه بشكل كامل في أداة القياس وبيانات العينة الاستطلاعية، والجدول التالي يوضح قيم معاملات الثبات بالطرق المختلفة:

جدول (٢)  
يوضح ثبات اختبار المواقف

الاختبارات	إعادة التطبيق	معامل ثبات أوميجا (ω)
بطاقة الملاحظة	٠.٧٨٦	٠.٧٨٢

ويتضح من خلال الجدول السابق أنّ بطاقة الملاحظة لتقييم أداء الطالب المعلم تتصف بثبات عالي على إختلاف الطرق المستخدمة في حساب الثبات، مما يدعو إلى الثقة في النتائج التي يمكن التوصل إليها عند أستخدام بطاقة الملاحظة وصلاحيتها للتطبيق.

- تصميمات الدراسة لإمكانية التعميم:

١- تصميمات الملاحظة والقياس: أشتمل تصميم الملاحظة على أربعة أبعاد متقاطعة في موقف القياس وهما "طالب×بُعد×مصحح×فترة"، وأُعتبر الطلاب موضوعاً للقياس (بعداً للتمييز)، وأُعتبرت باقي المصادر وهي الابعاد والمصححين والفترات مصادر للتباين (أبعاد الادائية).

جدول (٣)

تصميم الملاحظة ثلاثي البُعد لأختبار المواقف لتقييم أداء الطالب المعلم

الابعاد	رمز البعد	تصنيف البعد	المستويات الملاحظة	عدد المستويات
الطلاب	P	بعداً للتمييز	من الطالب ١ إلى الطالب ٢٧١	٢٧١
الأبعاد	T	أبعداً للادائية	بطاقة الملاحظة من بعد ١ إلى بعد ١٣	١٣
المُصححين	R		من المصحح ١ إلى المصحح ٢	٢
الفترات	O		من الفترة ١ إلى الفترة ٢	٢

٢- تصميمات دراسات القرار: تهتم هذه المرحلة بمحاولة زيادة معاملات إمكانية التعميم من خلال زيادة أو خفض أو حذف أو إضافة أبعاد لموقف القياس سواء بزيادة عدد الاختبارات أو زيادة عدد المصححين أو زيادة عدد الفترات أو خفضهم أو حذف بُعد أو إضافة بُعد ، وتسمى تلك المرحلة أيضاً بمرحلة التحسين.

ثالثاً: الاساليب الاحصائية:

١- معامل ثبات أوميغا، ومعاملات الارتباط والتحليل العاملي الاستكشافي بأستخدام برنامج

SPSS

٢- معاملات إمكانية التعميم النسبية والمطلقة ومكونات التباين المختلفة بأستخدام برنامج EduG

رابعاً: عرض نتائج البحث ومناقشتها:

(١) عرض نتائج دراسة إمكانية التعميم

جدول (٤)

تحليل التباين للتصميم الثلاثي البعد "طالب×بُعد×مُصح×فترة" لبطاقة الملاحظة

المكونات								
مصدر التغير	مجموع المربعات	درجات الحرية	متوسط المربعات	العشوائية	المختلطة	المُصححة	النسبة (%)	الخطأ المعياري
طالب (P)	٧٥٧٦.٠٧٩٩	٢٧٠	٢٨.٠٥٩٦	٠.٤٥٢٣	٠.٤٥٢٣	٠.٤٥٢٣	٣٦.٩	٠.٣٤٦٦
بُعد (T)	١٧٠٠٧.٩١٦١	١٢	١٤١٧.٣٢	١.٢٩٠٧	١.٢٩٠٧	١.٢٩٠٧	٢٣.٩	٠.٢٩٤٢
مُصح (R)	٠.٦٩٥٥	١	٠.٦٩٥٥	٠.٠٠١٠	٠.٠٠١٠	٠.٠٠١٠	٠.٠	٠.٠٠٠٠٨
فترة (O)	١٧.٥٢٨٣	١	١٧.٥٢٨	٠.٠٠٠١	٠.٠٠٠١	٠.٠٠٠١	٠.٠	٠.٠٠٠٢٣
طالب-بُعد (P×T)	١٢٩٢٢.٢٧٦٢	٣٢٤٠	٣.٩٨٨٤	٠.٦٣٥٤	٠.٦٣٥٤	٠.٦٣٥٤	١١.٧	٠.٠٢٨٤
طالب-مُصح (P×R)	٣٥٦.٤٠٠٧	٢٧٠	١.٣٢٠٠	٠.٠٠١٦	٠.٠٠١٦	٠.٠٠١٦	٠.٠	٠.٠٠٠٦٧
طالب-فترة (P×O)	٥٧٥.٢٦٠١	٢٧٠	٢.١٣٠٦	٠.٠١٢٦	٠.٠١٢٦	٠.٠١٢٦	٠.٣	٠.٠٠٠٨٧
بُعد-مُصح (T×R)	٥٤.٤٨٠٧	١٢	٤.٥٤٠١	٠.٠٠٢٥	٠.٠٠٢٥	٠.٠٠٢٥	٠.١	٠.٠٠٠٣٩
بُعد-فترة (T×O)	١٧٢.٧٧٧٠	١٢	١٤.٣٩٨١	٠.٠١٩٩	٠.٠١٩٩	٠.٠١٩٩	٠.٥	٠.٠١٠٣
مُصح-فترة (R×O)	٦.٥١٥٠	١	٦.٥١٥٠	٠.٠٠٠٨	٠.٠٠٠٨	٠.٠٠٠٨	٠.٠	٠.٠٠٠١٦
طالب-بُعد-مُصح (P×T×R)	٣٥٥١.١٧٣١	٣٢٤٠	١.٠٩٦٠	٠.٠٤٤١	٠.٠٤٤١	٠.٠٤٤١	٠.٠	٠.٠٢٠٠
طالب-بُعد-فترة (P×T×O)	٤٩٧٤.١٨٤٥	٣٢٤٠	١.٥٣٥٢	٠.١٧٥٥	٠.١٧٥٥	٠.١٧٥٥	٤.٦	٠.٠٢٤١
طالب-مُصح-فترة (P×R×O)	٣٩١.٨١٢٠	٢٧٠	١.٤٥١٢	٠.٠٢٠٥	٠.٠٢٠٥	٠.٠٢٠٥	٠.٥	٠.٠٠٠٩٨
بُعد-مُصح-فترة (T×R×O)	٣٩.١٧٤١	١٢	٣.٢٦٤٥	٠.٠٠٧٧	٠.٠٠٧٧	٠.٠٠٧٧	٠.٢	٠.٠٠٠٤٦
طالب-بُعد-فترة-مُصح (P×T×R×O, e)	٣٨٣٧.٢٤٨٩	٣٢٤٠	١.١٨٤٣	١.١٨٤٣	١.١٨٤٣	١.١٨٤٣	٢١.١	٠.٠٢٩٤
المجموع	٥١٤٨٣.٥٢٢٢	١٤٠٩١					%١٠٠	

نلاحظ من الجدول السابق أن نسبة التباين الحقيقي (التباين المنتظم) الذي يصف الفروق الحقيقية بين الافراد في أختبارات صيغة السمات الشخصية من بطارية تقييم أداء الطالب/المعلم بلغت (٣٦.٩%) من التباين الكلي، وتشير تلك النسب إلى أنه بالأخذ في الاعتبار متوسط المُصححين والاختبارات وعدد مرات التطبيق فإن هناك فروقاً حقيقية بين الطلاب في أبعاد بطاقة الملاحظة لتقييم أداء الطالب المعلم، ويوضح الجدول التالي معاملات إمكانية التعميم النسبية والمطلقة، ومكونات التباين لكل مصدر من مصادر الاخطاء المختلفة المؤثرة على دقة القياس على النحو التالي:

## جدول (٥)

## تحليل إمكانية التعميم لتصميم القياس طالب/ بعد مُصحح فترات (P/ TRO)

مصدر التباين	تباين التمييز	تباين الخطأ النسبي	نسبة تباين الخطأ النسبي	تباين الخطأ المطلق	نسبة تباين الخطأ المطلق
طالب (P)	٢.٤٥٢٣	.....	.....	.....	.....
بُعد (T)	.....	.....	.....	.....	٥٢.٢
مُصحح (R)	.....	.....	.....	.....	٠.٠
فترة (O)	.....	.....	.....	.....	٠.٠
طالب- بُعد (PT)	.....	.....	٠.٠٤٨٩	.....	٢٥.٧
طالب- مُصحح (PR)	.....	.....	٠.٠٠٠٠	.....	٠.٠
طالب- فترة (PO)	.....	.....	٠.٠٠٦٣	.....	٣.٣
بُعد - مُصحح (TR)	.....	.....	.....	.....	٠.١
بُعد - فترة (TO)	.....	.....	.....	.....	٠.٤
مُصحح- فترة (RO)	.....	.....	.....	.....	٠.١
طالب- بُعد - مُصحح (PTR)	.....	.....	٠.٠٠٠٠	.....	٠.٠
طالب- بُعد - فترة (PTO)	.....	.....	٠.٠٠٦٧	.....	٣.٥
طالب- مُصحح- فترة (PRO)	.....	.....	٠.٠٠٥١	.....	٢.٧
بُعد - مُصحح- فترة (TRO)	.....	.....	.....	.....	٠.١
طالب- بُعد - مُصحح - فترة (PTRO, e)	.....	.....	٠.٠٢٢٨	.....	١٢.٠
مجموع التباينات	٢.٤٥٢٣	.....	٠.٠٨٩٨	.....	%١٠٠
الانحراف المعياري	.....	.....	٠.٦٧٢٥	.....	.....
الخطأ المعياري النسبي	.....	.....	٠.٢٩٩٧	.....	.....
الخطأ المعياري المطلق	.....	.....	٠.٤٣٦٣	.....	.....
معامل إمكانية التعميم النسبي	.....	.....	٠.٨٣	.....	.....
معامل إمكانية التعميم المطلق	.....	.....	٠.٧٠	.....	.....

يتضح من الجدول السابق أن أكبر مصدر تباين في القياس النسبي يرجع إلى تفاعل طالب- بُعد والذي بلغت قيمة (٥٤.٤%) وهو يوضح الاختلاف بين متوسط أداء الطلاب من بُعد إلى آخر في صيغة بطاقة الملاحظة، وثاني أكبر مصدر لتباين الخطأ يرجع إلى تفاعل طالب- بُعد - مُصحح- فترة الممزوج بالاختلاف العشوائية والذي بلغ قيمة (٢٥.٣%) والذي يعبر عن اختلاف ترتيب أداء الطلاب عبر مختلف الأبعاد والمُصححين والفترات المختلفة.

في حين جاء مكونات تباين الخطأ الأخرى من متوسطة إلى منخفضة بحيث كان مكون التباين الناتج من تفاعل طالب- بُعد - فترة بلغت قيمة (٧.٥%)، والذي يفسر تباين ترتيب الطلاب بين الأبعاد والفترات المختلفة، ومكون التباين الذي يعبر عن التفاعل بين طالب- فترة بلغ قيمة (٧%)، ومكون التباين الذي يوضح التفاعل بين طالب- مُصحح- فترة الذي بلغ قيمة (٥.٧%) وهي قيم منخفضة لا تؤثر على ثبات درجات الطلاب أو إمكانية التعميم. في حين جاءت باقي مكونات التباين الأخرى منخفضة جداً لا تؤثر على إمكانية التعميم.

كما يتضح من الجدول السابق أن أكبر مصدر تباين في القياس المطلق يرجع إلى الأبعاد والذي بلغت قيمة (٥٢.٢%). ثم يأتي مكون التباين الناتج من تفاعل طالب- بُعد في المرتبة الثانية من

مصادر تباين الاخطاء حيث بلغ قيمته (٢٥.٧%) وهو يوضح الاختلاف بين متوسط أداء الطلاب من بعد إلى آخر في صيغة بطاقة الملاحظة ومن فترة إلى أخرى.

في حين جاء مكون التباين الناتج من تفاعل طالب- بعد- مُصحح- فترة الممزوج بالاطء العشوائية بقيمة متوسطة بلغت (١٢%). بينما جاءت باقي مكونات التباين الأخرى منخفضة ولا تؤثر على إمكانية التعميم فقد تراوحت بين (٠.٠% و ٣.٥%) وكانت ترجع إلى المُصححين، والفترات، وتفاعل طالب- مُصحح، وتفاعل طالب- بعد- مُصحح، وتفاعل طالب- بعد- مُصحح- فترة، وتفاعل طالب- مُصحح- فترة، وتفاعل طالب- فترة، وتفاعل طالب- بعد- فترة.

كما كانت معاملات إمكانية التعميم جيدة إلى مقبولة حيث بلغ معامل إمكانية التعميم النسبي (٠.٨٣) وهي قيمة جيدة، وبلغ معامل إمكانية التعميم المطلق (٠.٧٠) وهي قيمة مقبولة مما يدل على أن درجات الطلاب في صيغة بطاقة الملاحظة تتسم بالثبات أو إمكانية التعميم.

وتتفق تلك النتائج مع ما توصلت إليه نتائج دراسات تقييم الأداء مثل دراسة شافيلسون وآخرون (Shavelson et al., 1993)، ودراسة ويب وآخرون (Webb et al., 2000)، ودراسة سميث وكوليكوش (Smith & Kulikowich, 2004) خاصة فيما يتعلق بارتفاع مكون تباين تفاعل طالب- بُعد/أختبار، وتباين تفاعل طالب- بُعد/أختبار - فترة، وتتفق أيضاً مع دراسة ني وآخرون (Nie et al., 2007)، ودراسة جبريل (Gebriel, 2009)، ودراسة جولد وجيلبال (Guler & Gelbal, 2010)، ودراسة تيلور وباستور (Taylor & Pastor, 2013)، دراسة (صبري محمود، ٢٠١٦)، ودراسة (طباع فاروق، ٢٠٢٠) الذين أوضحوا ارتفاع كبير بمكون تباين تفاعل طالب-أختبار/بُعد مقارنة بباقي مكونات التباين الأخرى التي جاءت ضعيفة.

وفيما يتعلق بمعاملات إمكانية التعميم النسبية والمطلقة فقد جاءت مقبولة إلى جيدة وتتفق تلك النتائج مع بعض الدراسات التي توصلت إلى بلوغ معاملات إمكانية تعميم تصل للحد الأدنى المطلوب مثل دراسة كلاً من ويب وآخرون (Webb et al., 2000)، ودراسة سميث وكوليكوش (Smith & Kulikowich, 2004) ودراسة لي وكانتور (Lee & Kantor, 2007)، ودراسة جولد وجيلبال (Guler & Gelbal, 2010)، ودراسة (صبري محمود، ٢٠١٦).

في حين لا تتفق جزئياً مع نتائج دراسة ميسبي وبارينز (McBee & Barenz, 1998)، ودراسة شيفيلسون وآخرون (Shavelson et al., 1999)، ودراسة ني وآخرون (Nie et al., 2007)، ودراسة باستور (Pastor, 2013)، ودراسة (طباع فاروق، ٢٠٢٠) التي أظهرت ضعفاً في معاملات إمكانية التعميم النسبية والمطلقة بالرغم من اقترابها للحد الأدنى المقبول وهو (٠.٨٠) كما ذكر (Bain & Pini, 1996)، حيث جاءت معاملات إمكانية التعميم المطلقة ضعيفة نسبياً.

ترجع منطقية زيادة مكونات تباين تفاعل طالب- بُعد إلى وجود اختلاف في متوسط أداء الطلاب من بُعد لآخر داخل بطاقة الملاحظة لتقييم الاداء، وذلك بسبب: إختلاف كل بُعد عن الآخر في الطبيعة

وفى التكوين وتباينها من حيث خصائصها، ودرجة تركيبها، وقد يرجع أيضا ارتفاع مكون تباين تفاعل طالب- بُعد إلى تعدد وتنوع الأبعاد داخل بطاقة الملاحظة حيث تشمل بطاقة الملاحظة على مجموعة من السمات والقدرات والمهارات التدريسية المتحصل عليها من تحليل عمل مهنة المعلم.

كما يمكن تفسير ذلك المكون إلى وجود مشكلة معرفية ناتجة عن أنتقال أثر التعلم لدى الطلاب من اختبارات معينة إلى أخرى كما ذكرت دراسة باركيز (Parkes, 2001)، ودراسة سكالون (Scallon, 2004)، ودراسة هوانج (Huang, 2009) وذلك بسبب عدم قدرة الطلاب على نقل معارفهم من سياق الدراسة داخل الجامعة إلى السياق التطبيقي والعمل.

وفيما يخص ارتفاع مصدر تباين تفاعل طالب- بُعد - فترة، وتفاعل طالب-فترة فهذه النتيجة قد ترجع بشكل كبير إلى تغير في أداء الطلاب وترتيبهم عبر الأبعاد المختلفة داخل بطاقة الملاحظة وخلال عدد مرات التطبيق بحيث أن الطلاب انجزوا بعض الأبعاد بشكل أفضل في الفترة الثانية ولكن أنجزوا بعض الأبعاد الأخرى بشكل أقل في الفترة الأولى؛ مما يدل على تغير الطلاب في استراتيجيات حل الأبعاد خلال مرات التطبيق ومحاولة تحسين أدائهم في المرة الثانية.

## ٢) عرض نتائج دراسة القرار

وفيما يتعلق بدراسات القرار التي تتناول زيادة عدد مستويات الأبعاد لزيادة معاملات إمكانية التعميم النسبية والمطلقة، وذلك من خلال تجريب أفضل مستويات الأبعاد التي تزيد من نطاق الدرجة الشاملة وخفض تباينات الخطأ للحصول على درجات قابلة للتعميم على مواقف قياس أخرى. ولتحقيق ذلك سوف نقدم دراستين للقرار أهتمت الأولى بزيادة عدد الأبعاد وعدد المُصححين بينما أهتمت الدراسة الثانية من دراسات القرار بزيادة عدد مرات التطبيق.

- الدراسة الأولى من دراسات القرار والتي أهتمت بزيادة عدد الأبعاد وعدد المُصححين:

### جدول (٦)

دراسات القرار للتصميم الثلاثي البعد "طالب × بُعد × مُصحح × عدد مرات التطبيق"

بطاقة الملاحظة									
٢٠	٢٠	٢٠	١٦	١٦	١٦	١٣	١٣	١٣	عدد الأبعاد
٤	٣	٢	٤	٣	٢	٤	٣	٢	عدد المُصححين
٢	٢	٢	٢	٢	٢	٢	٢	٢	عدد مرات التطبيق
٠.٩٠	٠.٩٠	٠.٨٩	٠.٨٨	٠.٨٨	٠.٨٨	٠.٨٦	٠.٨٦	٠.٨٦	إمكانية التعميم النسبي
٠.٨٠	٠.٧٩	٠.٧٩	٠.٧٧	٠.٧٦	٠.٧٦	٠.٧٣	٠.٧٢	٠.٧٢	إمكانية التعميم المطلق

– الدراسة الثانية من دراسات القرار والتي أهتمت بزيادة عدد الفترات وكانت نتائجها كالتالي:

### جدول (٧)

دراسات القرار للتصميم الثلاثي البعد "طالب × بُعد × مُصحح × عدد مرات التطبيق"

بطاقة الملاحظة									
٢٠	٢٠	٢٠	١٦	١٦	١٦	١٣	١٣	١٣	عدد الأبعاد
٢	٢	٢	٢	٢	٢	٢	٢	٢	عدد المُصححين
٤	٣	١	٤	٣	١	٤	٣	١	عدد مرات التطبيق
٠.٩١	٠.٩٠	٠.٨٥	٠.٨٩	٠.٨٩	٠.٨٣	٠.٨٧	٠.٨٧	٠.٨٠	إمكانية التعميم النسبي
٠.٨١	٠.٨٠	٠.٧٦	٠.٧٧	٠.٧٦	٠.٧٢	٠.٧٣	٠.٧٣	٠.٦٨	إمكانية التعميم المطلق

ويتضح من الجدول (٦) أن بطاقة الملاحظة تتطلب (١٣) بعد مع عدد (٣) مُصححين للوصول لمعامل إمكانية التعميم النسبي (٠.٨٦) ومعامل إمكانية التعميم المطلق (٠.٧٢)، وفي حالة زيادة عدد الأبعاد إلي (١٦) بعد يرتفع معامل إمكانية التعميم النسبي إلى (٠.٨٨) ومعامل إمكانية التعميم المطلق إلى (٠.٧٦)، وفي حالة ارتفاع عدد الأبعاد إلى (٢٠) بعد يرتفع معامل إمكانية التعميم النسبي بمقدار ضئيل جداً إلى (٠.٨٩) ومعامل إمكانية التعميم المطلق يصل إلى (٠.٧٩) لعدد (٢) مُصحح وبزيادة عدد المُصححين من (٢) مُصححين إلى (٤) مُصححين يرتفع معامل إمكانية التعميم بمقدار ضئيل جداً ليصل لمعامل إمكانية تعميم نسبي مقدارة (٠.٩٠) ومعامل إمكانية تعميم مطلق مقدارة (٠.٨٠).

كما أوضحت نتائج جدول (٧) ان زيادة عدد مرات التطبيق يزيد من معاملات إمكانية التعميم مع زيادة عدد الأبعاد، فزيادة عدد الفترات من (١) فترة إلى (٤) فترات يزيد من معامل إمكانية التعميم من (٠.٨٠) إلى (٠.٨٧) في حالة وجود (١٣) بعد، ويزيد معامل إمكانية التعميم المطلق من (٠.٦٨) إلى (٠.٧٣)، و زيادة عدد الفترات من (١) فترة إلى (٤) فترات يزيد من معامل إمكانية التعميم من (٠.٨٣) إلى (٠.٨٩) في حالة وجود (١٦) بعد، ويزيد معامل إمكانية التعميم المطلق من (٠.٧٢) إلى (٠.٧٧)، وفي حالة وجود (٢٠) بعد وزيادة عدد الفترات من (١) فترة إلى (٤) فترات يزيد من معامل إمكانية التعميم من (٠.٨٥) إلى (٠.٩١)، ويزيد معامل إمكانية التعميم المطلق من (٠.٧٦) إلى (٠.٨١).

وبالتالي توصلت نتائج دراسات القرار إلى أن زيادة عدد الأبعاد وزيادة عدد مرات التطبيق تؤدي إلى زيادة معاملات إمكانية التعميم النسبية والمطلقة أفضل من زيادة عدد المُصححين.

وتتفق تلك النتائج مع عدد من الدراسات السابقة التي تؤكد أن زيادة عدد الاختبارات والأبعاد وزيادة عدد الفترات تساهم بشكل كبير في زيادة معاملات إمكانية التعميم النسبية والمطلقة، وخفض مكونات التباين التي تؤثر على التباين الكلي لدرجات أداء الطلاب مثل دراسة شافيلسون وآخرون (Shavelson et al., 1993)، ودراسة رولز-بريمو وآخرون (Rulz-Primo et al., 1993)،

ودراسة شافيلسون وآخرون (Shavelson et al., 1999)، ودراسة ويب وآخرون (Webb et al., 2000)، ودراسة ني وآخرون (Nie et al., 2007)، ودراسة لي وكانتور (Lee & Kantor, 2007)، ودراسة هوانج (Huang, 2009)، ودراسة (طباع فاروق، ٢٠١٦)، ودراسة تيلور وباستور (Taylor & Pastor, 2013)، ودراسة (طباع فاروق، ٢٠٢٠)، كما أوضحت دراسة جبريل (Gebril, 2009) أن زيادة عدد المُصححين لا يؤدي إلى زيادة معاملات إمكانية التعميم، وأوضحت دراسة جولر وجيلبال (Guler & Gelbal, 2010) أن زيادة عدد الاختبارات أفضل من زيادة عدد المصححين لزيادة معاملات إمكانية التعميم النسبية والمطلقة.

وتختلف تلك النتائج مع ما توصل إليه دراسة كازانوف وديميوس (Casanova & Demeuse, 2011) التي هدفت إلى فحص مصادر الخطأ المؤثرة على ثبات اختبار التعبير الكتابي في اللغة الفرنسية باستخدام نظرية إمكانية التعميم ونموذج راش متعدد الأبعاد، حيث طبق الاختبار على (٣٣) طالباً، وأوضحت الدراسة أن زيادة عدد المُصححين تزيد من معامل إمكانية التعميم المطلق أفضل من زيادة عدد الاختبارات.

وتختلف جزئياً مع دراسة (عصام الدسوقي، ٢٠٠٥) التي أوضحت أن زيادة عدد المُصححين وعدد المفردات تؤدي إلى زيادة معاملات إمكانية التعميم أفضل من زيادة عدد الفترات، ولكن تتأثر معاملات إمكانية التعميم بزيادة عدد البنود مقارنة بزيادة عدد المُصححين، كما أوضحت دراسة (صبرى محمود، ٢٠١٦) أن زيادة عدد المقيمين إلى ثلاث مقيمين وزيادة عدد مرات التطبيق إلى ثلاث مرات يزيد من معاملات إمكانية التعميم للمقياس.

وترجع منطقية ذلك إلى إنة عند الأجابة على التساؤل الأول أوضحت النتائج أن مكونات التباين الأكثر تأثيراً على ثبات درجات أداء الطلاب ترجع إلى مكون تباين تفاعل طالب-بُعد، وتباين الأبعاد، وتباين تفاعل طالب-بُعد-فترة، وبدرجة أقل تفاعل طالب-فترة ولا ترجع إلى المُصحح وتفاعلاته مع الأبعاد الأخرى، ومن أجل خفض مقادير مكونات تباين طالب-بُعد، وتباين الأبعاد، وتباين تفاعل طالب-بُعد-فترة، وتباين تفاعل طالب-فترة يتطلب ذلك زيادة عدد الأبعاد وعدد مرات التطبيق كطرق أكثر فاعلية من زيادة عدد المُصححين للوصول إلى معاملات إمكانية تعميم جيدة تساعد في تقييم أداء الطالب/المعلم إلا أن زيادة عدد مرات التطبيق تنتج عنها تكاليف وجهود إضافية قد لا تصل إليها تكاليف وجهود زيادة عدد الأبعاد.

وبالنظر إلى نتائج الدراسة الحالية، ونتائج الدراسات السابقة فإن عملية زيادة عدد الأبعاد وزيادة عدد مرات التطبيق تمثل الطرق الأفضل لزيادة معاملات إمكانية التعميم النسبية والمطلقة، ولكن ذلك يؤدي إلى مزيد من الجهد والوقت والتكلفة عند تقييم الاداء.

### توصيات البحث:

في ضوء ما توصلت إليه الدراسة الحالية من نتائج يمكن اقتراح التوصيات الآتية:

- ١) إدراج مفاهيم ومبادئ نظرية إمكانية التعميم ببرامج الدراسات العليا لباحثين الماجستير والدكتوراه لإستخدامها في تقدير دقة القياسات وإجراءات تحسين عمليات القياس.
- ٢) التنمية المهنية والمعرفية للمعلمين وطلاب كليات التربية وتغيير تفكيرهم للانتقال من أساليب التقويم التقليدي إلى أساليب تقييم الأداء المختلفة، ومعرفة مصطلحاته وأدواته بأعتبارة مجالاً مستحدثاً.
- ٣) أستخدام بطاقة الملاحظة كأختبارات دورية للمعلمين في المدارس لضمان الحد الأدنى من الأداء المطلوب.

## المراجع

- تهانى عبدالرحمن المزينى (٢٠١٥). تصور مقترح للكفايات البحثية لمعلمات العلوم المرحلة المتوسطة فى ضوء أدوار معلم القرن ٢١ بمدينة الرياض، المجلة العربية للعلوم التربوية والنفسية، ٥(٢٣)، ٥٣٥ - ٥٦٤.
- حجاج غانم أحمد (٢٠٠٧). المشكلات المرتبطة بتطبيق أسلوب التقويم التربوي الشامل في المدارس الابتدائية من وجهة نظر المعلمين ذوي المعارف الكافية عن هذا الاسلوب وعلاقة ذلك ببعض المتغيرات الاسمية. مجلة البحث في التربية وعلم النفس - كلية التربية جامعة المنيا، ٢١(١)، ١ - ٧٣.
- حسن حسين زيتون (٢٠٠٨). أصول التقويم والقياس التربوي "المفاهيم والتطبيقات". الرياض: الدار الصولتية.
- صبري محمود عبدالفتاح أمين (٢٠١٦). ثبات تقييمات الكفاءة الاجتماعية لدي الاطفال بأستخدام النسخة العربية لمقياس الكفاءة الاجتماعية المعدل في ضوء نظرية إمكانية التعميم. مجلة العلوم التربوية، (٢)، ١٣١ - ١٨١.
- صفوت توفيق هنداوي (٢٠٢١). برنامج تدريبي قائم علي التفكير التأملية لتنمية مهارات استخدام اساليب التقويم اللغوي البديل لدي معلمي اللغة العربية في المرحلة الاعدادية وأثرة علي تنمية مهارات الانتاج اللغوي لدي طلابهم. جامعة عين شمس. مجلة البحث العلمي في التربية، ١٢(٢٢)، ١٣٨ - ١٨١.
- صلاح الدين محمود علام (٢٠٠٠). القياس والتقويم التربوي والنفسي - أساسياته وتطبيقاته وتوجهاته المعاصرة. القاهرة: دار الفكر العربي.
- صلاح الدين محمود علام (٢٠٠٤). التقويم التربوي البديل أسسه النظرية والمنهجية وتطبيقاته الميدانية. القاهرة: دار الفكر العربي.
- طباع فاروق (٢٠١٩). من النظرية الكلاسيكية للاختبارات إلي نظرية إمكانية التعميم. جامعة مولود معمري - تيزي وزو.
- طباع فاروق (٢٠٢٠). استخدام نظرية إمكانية التعميم فى تقدير ثبات اختبار تقييم كفاءة الرياضيات لدي طلاب السنة الرابعة ابتدائي. المجلة الاردنية في العلوم التربوية. ١٦ (١)، ١٨-١.
- عادل سرايا (٢٠٠١). التقويم الحقيقى، مجلة التدريب والتقنية، ٤٠ (٧٤).
- عبد العزيز محمد (٢٠١٩). تقويم ممارسات التقويم لدي أعضاء هيئة التدريس بجامعة المنيا في ضوء استراتيجيات التقويم البديل. جامعة أسيوط. مجلة كلية التربية، ٣٥ (٦)، ٥١٩ - ٥٩٦.

- عبدالله السعدوى (٢٠١٠). دليل المعلم للتقويم المعتمد على الاداء من النظرية إلى التطبيق. الرياض: مكتبة التربية لدول الخليج العربي.
- عصام الدسوقي (٢٠٠٥). استخدام نظرية إمكانية التعميم في تقدير ثبات بطاقة ملاحظة أداء الطلاب في التربية الميدانية. مجلة كلية التربية بدمياط، ١(٤٨).
- علي عبد العظيم علي (٢٠١٥). التقويم البديل: مدخل للارتقاء بأداء المعلم وتطوير برامج إعدادة. المؤتمر العلمي الرابع والعشرون للجمعية المصرية للمناهج وطرق التدريس تحت عنوان: برامج إعداد المعلمين في الجامعات من أجل التميز، ٩١ - ١١٩.
- الفريق الوطني للتقويم (٢٠٠٤). استراتيجيات التقويم وأدواته (الإطار النظري). المملكة الأردنية الهاشمية: وزارة التربية والتعليم.
- فؤاد أبو حطب (١٩٩٢). دليل المعلم في تقويم الطالب. المركز القومي للامتحانات والتقويم التربوي بالاشتراك مع وزارة التربية والتعليم. القاهرة: دار غريب للطباعة.
- فوزي عزت علي (١٩٩٨). البنية العاملية لأختبار الاستعداد للقبول بكلية التربية. المجلة المصرية للدراسات النفسية، ٨ (١٩)، ٧٩ - ١٠٣.
- فوزي عزت علي (١٩٩٧). أختبار الاستعداد للقبول بكليات التربية. القاهرة، مطبوعات المركز القومي للامتحانات والتقويم التربوي.
- فوزية محمد الدوسري (٢٠٠٥). تقويم مهارات استخدام الخرائط لدى تلميذات الصف الثالث من المرحلة المتوسطة بالمملكة العربية السعودية، مجلة الجمعية التربوية للدراسات الاجتماعية، (٤).
- محمد بن راشد عبدالكريم (٢٠٢١). واقع التقويم البديل في التعليم العام بمحافظة القنفذة من وجهه نظر المعلمين. رابطة التربويين العرب. السعودية، ١٣٧(١٣٧)، ١٤١ - ١٦٢.
- محمد موسى نصر الله (٢٠١٨). واقع استخدام التقويم الواقعي وعلاقته بالتطور المهني الذاتي لدى معلمي المرحلة الاساسية الدنيا في محافظة جنين. رسالة ماجستير غير منشورة، كلية العلوم التربوية، جامعة القدس.
- محمد يوسف أحمد (٢٠١٧). التقويم التربوي البديل ودوره في تنمية كفايات الطالب/ المعلم بكليات التربية. كلية التربية جامعة بيشة. المملكة العربية السعودية. المجلة العربية للعلوم ونشر الابحاث، ١(١).
- ياسر عبدالحافظ و عبدالله سالم العازمي (٢٠١٣). تقويم اداء معلم التعليم الاساسي في ضوء معايير المجتمع المدرسي والمحلي. مجلس النشر العلمي بجامعة الكويت، ٢٧(١٠٦)، ٥٥ - ٩٩.

يحي عبد الخالق يوسف (٢٠١٨). المعينات التي تواجه تطبيق التقويم الحقيقي في تعليم وتعلم مقررات التربية الإسلامية بمدارس منطقة تبوك التعليمية. المجلة الدولية للدراسات التربوية والنفسية، ٢ (٣)، ٢٩٢ - ٣١٦.

يسري زكي عبود (٢٠١٦). التقويم البديل كاتجاه حديث في تقويم أداء الطلاب. المؤتمر الدولي لكلية التربية بجامعة الملك خالد تحت عنوان: المعلم وعصر المعرفة الفرص والتحديات. كلية التربية بجامعة الملك خالد - السعودية، ١٨٠٣ - ١٨٣٨.

Atila, Y. & Ezel, T. (2010). The Examination of Reliability According to Classical Test and Generalizability on a Job Performance Scale. Educational Sciences: Theory & Practice, 10 (3), 1847- 1854.

Bain, D. & Pini, G. (1996). To rate your ratings. Generalizability: a user's guide. Geneva: psych pedagogical research center. General Management of the Orientation cycle.

Bakker, M., Sanders, P., Beijaard, D., Roelofs, E., Tigelaar, D., & Verloop, N. (2008). Reliability and generalizability of performance judgments based on a video portfolio. Pedagogische Studiën 85(4), 240-260

Bell, C. A., Gitomer, D. H., McCaffrey, D. F., Hamre, B. K., Pianta, R. C., & Qi, Y. (2012). An Argument Approach to Observation Protocol Validity. Educational Assessment, 17(2-3), 62-87.

Bertrand, R. & Blais, J. (2004). Models de mesure. Canada: presses de L'Universite du Quebec.

Brennan, R. L. (2001). Generalizability Theory. New York: Springer.

Brennan, R. L. (2010). Generalizability Theory. Educational Measurement, 61-68. At: <http://www.education.uiowa.edu>.

Broger, T., Saintmeen, C., & Keleven, W. (2012). The generalizability of systematic rating scales across time: A preliminary analysis using generalizability theory. Journal of educational psychology, 3, 75-92.

Cardinet, J., Sandra, J., & Pini, G. (2010). Applying Generalizability Theory using EdUG. New York: Routledge.

Carminla, S. (2010). Social function assessment tools for children within the framework of the generalizability theory. Clinical psychology review, 4, 52-65.

Carrie, L. S. (2013). Using Generalizability theory to measure sources of variance on a special education teacher observation tool. A dissertation of doctor of education in curriculum and instruction, Boise state university.

Cary, Maranell.M. (1994). Scalling - Chicago: Aldine publishing company.

Casanova, D., & Demeuse, M. (2011). Analysis of different facets influencing the reliability of the written expression test of French as

- a foreign language test. *Measurement and Evaluation in Education*, 34(1), 25-53.
- Christ, G. (2012). An Application of the generalizability theory to the social skills improvement system- Rating scales. *North American journal of psychology*, 3, 25-37.
- Chen, E., Niemi, D., Wang, H., & Mirocha, J. (2007). Examining the generalizability of direct writing assessment tasks. University of California. Los Angeles: CRESST (CSE Technical Report N° 718). Available online at: [cresst.org/publications/cresst-publication-3089](http://cresst.org/publications/cresst-publication-3089).
- Crocker, L. & Algina, J. (2006). *Introduction to classical and modern test theory*. Belmont, CA: Wadsworth Publishing Company.
- Demet, H. & Erkut, A. (2018). Factors that affect the performance of teachers working in secondary-level education. *Journal of Academy of Educational Leadership journal*, 22(1), 1- 19.
- Everitt, B. & Skrondal. J. (2010). *The Cambridge dictionary of statistics*. London: Cambridge university press.
- Gao, X., & Brennan, R. (2001). Variability of estimated variance components and related statistics in a performance assessment. *Applied Measurement in Education*, 14(2), 191–203.
- Gebril, A. (2009). Score generalizability of academic writing tasks: Does one test Methods fit it all? *Language Testing*, 26(4), 507-531.
- Guler, N. & Gelbal, S. (2010). Studying Reliability of Open Ended Mathematics Items According to the Classical Test Theory and Generalizability Theory. *Educational Sciences: Theory & Practice*, 10 (2), 1011-1019.
- Haladyna, T. (2004). *Developing and validating multiple- choice test items*. New jersey: Lawrence Erlbaum Associates publishers.
- Hill, H. C., Charalambous, C. Y., & Kraft, M. A. (2012). When Rater Reliability Is Not Enough Teacher Observation Systems and a Case for the Generalizability Study. *Educational Researcher*, 41(2), 56–64.
- Huang, C. H. (2009). Magnitude of task-sampling variability in performance assessment: a meta-analysis. *Educational and Psychological Measurement*, 69 (6), 887- 912.
- Keller, L., Clauser, B., & Swanson, D. (2010). Using multivariate generalizability theory to assess the effect of content stratification on the reliability of a performance assessment, *Advanced in Health Science Education*, (15), 717-733.
- Kerry, E. (2015). Teacher Performance Appraisal: More about performance or Development. *Australian Journal of Teacher Education*, 40(9), pp 102- 116.
- Lane, S., Liu, M., Ankenmann, R., & Stone, C. (1996). Generalizability and validity of a mathematics performance assessment. *Journal of Educational Measurement*, 33(1), 71-92.

- Lee, Y., & Kantor, R. (2007). Evaluation prototype tasks and alternative rating schemes for a new ESL writing test through G-theory. *International journal of Testing*, 7(4), 353-385.
- Linn, R., Baker, E. & Dunbar, S. (1991). Complex Performance- Based Assessment: Expectations and Validation Criteria. *Educational Researcher*, 20(8), 15-21.
- Lueng, T., Lui, R. & Yim, H. (2010). Advances in the assessment of social competence: Finding from a preliminary dependability investigation in a Japanese sample. *Japanese Journal of educational psychology*, 2, 56-79.
- Maria, L. & Maria, H. (2013). Approaches to Teachers' Performance Assessment for Enhancing Quality of Education at Universities. *Procedia - Social and Behavioral Sciences*. (106), pp 476-484.
- Mcbee, M., & Barnes, L. (1998). The generalizability of a performance assessment measuring achievement in eight-grade mathematics. *Applied measurement in educations*, 11(2), 179-194.
- Meyer, J. P. (2010). *Reliability*. New York: Oxford University Press.
- Mohamed, A. A. (2010). Teaching licensure as a requirement to the teacher's performance quality. *مجلة كلية التربية جامعة الأزهر*, 144(2), pp 585- 606.
- Nicholas, T. (2016). *Generalizability of universal screening measures for behavioral and emotional risk. A Dissertation of Doctor of Philosophy*, university of Arizona, ProQuest.
- Nie, Y., Yeo, S., & Lau, S. (2007). Application of generalizability theory in the investigation of the quality of journal writing in mathematics. *Studies in Educational Evaluation*, 33(3-4), 371-383.
- Parkes, J. (2001). The role of transfer in the variability of performance assessment scores. *Educational Assessment*. 7(2), 143- 164.
- Ruiz-Primo, M., Baxter, G., & Shavelson, R. (1993). On the stability of performance assessment. *Journal of educational measurement*, 30(1), 41-53.
- Salkind, N. J. (2008). *Encyclopedia of educational psychology*. Los Angeles: Sage publications.
- Scallon, G. (2004). *The evaluation of learning in a competency-based approach*. Brussels: De Boeck.
- Shavelson, R., Baxter, G., & Gao, X. (1993). Sampling variability of performance assessments. *Journal of Educational Measurement*, 30(3), 215- 232.
- Shavelson, R. & Webb, N. (2009). Generalizability Theory and its contribution to the discussion of the generalizability of research finding. In K. Erichan., & W. Roth (Eds.), *Generalizability from Educational Research* (pp13-32). New York: Routledge.
- Shavelson, R.J., & Webb, N.M. (1991). *Generalizability theory: A primer*. Newbury Park, CA: Sage.

- Smith, E., & Kulikowich, j. (2004). An application of generalizability theory and many-facet measurement using a complex problem-solving skills. *Educational and Psychological Measurement*, 64(4), 617- 639.
- Suen, H.K. & Lei, P.W. (2007). Classical versus Generalizability theory of measurement. *Educational Measurement*, (4), 1-13.
- Taylor, A., & Pastor, D. (2013). An application of generalizability theory to evaluate the technical quality of an alternate assessment. *Applied Measurement in Education*, 26(4), 279–297.
- Urbina, S. (2004). *Essentials of psychological testing*. New Jersey: John Wiley & sons.
- Webb, N., Shavelson, R. & Haertel, H. (2006). Reliability coefficient and generalizability theory. *Handbook of statistics*, 26, 4-44.
- Webb, N., Schlackman, J., & Sugrue, B. (2000). The dependability and interchangeability of assessment methods in science. *Applied Measurement in Education*, 13(3), 277–301.
- Yang, Y., & Green, S. (2011). Coefficient alpha: A reliability coefficient for the 21<sup>st</sup> century?. *Journal of Psychoeducational Assessment*, 29, 377- 392.
- Yılmaz, N., Başbaşa, N. (2015). Assessment of sewing and picking skills station reliability with generalizability theory. *Journal of Measurement and Evaluation in Education and Psychology*, 6(1), 107-116.
- Yılmaz, N., Gelbal, S. (2011). Comparison of different patterns with generality theory in the case of communication skills station. *Hacettepe University Journal of Education Faculty*, 41, 509-518.