

# 学習者評価の信頼性とは

大西 弘高

東京大学医学系研究科医学教育国際研究センター

## 抄録

評価手法が良いか否かを議論するためには妥当性が重要だが、総括評価に客観試験以外の評価手法を用いる際、信頼性検証が不可欠である。OSCE (objective structured clinical examination) における信頼性指標は、各課題に挙げられている項目毎の評価者間信頼性と、ステーション毎にまとめた項目合計点をステーション間でみたときの信頼性に大別される。評価者間の信頼性指標には、評価者間一致率、 $\kappa$  値、級内相関係数が挙げられる。ステーション間の信頼性は  $\alpha$  係数で表すことが可能だが、一般化可能性理論を用いることで、誤差要因同士の分散の比較、一般化可能性係数や信頼度指数の計算、決定解析が可能となる。

**キーワード**：学習者評価，信頼性，テスト理論，一般化可能性理論

## Reliability in learner assessment

Hiroataka Onishi

International Research Center for Medical Education, Graduate School of Medicine, The University of Tokyo

### Abstract

Validity is important to discuss whether the assessment tool is good or not, evaluation of the validity is important. However, identifying reliability is indispensable when non-objective assessment tools are used for summative purposes. Two types of reliability indices are used for OSCE (objective structured clinical examination); one is the inter-rater reliability for each item of the checklist in a station task, and the other is inter-station reliability using the total item scores from all the stations. Indices of inter-rater reliability consist of inter-rater agreement, kappa value, and intra-class correlation coefficients. Inter-station reliability can be expressed by the alpha coefficient. Utilization of the generalizability theory in OSCE enables comparison of variances among different error factors, calculation of the generalizability and dependability indices, and decision study.

**Key words** : learner assessment, reliability, test theory, generalizability theory

## はじめに

学習者を評価する際、精神運動領域、情意領域、高次認知領域の評価等においては評価者のバラツキの懸念が大きくなる。総括評価、特に high-stake test に

おいては、一定以上の信頼性が達成されていなければ、合否判定等に関する誤差が大きくなるため、信頼性を検証し、改善することは不可欠といえるだろう。

本稿では、OSCE による評価をイメージしながら信頼性指標の持つ意味、様々な信頼性指標の意義に分け

【連絡先】東京大学医学系研究科医学教育国際研究センター  
〒113-0033 東京都文京区本郷7-3-1 医学部総合中央館2階  
受理日：2018年2月9日 採録決定日：2018年3月6日

表1 妥当性を証明するための情報源

内容	テストで得たスコアの処理プロセス	テストの内的構造	他の変数との関係	結果
<ul style="list-style-type: none"> <li>試験のブループリント, 達成すべき目標領域, テスト内容が一貫している</li> <li>テストが目標領域を代表している</li> <li>テストで問われている内容が現場で必要となる内容と一致している</li> <li>テストの問題の質が高い</li> <li>問題作成者が一定レベルの質に達している</li> </ul>	<ul style="list-style-type: none"> <li>学習者がテストのフォーマットに慣れているか</li> <li>採点が正確かどうかの確認</li> <li>異なったフォーマットの点数を正確に合計しているか</li> <li>スコアや評価の正確性に関する質管理</li> <li>スコアと合否判定に関する質管理</li> <li>学習者や教育者への報告に関する質管理</li> <li>スコアの記述が分かりやすく, 正確か</li> </ul>	<ul style="list-style-type: none"> <li>項目分析データ (項目の難易度, 識別係数, 項目特性曲線, 項目間相関, 項目と合計点との相関)</li> <li>項目の信頼性, 測定の標準誤差, 一般化可能性分析, 因子分析</li> <li>心理測定モデル</li> </ul>	<ul style="list-style-type: none"> <li>他の関連した変数との相関</li> <li>類似したテストとの収束的相関 (内的, 外的)</li> <li>異なったテストとの弁別的相関 (内的, 外的)</li> <li>テストと規準 (クライテリア) との相関</li> <li>エビデンスの一般化可能性 (外挿可能性)</li> </ul>	<ul style="list-style-type: none"> <li>テスト結果が学習者や社会に与える影響</li> <li>将来の学習者に与える影響</li> <li>正の結果が予測しなかった負の結果を上回っているか</li> <li>合否基準となるスコアを決定する方法が納得できるものか</li> <li>合否決定による分類の正確性</li> <li>能力が不十分なのに合格になったり, 能力が十分なのに不合格になったりする者の数</li> <li>合否判定基準スコアの条件付き標準誤差</li> <li>評価が学習者に与える影響</li> </ul>

て概説し, 例を挙げてその有用性にも触れてみたい。

### 信頼性指標の意味

評価には常に信頼性, 妥当性の問題が付いて回る。妥当性とは, その評価手法が真に測定したいものを測定できているかどうか, 信頼性とは, その評価手法の測定が安定しているかと定義される。学習者評価が学習者の能力を反映しているかどうかは, 妥当性で判断されるべきであり, 信頼性の問題はさほど大きくないはずである。

妥当性は, Messick の再定義により, 構成概念妥当性に一元化されると共に内容がやや拡がり, テストデータの処理, 結果が社会に与えるインパクト等も考慮することになった<sup>1)</sup>。Downing はこの新しい定義に基づき, 妥当性を証明するための情報源を表1のように列挙した<sup>2)</sup>。特に, 合否判定に関する内容は以前よりもはるかに重視されるに至った。

さて, 再度各学生の評点の持つ意味を考えてみよう。学生の評点が60点であった場合, これはいくらかの誤差を含んだ点推定値であるとみなすことも可能である。ここで, 測定値の標準誤差 (standard error of measurement) が5点で, 評点が正規分布に従うと仮定すれば, [50.2 (60-5 × 1.96) ~ 69.8 (60+5 × 1.96)]

が点推定値の95%信頼区間となる。体重計に乗って60kgを指したとき, その95%信頼区間が50~70kgであったなら, その体重計は全く信頼できないと呆れるだろうが, 学習者評価の評点というのはそのような性質の数値であると認識しておくべきだろう。

ところが, 合否を判定しようとすれば, 60点の学生が合格で59点の学生が不合格, ということも起こりうる。このとき, 点推定値としての評点の信頼区間があまりに広いと, 特にぎりぎりでは不合格になった受験者から試験に疑義を挟まれる可能性はある。そのような問題に可能な限り対処するために, 学習者評価の信頼性に配慮が必要なのである。

まずは, 信頼性指標とは何かを明確にしていこう。概念的には, 同じ人が違うタイミング, 違う視点で何回も評点を付ける, あるいは異なる人が評点を付けるといった場合に, その評点がどの程度安定しているかを意味する。

$$\text{信頼性係数} = \frac{\text{真値の分散}}{\text{真値の分散} + \text{誤差分散}} \dots \text{式①}$$

式①より, 誤差分散が小さくなると信頼性係数が大きくなる関係があることが分かる。各施設レベルの試験なら0.7, 全国レベルの試験なら0.8という信頼性係数が目安として示されている<sup>3)</sup>。

表2 信頼性指標の種類

1. 評価者間信頼性指標
(ア) 項目毎の信頼性指標
① 評価者間一致率
② $\kappa$ 値
③ 級内相関係数 (通常は ICC (2, 1))
(イ) ステーションの項目合計点を用いた信頼性指標
① Pearson の相関係数
② Spearman の相関係数
2. ステーション間の信頼性指標
(ア) Cronbach の $\alpha$ 係数
3. 評価者とステーションの要因を併せ持つ信頼性指標
(ア) 一般化可能性係数 (generalizability index, G 係数)
(イ) 信頼度指数 (dependability index, $\Phi$ 係数)

表3 2名の評価者によるチェックリスト評点の例

	評価者 2			計
	0	1		
評価者 0	0	0	0	0
1	1	1	49	50
計	1	1	49	50

注1：観察された一致率 =  $P_o$ 、一致率の期待値 =  $P_e$  とすると、 $P_o = 49/50 = 0.98$ 、 $P_e = 49/50 = 0.98$  であり、 $\kappa$  値 =  $(P_o - P_e) / (1 - P_e) = 0$

注2：被評価者数  $n$ 、評価者数  $k$  と置き、BMS (between-subject mean squares), JMS (between measures (judges) mean squares), EMS (error mean squares) を用いて表すと、 $ICC(2, 1) = (BMS - EMS) / \{BMS + (k - 1) EMS + k (JMS - EMS) / n\}$  となる。文献<sup>4)</sup>の定義によって計算すると、 $BMS = 0.01$ 、 $EMS = 0.01$  であり、 $ICC(2, 1) = 0$  となる。

## 信頼性指標の種類やその意義

様々な信頼性指標を表2に一括して挙げた。このうち、評価者間の信頼性指標については、項目毎のものとステーションの項目合計点を用いたものに分けられる。ステーション間の信頼性指標を求める際には各ステーションの項目合計点が用いられる<sup>3)</sup>。

このうち、利用の意義が大きいのは項目毎の評価者間信頼性指標と信頼度指数 (dependability index) である。これらについて、下に詳述する。

### 項目毎の評価者間信頼性指標

項目毎の評価者間信頼性指標は、各項目がどのような性質を持ち、評価基準が明確になっているかどうかを確認するのに役立つことがある。例えば、医療面接ステーションで「患者の姓を確認した」という項目は評価者間のズレが生じにくく、「患者の話を共感的に聴いた」という項目はズレが生じやすいという違いが

明確になるかもしれない。後者に対し「患者の話を聴きながら5回以上相づちを打った」という評価基準を設ければ評価者間信頼性は向上するかもしれないが、マニュアル的な対応を増やして妥当性を減じるかもしれない。

$\kappa$  値や級内相関係数は項目の得点率が50%付近では信頼区間が小さいが、0%や100%に近づくと信頼区間が大きくなるため注意が必要である。例えば表3では「患者に挨拶した」という項目に対し、ほぼ全員が正しく実施できているが、ただ1人の学生に対する評価がずれている。このとき、評価者間一致率は0.98だが、 $\kappa$  値、級内相関係数 (ここでは評価者間一致度を表すときに頻用される  $ICC(2, 1)$ <sup>4)</sup> を用いる) のいずれもが0となる。もし、この評価がずれた学生の真の評価は0点であって、評価者1がそれを正しく0点としていたとすれば、全ての指標はいずれも1となる。これらの観察から、評価者間一致率と  $\kappa$  値、あるいは級内相関係数の併記が奨められている<sup>5)</sup>。

なお、評価者が3名以上のときに評価者間信頼性指

表4 OSCE 結果による一般化可能性解析の例

要因	偏差平方和	自由度	平均平方	分散推定値	分散の比率
学生 (i)	4.4813	91	0.04924	0.004422	49.5%
ステーション (j)	0.6672	4	0.16680	0.000804	9.0%
評価者 (k:j)	0.0801	5	0.01602	0.000151	1.7%
学生×ステーション (ij)	1.8273	364	0.00502	0.001463	16.4%
学生×評価者 (i (k:j))	0.9527	455	0.00209	0.002094	23.4%
Sum	8.0086	919		0.008934	

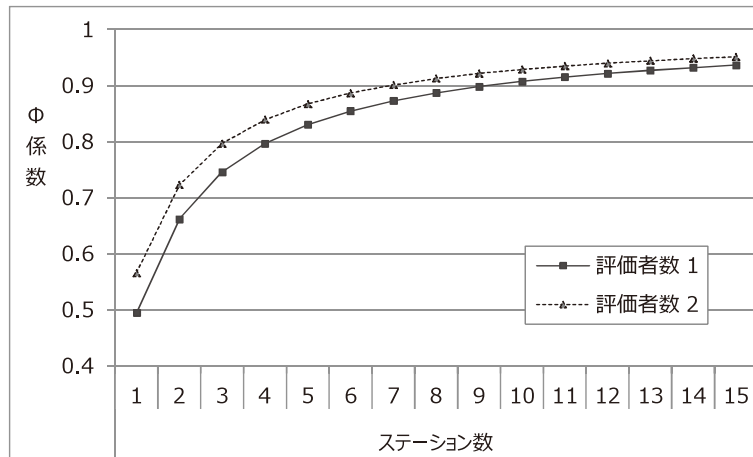


図1 表4と同じデータを用いた決定解析の例

標を求めようとするれば、級内相関係数 (intraclass correlation coefficient : ICC) が最も簡便に求められる。ICC にはどの分散を用いるかによって6種類のデータが得られるが、複数評価者を変量モデルとした場合には ICC (2, 1) が用いられる<sup>4)</sup>。

$\kappa$  値は、どの程度の数値なら評価者間信頼性が保たれているかに関するガイドラインがあることでも知られる。Landis and Koch は、 $\kappa \leq 0.2$  なら Poor,  $0.2 < \kappa \leq 0.4$  なら Fair,  $0.4 < \kappa \leq 0.6$  なら Moderate,  $0.6 < \kappa \leq 0.8$  なら Good,  $0.8 < \kappa \leq 1.0$  なら Very good とした<sup>6)</sup>。Fleiss は、 $0.75 \leq \kappa$  なら Excellent,  $0.40 < \kappa < 0.75$  なら Fair to good,  $\kappa \leq 0.40$  なら poor とした<sup>7)</sup>。このガイドラインは ICC にも用いられることがある。しかし、これらの基準には数値的な根拠は特になく、一部の統計学者はこのガイドラインの使用を憂慮している<sup>8)</sup>。

評価者間信頼性指標のうち、ステーションの項目合計点を用いたものはさほど意味を持たない。相関係数が低くても、どこを改善すべきかには項目毎の指標が必要になるからである。ステーション間の信頼性指標を求めれば、ステーション数を増やしたときにどの程度信頼性指標が上がるかが分かるため、参考になる。このような分析は  $\alpha$  係数、信頼度指数のいずれでも可

能だが、評価の主観性が問題になりがちな OSCE のような場合には、これも含めて解析できる後者がより有用であろう。

### 一般化可能性解析の概要

一般化可能性解析では各ステーションの項目合計点を用い、分散分析の手法によって式①における誤差分散を3つの要因 (学生の評点の差によるもの (i), ステーション (j), 評価者 ( $\kappa$ )) で分離することが主眼となる。その際、各要因の交互作用も得られるが、評価者はステーション毎に異なるため、評価者とステーションの交互作用を求めることはできない (要因のネスト構造と呼ばれ、評価者の分散は  $\sigma_{(\kappa:j)}$  と表記される)。よって、学生、ステーション、評価者の分散と、学生とステーションの交互作用の分散、学生と評価者の交互作用の分散の5つが分離されることになる。

信頼度指数  $\Phi$  は、学生の分散と、それ以外の分散成分から、式②の形

$$\Phi = \frac{\hat{\sigma}_i^2}{\hat{\sigma}_i^2 + \frac{\hat{\sigma}_{ij}^2 + \hat{\sigma}_j^2}{n} + \frac{\hat{\sigma}_{i(\kappa:j)}^2 + \hat{\sigma}_{(\kappa:j)}^2}{nr}} \dots \text{式②}$$

( $n, r$  はステーション数及び評価者数) で求まる。なお、

式②において学生要因である  $i$  を含まない項を削除したものが一般化可能性係数 (generalizability index) となる。一般化可能性係数は相対評価、信頼度指数は絶対評価に利用すべきとされる。詳細については、文献を参照していただきたい<sup>9-11)</sup>。

## OSCE 評点を用いた一般化可能性解析の例

表4には、受験者92名、5ステーション、各ステーション評価者2名の形で実施されたOSCEの評点を用いた一般化可能性解析 (generalizability study) の例を挙げた。ここには、5種類の分散推定値やその割合が示されている。絶対評価に用いることを想定して式②に代入すると、 $\Phi = 0.004422 / (0.004422 + 0.000804 / 5 + 0.001463 / 5 + 0.000151 / 10 + 0.002094 / 10) = 0.867$ となる。図1には、これと同じデータを用いた決定解析 (decision study) のグラフを示した。決定解析はステーション数や評価者数を変更したときに、信頼度指数がどう変化するかを予測するものである。例えば、ステーション数7、評価者数2とすれば、 $\Phi = 0.004422 / (0.004422 + 0.000804 / 7 + 0.001463 / 7 + 0.000151 / 14 + 0.002094 / 14) = 0.901$ となる。

## まとめ

評価の信頼性に関する一般的な内容を、OSCEの例を中心にまとめた。特に、項目毎の評価者間信頼性指

標、複数の相関要因を含めた一般化可能性解析については、OSCEの信頼性向上に対して有用であり、詳述した。

## 文献

- 1) Messick, S. (1989). Validity. In : Educational Measurement, 3rd edn. Ed : Linn RL. New York : American Council on Education and Macmillan, 13-104.
- 2) Downing, SM. (2003). Validity : on the meaningful interpretation of assessment data. Med Educ, 37 : 830-837.
- 3) Downing, SM. (2004). Reliability : on the reproducibility of assessment data. Med Educ, 38 : 1006-1012.
- 4) Hasnain, M, Onishi, H, Elstein, A.S. (2004). Inter-rater agreement in judging errors in diagnostic reasoning. Med Educ, 38 : 609-616.
- 5) Shrout, P.E, Fleiss, J.L. (1979). Intraclass Correlations : Uses in Assessing Rater Reliability. Psychol Bull, 86 : 420-428.
- 6) Landis, J.R, Koch, G. (1977). The measurement of observer agreement for categorical data. Biometrics, 33 : 159-174.
- 7) Fleiss, J.L. (1981). Statistical Methods for Rates and Proportions, 2nd edn. New York : Wiley, : 212-236.
- 8) Uebersax, J. (2002). Kappa Coefficients. <http://our-world.compuserve.com/homepages/jsuebersax/kappa.htm> (accessed 11 October 2006).
- 9) Shavelson, R. (1991). Webb, N. Generalizability Theory ; A Primer. Sage Publications, Thousand Oaks.
- 10) 池田央. (1994). 現代テスト理論, 東京 : 朝倉書店.
- 11) Brennan, R.L. (2005). Generalizability theory, New York : Springer.