# Cambridge English

## Multiple Standard Setting Outcomes: Empirical Exercises Informing Theory

## How Should We Think About Divergent Standard Setting Outcomes?

It has been observed that standard setting sometimes results in different outcomes. How should we think about such results?

Standard setting theory argues that standard setting is not an exercise in homing in on a "true", objectively-existing cut score, but a values-driven enterprise arrived at for some policy-related end (Camilli, Cizek & Lugg, 2001; Cizek, 1993; Kane, 1998; Zieky, 2001), so disparate outcomes should be acceptable.

However, standard setting is sometimes concerned with finding levels on standards, frameworks and benchmarks with objective criteria, e.g. the Common European Framework of Reference, the Canadian Language Benchmarks (Figures 1 and 2). Under these circumstances, standard setting should be seen as expert informed but data constrained exercises where convergent outcomes should be expected.



## An Analogy

If the question about hotels is about "good enough", people can be expected to disagree and that is all right. On the other hand, if the question is about the star rating for a hotel, the objective criteria limit the correct answer that the informant can arrive at.

| Is this hotel good enough for us to stay in on our trip? | What is the class / star rating for this hotel? |
|---|---|



| Stars | Class | Excerpt of Criteria |
|---|---|---|
| ★ | Tourist | • All rooms with table and chair<br>• All rooms with colour TV and remote control<br>• All rooms with shower/WC or bath tub/WC |
| ★★ | Standard | In addition to ★:<br>• Linen shelves<br>• Bath essence of shower gel<br>• Bath towels<br>• Credit cards accepted |
| ★★★ | Comfort | In addition to ★★:<br>• Place to put luggage/suitcase<br>• Dressing mirror<br>• Telephone in room |
| ★★★★ | First Class | In addition to ★★★:<br>• Upholstered chair/couch with side table<br>• Minibar or 24 hour beverages via room service<br>• Cosmetic products, bath robe and slippers on demand |
| ★★★★★ | Luxury | In addition to ★★★★:<br>• Fresh flowers or a present in the room<br>• Internet PC in the room<br>• Safe in the room<br>• Turndown service |

## The Argument

When multiple standard setting exercises are conducted in relation to criterion standards, convergent outcomes should be expected and divergent outcomes investigated. This makes standard setting outcomes subject to validation/falsification, making them more useful, and leading to better supported validity arguments. They should be seen as repeated measures (with a certain measurement error attached to them) which together leads to a more accurate outcome.

## Empirical Illustration 1

Standard setting exercises were conducted in 2005 and 2013 relating IELTS, a test of English language proficiency, to the Canadian Language Benchmarks. The representative excerpt in Table 1 shows that results generally converge, and where they do not, the standard error informs what the correct cut score should be.

*Table 1. IELTS Band Scores Relating to CLB Benchmark 7*

| | | – 1 SE | Ave | + 1 SE | | | – 1 SE | Ave | + 1 SE |
|---|---|---|---|---|---|---|---|---|---|
| Listening | 2005 | 6.0 | 6.0 | 6.5 | Speaking | 2005 | 5.5 | 6.0 | 6.5 |
| | 2013 | 5.5 | 6.0 | 6.0 | | 2013 | 6.0 | 6.0 | 6.5 |
| Reading | 2005 | 5.5 | 6.0 | 6.5 | Writing | 2005 | 5.5 | 6.0 | 6.5 |
| | 2013 | 5.5 | 5.5 | 6.0 | | 2013 | 5.5 | 6.0 | 6.0 |

## Empirical Illustration 2

Table 2 shows four studies (A: Pearson, 2009; B: Tannenbaum & Wiley, 2004, ETS, 2006, 2010; C: Buckendahl & Khalifa, 2009; D: Lim, Geranpayeh, Khalifa & Buckendahl, 2013) relating IELTS to level C1 on the Common European Framework of Reference appear to be divergent, making them unhelpful.

However, further thought about the studies reveal the following:

- Study A related IELTS to the CEFR via another test, asking test takers to sit both exams. But data did not account for test-takers' motivation, (participants took IELTS test for high-stakes purpose, took non-bearing field test of other exam), providing an inaccurate picture of the relationship between the two exams.

- Study C was influenced by standard setting method effects (Jaeger, 1989). The Analytic Judgment method (Plake & Hambleton, 2000) was used for the Writing and Speaking tests, where panellists select a performance deemed at a level, and the score for that is the cut score. IELTS reports wide band scores with half bands (4.0, 4.5, 5.0, 5.5…). Thus, if the true cut point is, say, 3.13, panellists would logically choose a 3.5 response, inflating the cut score. The Yes-No Tucker-Angoff (Angoff, 1971; Impara & Plake, 1997) was used for the discrete-point Reading and Listening tests, which produces overestimates at higher levels where items are easy for the target test taker (Cizek & Bunch, 2007; Council of Europe, 2009).

- After divergent outcomes are investigated and cut scores validated, a clearer picture emerges (Table 3), and results become much more useful.

*Table 2. IELTS Band Scores Relating to CEFR Level C1, Before Validation*

| Study | C1 Cut Score |
|---|---|
| A | 7.5 |
| B | ~6.6 |
| C | ~7.3 |
| D | ~6.6 |

*Table 3. IELTS Band Scores Relating to CEFR Level C1, After Validation*

| Study | C1 Cut Score |
|---|---|
| A | 7.5 |
| B | ~6.6 |
| C | ~6.8 |
| D | ~6.6 |

## Conclusion

Situations exist where multiple standard setting exercises should produce convergent outcomes. Standard setting theory needs to account for these situations involving objective criteria, seeing them as expert informed but data constrained exercises subject to validation and falsification, leading to better supported validity arguments. Examples from empirical studies illustrated approaches for validating cut scores.

## Authors

Gad S Lim, Cambridge English Language Assessment

Chad W Buckendahl, Alpine Testing Solutions

**CAMBRIDGE ENGLISH** Language Assessment Part of the University of Cambridge

**Alpine** Testing Solutions