# Items and Item Writing

PSY 600K
Dr. Alyssa Gibbons
January 30, 2018

# Agenda

O Ways to approach scale development.

O Key choices in the development process.

O Redundancy.

O Characteristics of good items.

O Cognitive process of item responses.

# Scale Development

0 Burisch (1984) identifies 3 possible approaches:
  0 External approach
    0 Find items that differentiate between people who differ on the trait you are interested in.
    0 Similar: criterion-referenced approach – find items that are correlated with the thing you want to predict.
    0 Explicitly atheoretical.
  0 Inductive approach
    0 Start with a bunch of items.
    0 Factor analyze them ("matrix storing").
    0 Interpret the results as describing the true structure of the construct.

# Scale Development

0 Deductive approach:
  0 Start with a theory (elaborate or common-sense).
  0 Use theory to decide how many factors there should be.
  0 Write items that relate to these factors.

0 Which approach does Burisch (1984) clearly prefer? Why?

0 Is it possible to develop a test without a theory?

# Be Intentional

0 If you have no idea whatsoever about the probable structure of your measure, you're probably not ready to write it.
  0 DeVellis: "The point is that scale developers should make this determination as an active decision and not merely generate a set of items and then see what they look like after the fact." (p. 75)
  0 If there's not much theory, you need to start on one.
0 Know the literature.
  0 Those who do not remember the past are destined to republish it.
  0 Learn from others' mistakes.

# Key Choices

0 Level of specificity:
  0 Many constructs can be unidimensional or multidimensional… depending on your purpose.
  0 Example: How many factors of intelligence are there?
    0 Spearman: 2
    0 Thurstone: 7
    0 Guilford: 150!
      Carroll: 9, in a hierarchical model
0 Should we be trying to "discover" the "true" structure of our measures?

# Key Choices

0 What is and is not relevant?
   0 Important to identify what does **not** belong in your scale.
   0 DeVellis example: depression & physical health.
   0 Other examples?
0 Purpose
   0 What do you want to be able to say about people after they complete your measure?
   0 What part of the latent continuum are you most interested in?
0 Population
   0 Who will take the test? Under what circumstances?

# Every Test Needs:

0 Item **stems:** the stimuli a test-taker responds to.
0 Response options: what the test-taker can do in response to an item stem.
   0 Minimum: 2
   0 Maximum: infinity (or constrained by the scale developer).
0 A method for assigning item scores (numbers) to those responses.
   0 **Objective** test: means that human judgment is not needed at the *scoring* stage (there is plenty of judgment at other stages!)
   0 Of course, not all tests are objective.
0 A method for combining the item scores into test scores.

# Terminology

0 **Objective** test: means that human judgment is not needed at the scoring stage.
  - 0 Straightforward link from responses to scores.
  - 0 Lots of judgment needed in **other** stages…
  - 0 We're mostly going to focus on this type.
- 0 **Selected-response** vs. **constructed-response**.
  - 0 How much structure is provided to tell the test-taker how to respond?

# Kinds of Items

0 Depends (as always) on your purpose – what sort of information do you want to end up with?
- 0 Key distinction:
  - 0 **Dichotomous** item responses are scored 0 or 1.
    - 0 Correct or incorrect, present or absent, yes or no.
  - 0 **Continuous** item responses have a wider range of possible scores.
- 0 You can't always tell whether an item is scored dichotomously or continuously just by looking at it.

# Redundancy

0 Do we really need to ask the same thing over and over?

0 Depends on what we mean by "same thing."

0 **Useful** redundancy is repeating the same *idea* in a different way.

- 0 Item uniqueness - idiosyncrasies of the items – quirks of wording, interpretation, etc.
- 0 Across several items, these cancel out – common variance dominates unique variance.
- 0 Allows us to capture the construct more fully.

# Redundancy

0 **Useless** redundancy is repeating the same idea in pretty much the same way.

- 0 Example: "I like social occasions" and "I enjoy social occasions."
- 0 Cannot tell whether common variance is really due to the idea or the wording.

0 When 2 items have "something extra in common" with one another, over and above the common construct, they wreak havoc on factor structure.

- 0 Example: DeVellis' African grey parrot items.
- 0 **Doublet** factor (or, if more than 2 items, a subfactor).

# How Many Items?

0 Absolute minimum: 3 per subscale
  0 Mathematical necessity to have a stable factor structure.
0 Having more items:
  0 Increases reliability (not always for the right reasons).
  0 Allows more thorough coverage of the construct.
0 Having fewer items:
  0 Is usually practical .
  0 Prevents respondent fatigue.
0 Your **initial** item pool should be 2x – 4x the number of items you hope to end up with.

# Characteristics of Good Items

0 #1. Short.
  0 The more words in an item, the more opportunities there are for a respondent to misread or misinterpret one.
0 #2. Readable.
  0 Can evaluate reading level if you like.
  0 Use simple, everyday language.
  0 Avoid jargon.
    0 Especially if you don't understand it yourself.

1/15/2018

# Good Items

- #3. Grammatically correct.
  - Avoid double negatives.
    - Or perhaps avoid negative words altogether.
  - Use adverbs & adjectives correctly.
  - Read for unintended alternative interpretations.
  - Complete sentences are often best.
- #4. Consistent.
  - Use the same referent ("I", "you", etc.) and general structure in all items.
  - Not:
    - "I am a warm and outgoing person."
    - "You like parties."
    - "Sociable."

---

# Good Items

- #5. Straightforward
  - Interpreting the question may be obvious to you, but is it obvious to the test-taker?
  - **Double-barreled** items: two questions in one.
    - "I believe graduate students are underpaid and overworked."
    - Does an "agree" response you agree with both parts of the statement? Or only one?
  - Clear relationship to the construct.
    - "I am confident that I am ready to be a parent."
    - Two possible response processes here...

# Communicability

o Burisch's idea: to what degree do the items tell us something clear about a person?
  o And is it the something we want to know?
o Burisch notes that these don't guarantee honesty.
  o Not really an issue of transparency vs. subtlety to the *respondent*.
o "Defining" vs. "correlating" characteristics.
  o Items that hit the center of the construct vs. the outside.
  o Also called "prototypical" items.

# The Law of Simplicity

o Burisch also points out that (in personality) simple trait rating scales perform as well or better than more complex instruments.
  o "I am outgoing" vs "I enjoy interacting with other people."
o Simple formats also tend to perform as well as or better than more complex ones.
o Do you think this is true?
  o Does it hold across other areas of psychology?
  o Why?

# Cognitive Psychology of Item Responses

0 Test-takers go through a 4-stage process:
  0 1. Comprehension – what is this question asking me?
  0 2. Retrieval – thinking about relevant information – facts, behaviors, etc.
  0 3. Judgment – choosing which response is most appropriate.
  0 4. Response communication – recording or conveying that response.
0 Error could occur at any step!
0 Think through your items as a respondent would – troubleshoot.

# Good Process

0 Think carefully through your theory.
  0 Including expected dimensionality, even if you are "exploring."
0 Write a good definition & get feedback on it.
0 Write more items than you expect to need.
  0 4-5 per dimension at a minimum.
0 Proofread.
0 Proofread again.
0 Pilot test and invite comments from respondents.

# Evidence of Appropriate Content

0 Appropriate content = do your items represent the whole domain of the construct you want to cover?

  0 Is anything included that shouldn't be?

  0 Is anything missing that should be included?

  0 Is there balance across all the important aspects of the construct?

0 We need to **support** our argument here.

  0 Expert judgment – SMEs.

  0 We're going to do this informally, but we can do it quite formally indeed (more on this later).

# Expert Judgment

0 Ask your SMEs:

  0 Is my definition appropriate for this construct?

  0 Do these items fit with my definition?

    0 Is anything here irrelevant? Only a little relevant?

    0 Is each item **essential** for measuring this construct?

  0 Is anything missing?

  0 Can these items be clarified or revised?

# For Lab on Friday

0 Please write 5 items measuring the construct **"Satisfaction with Graduate School"**
  0 Definition: "cognitive and affective evaluations of one's graduate education."
  0 Purpose: research.
0 Use any item style and response format you wish.
0 Submit your items on Canvas by Thursday night. I'll compile them so we can critique them anonymously in lab.
0 Try to write some good items and some bad items.

# Questions?

Project Plans Due Thursday!

For next time:
Response Formats & Scales
Read: DeVellis pp. 85 – 104