

A practical primer on processing semantic property norm data

Erin M. Buchanan¹, Simon De Deyne², & Maria Montefinese³

¹ Harrisburg University of Science and Technology

² University of Melbourne

³ University of Padua

Author Note

Add complete departmental affiliations for each author here. Each new line herein must be indented, like this line.

Enter author note here.

Correspondence concerning this article should be addressed to Erin M. Buchanan, 326 Market St., Harrisburg, PA 17101. E-mail: ebuchanan@harrisburgu.edu

Abstract

Semantic property listing tasks require participants to generate short propositions (e.g., <barks>, <has fur>) for a specific concept (e.g., dog). This task is the cornerstone of the creation of semantic property norms which are essential for modelling, stimuli creation, and understanding similarity between concepts. However, despite the wide applicability of semantic property norms for a large variety of concepts across different groups of people, the methodological aspects of the property listing task have received less attention, even though the procedure and processing of the data can substantially affect the nature and quality of the measures derived from them. The goal of this paper is to provide a practical primer on how to collect and process semantic property norms. We will discuss the key methods to elicit semantic properties and compare different methods to derive meaningful representations from them. This will cover the role of instructions and test context, property pre-processing (e.g., lemmatization), property weighting, and relationship encoding using ontologies. With these choices in mind, we propose and demonstrate a processing pipeline that transparently documents these steps resulting in improved comparability across different studies. The impact of these choices will be demonstrated using intrinsic (e.g. reliability, number of properties) and extrinsic measures (e.g., categorization, semantic similarity, lexical processing). Example data and the impact of choice decisions will be provided. This practical primer will offer potential solutions to several longstanding problems and allow researchers to develop new property listing norms overcoming the constraints of previous studies.

Keywords: semantic, property norm task, tutorial

A practical primer on processing semantic property norm data

1. Available feature norms and their format

- Property listing task original work: (???) (???) (???) (???)
- English: (???) (???) (???) (???) (???)
- Italian: (???) (???) (???)
- German: (???)
- Portuguese: (???)
- Spanish: (???)
- Dutch: (???)
- Blind participants: (???)

I'm sure there are more, here's what we cited recently.

Define concept, feature for clarity throughout - make sure you use these two terms consistently.

2. Pointers about how to collect the data

- a. instructions, generation, verification, importance

I really like the way the CSLB did it: <https://cslb.psychol.cam.ac.uk/propnorms>

They showed the concept, then had a drop down menu for is/has/does, and then the participant typed in a final window. That type of system would solve about half the problems I am going to describe below about using multi-word sequences. Might be some other suggestions, but for that type of processing, you could do combinations and have more consistent data easily.

3. Typical operations performed on features

55 Materials and Data Format

56 The data for this tutorial includes 9553 unique concept-feature responses for 104
 57 concepts from (???) that were included in (???), (???), and (???). The data should be
 58 structured in tidy format wherein each concept-feature observation is a row and each column
 59 is a variable (???). Therefore, the data includes a **word** column with the normed concept
 60 and an **answer** column with the participant answer.

word	answer
airplane	you fly in it its big it is fast they are expensive they are at an airport you have to be trained to fly it there are lots of seats they get very high up
airplane	wings engine pilot cockpit tail
airplane	wings it flies modern technology has passengers requires a pilot can be dangerous runs on gas used for travel
airplane	wings flies pilot cockpit uses gas faster travel
airplane	wings engines passengers pilot(s) vary in size and color
airplane	wings body flies travel

61 This data was collected using the instructions provided by (???), however, in contrast
 62 to the suggestions for consistency detailed above (???), each participant was simply given a
 63 large text box to include their answer. Each answer includes multiple embedded features, and
 64 the tutorial proceeds to demonstrate potential processing addressing the data in this nature.
 65 With structured data entry for participants, the suggested processing steps are reduced.

66 Spelling

67 Spell checking can be automated with the `hunspell` package in *R* (???), which is the
68 spell checking library used in popular programs such as FireFox, Chrome, RStudio, and
69 OpenOffice. Each `answer` can be checked for misspellings across an entire column of answers,
70 which is located in the `master` dataset. The default dictionary is American English, and the
71 `hunspell` vignettes provide details on how to import your own dictionary for non-English
72 languages.

```
## Install the hunspell package if necessary
#install.packages("hunspell")
library(hunspell)

## Check the participant answers
## The output is a list of spelling errors for each line
spelling_errors <- hunspell(master$answer, dict = dictionary("en_US"))
```

73 The result from the `hunspell()` function is a list object of spelling errors for each row
74 of data. For example, when responding to *apple*, a participant wrote *fruit grocery store*
75 *orchard red green yellowe good with peanut butter good with caramell*, and the spelling errors
76 were denoted as *yellowe caramell*. After checking for errors, the `hunspell_suggest()`
77 function was used to determine the most likely replacement for each error.

```
## Check for suggestions
spelling_suggest <- lapply(spelling_errors, hunspell_suggest)
```

78 For *yellowe*, both *yellow yell* were suggested, and *caramel caramels caramel l camellia*
79 *camel* were suggested for *caramell*. The suggestions are presented in most probable order,
80 and using a few loops with the `a substitute` function, we can replace all errors with the most

81 likely replacement in a new dataset `spell_checked`.

```
## Replace with most likely suggestion
spell_checked <- master

### Loop over the data.frame
for (i in 1:nrow(spell_checked)){
  ### See if there are spelling errors
  if (length(spelling_errors[[i]]) > 0) {
    ### Loop over all errors
    for (q in 1:length(spelling_errors[[i]])){
      ### Replace with the first answer
      spell_checked$answer[i] <- gsub(spelling_errors[[i]][q],
                                     spelling_suggest[[i]][[q]][1],
                                     spell_checked$answer[i])
    }
  }
}
```

82 a. stemming, spelling, normalization

83 b. exceptions dictionary

84 ii. bing spell checker service

85 b. Weighting

86 c. feature type ontologies

87 d. identify cut off for idiosyncratic features (should it be necessary?)

88 4. Specification of how this is automated (package description)

- 96 Discussion

References