

A practical primer on processing semantic property norm data

Erin M. Buchanan<sup>1</sup>, Simon De Deyne<sup>2</sup>, & Maria Montefinese<sup>3</sup>

<sup>1</sup> Harrisburg University of Science and Technology

<sup>2</sup> University of Melbourne

<sup>3</sup> 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. The choice of dictionary should also normalize between multiple varieties of the  
73 same language, for example, the `"en_GB"` would convert to British English spellings.

```
## 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"))
```

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

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

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

82 most likely replacement in a new dataset `spell_checked`. A specialized dictionary with  
83 precoded error responses and corrections could be implemented at this stage. Other paid  
84 alternatives, such as Bing Spell Check, can be a useful avenue for datasets that may contain  
85 brand names (i.e, *apple* versus *Apple*) or slang terms.

```
## 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])

    }

  }

}
```

## 86 Lemmatization

87 The next step approaches the clustering of word forms into their lemma or head word  
88 from a dictionary. The process of lemmatizing words involves using a lexeme set (i.e., all  
89 words forms that have the same meaning, *am*, *are*, *is*) to convert into a common lemma (i.e.,  
90 *be*) from a trained dictionary. In contrast, stemming involves processing words using  
91 heuristics to remove affixes or inflections, such as *ing* or *s*. The stem or root word may not

92 reflect an actual word in the language, as simply removing an affix does not necessarily  
93 produce the lemma. For example, in response to *airplane*, *flying* can be easily converted to  
94 *fly* by removing the *ing* inflection. However, this same heuristic converts the feature *wings*  
95 into *w* after removing both the *s* for a plural marker and the *ing* participle marker. Several  
96 packages for *R* include customizable stemmers, notably the `hunspell`, `corpus` (???), and `tm`  
97 (???) packages.

98 Lemmatization is the likely choice for processing property norms, and this process can  
99 be achieved by installing `TreeTagger` (???) and the `koRpus` package in *R* (???).  
100 `TreeTagger` is a trained tagger designed to annotate part of speech and lemma information in  
101 text, and parameter files are available for multiple languages. The `koRpus` package includes  
102 functionality to use `TreeTagger` in *R*. After installing the package and `TreeTagger`, we will  
103 create a unique set of tokenized words to lemmatize to speed computation.

```
lemmas <- spell_checked

## Install the koRpus package
#install.packages("koRpus")
#install.packages("koRpus.lang.en")

## You must load both packages separately
library(koRpus)
library(koRpus.lang.en)

## Install TreeTagger
#https://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/

## Find all types for faster lookup
all_answers <- tokenize(lemmas$answer, format = "obj", tag = F)
all_answers <- unique(all_answers)
```

104 The `treetag()` function calls the installation of `TreeTagger` to provide part of speech  
105 tags and lemmas for each token. Importantly, the `path` option should be the directory of the

106 TreeTagger installation.

```
## This function has both suppressWarnings & suppressMessages
## You should first view these to ensure proper processing
temp_tag <- suppressWarnings(
  suppressMessages(
    ## Note: the NULL option is to control for the <unknown> that appears
    ## to occur with the last word in each text
    treetag(c(all_answers, "NULL"),
      ## Control the parameters of treetagger
      treetagger="manual", format="obj",
      TT.tknz=FALSE, lang="en",
      TT.options=list(path=~"/TreeTagger", preset="en"))))
```

107 This function returns a tagged corpus object, which can be converted into a dataframe  
 108 of the token-lemma information. The goal would be to replace inflected words with their  
 109 lemmas, and therefore, unknown values, number tags, and equivalent values are ignored by  
 110 subsetting out these from the dataset.

```
## Remove all tags not using
replacement_lemmas <- temp_tag@TT.res
replacement_lemmas <- subset(replacement_lemmas,
  #ignore punctuation
  wclass != "punctuation" &
  #unknown values
  lemma != "<unknown>" &
  #numbers
  lemma != "@card@" &
  #token should change more than case
```



```
tolower(token) != tolower(lemma))
```

token	tag	lemma	lttr	wclass
is	VBZ	be	2	verb
are	VBP	be	3	verb
trained	VCN	train	7	verb
lots	NNS	lot	4	noun
seats	NNS	seat	5	noun
wings	NNS	wing	5	noun

111 From this dataset, you can use the `stringi` package [] to replace all of the original  
 112 tokens with their lemmas. This package allows for replacement lookup across a large set of  
 113 substitutions.

word	answer
airplane	you fly in it its big it be fast they be expensive they be at an airport you have to be train to fly it there be lot of seat they get very high up
airplane	wing engine pilot cockpit tail
airplane	wing it fly modern technology have passenger require a pilot can be dangerous run on gas use for travel
airplane	wing fly pilot cockpit use gas fast travel
airplane	wing engine passenger pilot(s) vary in size and color
airplane	wing body fly travel

114 b. Weighting

115 c. feature type ontologies

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

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

a. tests to see if things work: e.g. manual spell checks vs automated ones

5. Evaluation of the approach

a. internal (quality, size, consistency) - ?

b. feature size number of features work

ii. classifier for ontology, compare results to previous work

b. externally (categorization, similarity) – MEN dataset, Lapata categorization task

6. Challenges and opportunities

## Discussion

## References