

Ignorance and unawareness

February 2022

1 Hypothesis

The main hypothesis is that the conjoined sentiment values for “ignorant” are significantly lower (more negative) than for “unaware”.

2 Data

The data for this study consists of 5’822 Reddit comments, which were initially collected by using the Pushift API (Baumgartner et al., 2020). The data spans over a time period from 2022-01-01 back to 2011-03-09. Each comment contains a coordinating conjunction of two adjectives, one of which is either “ignorant” (treatment) or “unaware” (control). For simplicity’s sake, only *and*-conjunctions are considered, as conjunctions with *but*, *or*, or *yet* work differently (Elhadad and McKeown, 1990; Hatzivassiloglou and McKeown, 1997). Comments which include a negation of the adjectives (e.g. *not*, *hardly*, *barely*) or any other adverbial modifier (e.g. *very*, *rather*, *mostly*) were discarded.

We are only interested in non-personal attributions of ignorance (and unawareness), as in “stupid and [ignorant/unaware] **of the fact**” or “stupid and [ignorant/unaware] **opinion**”, rather than, for instance, “**he** is stupid and [ignorant/unaware]”. There are two ways to ensure that our data is limited to non-personal attributions: prepositional structures such as [IGNORANT/UNAWARE] **OF** X and [IGNORANT/UNAWARE] **THAT** X typically only take non-personal arguments for X. Another option is to annotate the entity type of the target of attribution—“ignorant **behaviour**” describes an event, “ignorant **idea**” an abstract entity, etc.—, using non-named entity recognition. As both options seem reasonable, we collected two distinct corpora based on these two criteria. For the prepositional structures-approach, we simply queried the Reddit API for A AND [IGNORANT/UNAWARE] **OF** X as well as A AND [IGNORANT/UNAWARE] **THAT** X and only retained those observations, for which A is an adjective. For the entity type-approach, on the other hand, we queried the Reddit API for “[ignorant/unaware] and A”, annotated the data, and only retained attributions targeting abstract entities, events, places, or time. To guarantee a balanced sample, we initially collected 2000 comments each for, on the one hand, “and ignorant of”, “and ignorant that”, “and unaware of”, “and unaware that”, and, on the other hand, “ignorant and” and “unaware and”. Lastly, we annotated the conjoined adjectives with sentiment values from the SentiWords dictionary (Baccianella et al., 2010; Esuli and Sebastiani, 2006; Gatti et al., 2016).¹ The dictionary codes a term’s sentiment intensity on a scale from $-1 \leq x \leq 1$. After preprocessing and filtering out target structures containing adverbial modifiers, the entity type corpus contains 709 observations, the preposition structures corpus contains 5’113.

3 Annotation of the entity type corpus

For the entity type corpus, the data annotation process is much more extensive. In this section, we detail the full procedure.

¹For the sentiment annotation we used the `quanteda`-package (v3.0.0) in R (v4.1.0).

Reddit comments typically span over multiple sentences and make a heavy use of coreferences, i.e. anaphora and kataphora. Since we are only interested in the sentences containing the aforementioned adjective conjunctions (rather than the complete comment respectively), coreferences can ultimately lead to a loss of semantic information, if left unresolved. Hence, we applied a coreference resolution algorithm by Zeldes and Zhang (2016).² The same algorithm also detects the animacy state (*animate* or *inanimate*) and the entity type (e.g., *abstract*, *person*, *object*, etc.) of named and non-named entities mentioned in the texts.³ Ultimately, this allows us to determine whether the adjective conjunction is attributed to (or predicated of) an animate or inanimate entity. Accordingly, we also resolved which phrase the adjectives are attributed to (or predicated of) using a custom function based on syntactic dependency trees.⁴ Due to malformed sentences, not all comments could be annotated. After the annotation step, the corpus retains 709 sentences which include our target structures.⁵

4 Results

Figure 1 shows that the sentiment values are not normally distributed. Hence, we use non-parametric tests. For the entity type corpus, a one-sided unpaired two-samples Wilcoxon test ($W = 33718$, $p\text{-value} < 0.001$) shows that “ignorant” has significantly lower conjoined sentiment values, on average. The same can be shown in the prepositional structure corpus ($W = 2482818$, $p\text{-value} < 0.001$). **Consequently, the hypothesis cannot be rejected.**

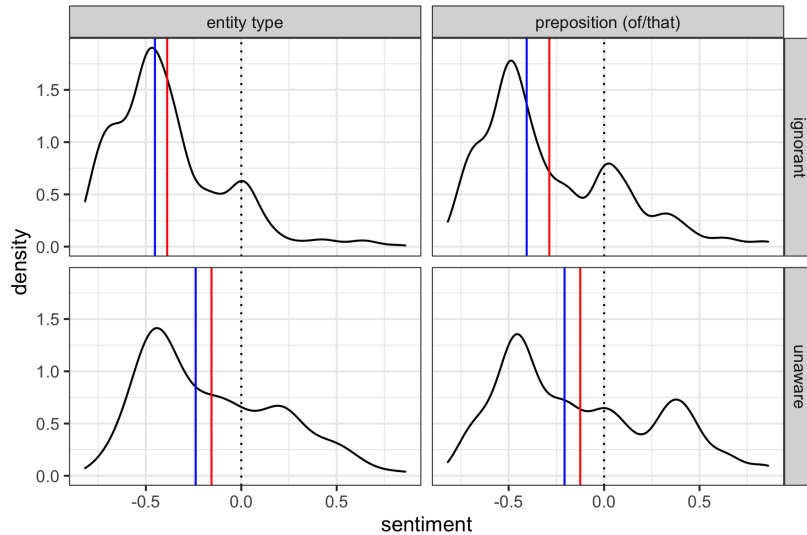


Figure 1: Sentiment distribution. The panels on the left shows the distribution for the *inanimate* subset; the one on the right for the *of/that* subset. The red line depicts the mean, the blue one the median.

We also controlled whether entity type and the kind of prepositional structure (*of* vs *that*) have significant effects, respectively. The factor levels for entity type show no significant differences. The prepositional structure, on the other hand, has a significant effect (linear model, $p\text{-value} < 0.001$):

²The classifier’s performance is reviewed in Sukthanker et al. (2020).

³Both the coreference resolution and the animacy detection are conducted with *xrenner* (v2.2.0.0) by Zeldes and Zhang (2016), based on the pretrained Electra model for GUM7, using Python (v3.7.11).

⁴The dependency parsing was conducted using the *stanza* toolkit (v1.3.0) provided by the Stanford NLP Group Qi et al. (2020) in Python (v3.7.11).

⁵The data, analyses, as well as the full selection of target adjectives can be found on the anonymized Open Science Framework repository at <https://github.com/lucienbaumgartner/ignorance>.

that-constructions have -0.034558 lower conjoined sentiment compared to *of*-constructions, on average, *ceteris paribus*. As the difference is marginal, this effect should not be overestimated.

References

- Baccianella, S., Esuli, A., and Sebastiani, F. (2010). SENTIWORDNET 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining. In *Proceedings of the 7th International Conference on Language Resources and Evaluation*, pages 2200–2204.
- Elhadad, M. and McKeown, K. R. (1990). Generating Connectives. In *Proceedings of the 13th conference on Computational linguistics*, pages 97–101.
- Esuli, A. and Sebastiani, F. (2006). SENTIWORDNET: A Publicly Available Lexical Resource for Opinion Mining. In *Proceedings of the Student Research Workshop at the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 417–422.
- Gatti, L., Guerini, M., and Turchi, M. (2016). SentiWords: Deriving a High Precision and High Coverage Lexicon for Sentiment Analysis. *IEEE Transactions on Affective Computing*, 7(4):409–421.
- Hatzivassiloglou, V. and McKeown, K. R. (1997). Predicting the semantic orientation of adjectives. pages 174–181. Association for Computational Linguistics (ACL).
- Qi, P., Zhang, Y., Zhang, Y., Bolton, J., and Manning, C. D. (2020). Stanza: A Python Natural Language Processing Toolkit for Many Human Languages. In *Association for Computational Linguistics (ACL) System Demonstrations*, pages 101–108.
- Sukthanker, R., Poria, S., Cambria, E., and Thirunavukarasu, R. (2020). Anaphora and coreference resolution: A review. *Information Fusion*, 59:139–162.
- Zeldes, A. and Zhang, S. (2016). When Annotation Schemes Change Rules Help: A Configurable Approach to Coreference Resolution beyond OntoNotes. In *Proceedings of the Workshop on Coreference Resolution Beyond OntoNotes (CORBON 2016)*, pages 92–101.