

# 11830 Report: Biases in Crowdsourced Annotations

Jiyang Tang

jiyangta@andrew.cmu.edu

## Abstract

This document is a supplement to the general instructions for \*ACL authors. It contains instructions for using the L<sup>A</sup>T<sub>E</sub>X style files for ACL conferences. The document itself conforms to its own specifications, and is therefore an example of what your manuscript should look like. These instructions should be used both for papers submitted for review and for final versions of accepted papers.

## 1 Introduction

The Stanford Natural Language Inference (SNLI) corpus<sup>1</sup> is a large crowdsourced natural language inference dataset (Bowman et al., 2015). For each premise sentence, annotators were asked to write a hypothesis that is either a contradiction, neutral statement, or entailment of the premise.

The data was collected from Amazon Mechanical Turk (MTurk) crowdsourcing service. Previous studies have shown that MTurk crowd-workers tend to have lower income, higher education levels, and lower average ages (Levy et al., 2016). We believe that biases are inevitably propagated from the dataset authors to the premise text, and from the crowd-workers to the hypothesis text. Therefore, we will analyze the biases in the data by performing word association tests in this report.

## 2 Method

### 2.1 Pointwise Mutual Information

Pointwise Mutual Information (PMI) is used to measure how much word  $w_i$  is associated with word  $w_j$  (Church and Hanks, 1990; Jurafsky and Martin, 2009). PMI is calculated as follows:

$$\text{PMI}(w_i, w_j) = \log_2 \frac{N \cdot c(w_i, w_j)}{c(w_i)c(w_j)}$$

<sup>1</sup><https://nlp.stanford.edu/projects/snli/>

where  $N$  is the total number of sentences in the corpus,  $c(w_i, w_j)$  is number of times  $w_i$  and  $w_j$  co-occur in a sentence,  $c(w_i)$  is the number of times  $w_i$  occurs in the corpus, and  $c(w_j)$  is the number of times  $w_j$  occurs in the corpus.

Note that if a pair of words  $w_i$  and  $w_j$  occurs multiple times in the same sentence,  $c(w_i, w_j)$  is counted as 1.

PMI ranges from negative infinity to positive infinity. Large PMI suggests high word association. Negative PMI implies two words co-occur less often than by chance and is unreliable in practice (Jurafsky and Martin, 2009).

## 3 Experiments

### 3.1 Data Preprocessing

As mentioned before, each data sample contains a premise and a hypothesis. Note that multiple hypotheses might be generated from the same premise, therefore duplicated premise text is removed. All words are converted into its lower case form and stop words are removed. Then we use spaCy (Honnibal et al., 2020) `en_core_web_sm` model to tokenize raw strings into lists of words and remove all punctuations.

Note that words that occurs less than 10 times in the corpus are removed.

### 3.2 Unigram PMI

We first perform unigram PMI analysis on the entire corpus, and then on the premise text and the hypothesis text individually.

For each analysis experiment, we focus on the most associated words with a set of identity labels (Rudinger et al., 2017).

## 4 Results and Discussion

### 4.1 Unigram PMI

Table 1 lists top associated unigrams with some of the identity labels in the entire corpus, in premise

Identity	Corpus	Premise	Hypothesis
women	saris, headscarves, burkas, notre, husbands, bikinis, gossip, dame, bras, port-a-potty, grains, kimonos, conversating, skirts, coverings, breakroom, headdresses, gypsy, sunhats, cher	saris, headscarves, bikinis, headdresses, coverings, skirts, kimonos, dresses, grains, purses, clap, indigenous, menus, canes, knits, scarves, derby, wigs, badges, jewelery	burkas, husbands, saris, kimonos, bikinis, gossip, conversating, textiles, skirts, gossiping, keyboards, sombreros, incense, gowns, grain, dates, ikea, bistro, congratulate, scarfs
men	turbans, briefcases, wives, cylinders, tuxedos, rickshaws, beards, cigars, jumpsuits, wrestling, wheelbarrows, sashes, canes, wetsuits, kilts, guitars, gutters, mining, supervise, settle	turbans, tuxedos, ladders, jumpsuits, wetsuits, hoses, guitars, hell, dominoes, sashes, solar, suits, beds, netting, ties, trumpets, cannon, kilts, fourth, chess	turbans, rickshaws, wives, cigars, tuxedos, rig, beards, southeast, kilts, gutters, jumpsuits, settle, speedos, guitars, fatigues, banjos, hardhats, naval, meditate, trumpets
africans	shawls, huts, source, hearts, armed, die, tribe, organized, forced, rights, chop, cloths, weapons, diamond, tap, tribal, homes, zebra, belongings, huddle	tribe, hearts, tap, huts, organized, belongings, huddle, clearing, begins, nature, ages, jungle, song, barber, appear, gather, taken, canoe, dinner, digging	armed, source, die, cloths, tribal, forced, rights, weapons, diamond, homes, zebra, america, native, guns, travel, shelter, rice, bananas, loud, supermarket
caucasian	blonds, graffited, slender, straddling, pallet, unpleasant, toronto, lockers, fleece, bared, handsome, non, hairdresser, sunflower, undergoing, bonding, index, leafs	lockers, handsome, explains, contemplates, straddling, pages, leafs, era, frosting, projects, postal, latino, southeast, slender, extending, diverse, addressing, sweaty, stationary, reviewing	slender, fleece, non, zip, festive, handsome, defending, iced, auditioning, heavysset, arrival, barrels, spar, collect, ipod, wildly, freezing, mardi, gras, await
muslims	terrorists, christians, channel, sponsored, celebrate, opening, phones, local, news, skateboards, great, speech, marching, bicycles, new, gather, restaurant, listening, body	NA	christians, terrorists, celebrate, opening, phones, skateboards, local, marching, gather, great, body, new, like, restaurant, shopping, store, talking, beach, near, riding
christians	praising, gospel, lord, muslims, impressed, pork, villagers, lobster, crazy, campfire, sing, church, woods, fun, gathered, selling, music, friends, having, outdoors	NA	muslims, gospel, impressed, pork, villagers, lobster, campfire, sing, church, gathered, selling, music, friends, outdoors, boys, play, running, group, playing, people

Table 1: Top associated unigrams with identity labels

text, and in hypothesis text.

## References

- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.
- Kenneth Ward Church and Patrick Hanks. 1990. [Word association norms, mutual information, and lexicography](#). *Computational Linguistics*, 16(1):22–29.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. [spaCy: Industrial-strength Natural Language Processing in Python](#).
- Dan Jurafsky and James H. Martin. 2009. *Speech and language processing: an introduction to natural language processing, computational linguistics, and speech recognition*, 2nd ed edition. Prentice Hall series in artificial intelligence. Pearson Prentice Hall, Upper Saddle River, N.J. OCLC: 213375806.
- Kevin E. Levay, Jeremy Freese, and James N. Druckman. 2016. [The demographic and political composition of mechanical turk samples](#). *SAGE Open*, 6(1):2158244016636433.
- Rachel Rudinger, Chandler May, and Benjamin Van Durme. 2017. [Social bias in elicited natural language inferences](#). In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, pages 74–79, Valencia, Spain. Association for Computational Linguistics.