

The importance of alternatives in scalar implicature

Benjamin Peloquin

bpeloqui@stanford.edu

Department of Psychology
Stanford University

Michael C. Frank

mcfrank@stanford.edu

Department of Psychology
Stanford University

Abstract

Successful communication regularly requires listeners to make pragmatic inferences - enrichments beyond the literal meaning of a speaker's utterance. For example, listeners routinely enrich the meaning of "some" to "some, but not all" when interpreting a sentence such as "Bob ate some of the cookies." A Gricean account of this phenomena assumes the presence of *salient alternatives* with varying degrees of informativity. "Some," in the example above, is enriched to "some, but not all" in the presence of the stronger alternative "all." Empirical evidence for the presence of such alternatives and accounts of their effect on implicature has been limited. Our current work explores the role different scale representations may play in scalar implicature using empirical measures of literal semantics and a Bayesian model of implicature generation. Comparisons with human judgments indicate that pragmatic inference may rely on fairly complete alternative sets rather than the typical entailment only scales assumed in most research to date.

Keywords: pragmatics; scalar implicature; bayesian modeling

Introduction

Humpty Dumpty Quote...

Successful communication regularly requires listeners to make inferences that go beyond the literal semantic content of a speaker's utterance. "Scalar implicature" is a hallmark of such pragmatic enrichment. For example, listeners routinely enrich the meaning of the scalar item "some" to "some, but not all" in sentences like "Bob ate some of the cookies." (citations) A Gricean account of this phenomena assumes listeners reason about the intended meaning of a speaker while incorporating knowledge about a) alternative scalar items a speaker could have used (such as "all") and b) the relative informativity of using such alternatives (Grice, 1975). Continuing with the cookie example, a listener will infer that the speaker must have intended that Bob did not eat "all" the cookies because it would have been *underinformative* to use "some" when "all" could have been used.

The presence of "salient" alternatives is critical to standard accounts of scalar implicature (Grice). Recent studies indicate implicature generation is affected by the set of alternatives (hereafter 'scale representations') available to listeners, both in adults (Degen, Franke) and is possibly a source of difficulty in computing implicatures for children (Barner). [Degen found ... Similarly, Franke demonstrated... (the importance of a strong negative valenced term such as "none" adopting a Bayesian model of implicature generation. Model performance improved significantly when "none" was included...)] Despite these signals, little empirical evidence points directly to the impact of scale representation in scalar implicature. The current work investigates the role different scale representations may play in scalar implicature.

While the notion of a set of alternatives is fairly intuitive in the context of scalar items (it's part of how their defined), measuring alternatives is not. Issues of introspection aside, simply querying a participant about the saliency of a particular alternative is problematic because the alternative must be made salient during the actual query. To avoid this issue we adopt a computational model, formalizing implicature generation as Bayesian inference. Using empirical measures we simulate the effects of different alternatives sets on model predictions. Model predictions are compared to human judgment data.

Bayesian models of pragmatic enrichment have accurately predicted human judgments in ad-hoc (Frank & Goodman, 2012) and embedded (Stulhummiller & Goodman, 2014) implicature settings. Our current model follows both these studies and is based on Rational Speech-act theory (RSA). RSA frames language understanding as a special case of social cognition (Stulhummiller & Goodman, 2014), in which listeners and speakers reason about one-another. In the following set of experiments, we outline a literal semantics task (Experiments 1a, 3a and 4) as well as pragmatic judgment task (1b, 3b). Data from the literal semantics task are used to approximate information a Speaker might assume when reasoning about a Listener who interprets semantics literally (without pragmatic enrichment). Likewise, pragmatic judgments are used to compare model posterior predictions and assess quality of fit.

The heart of this investigation is concerned with populating our computational framework with empirical data. In particular we use experimental tools to measure otherwise ambiguous literal semantics, to generate a set of plausible alternatives, and to obtain human pragmatic judgments for model comparison. Recent experimental investigations of scalar implicature have focused almost exclusively on a small subset of possible scales, most notably "some/all". This lack "Scalar Diversity" (van Tiel, 2013) makes generalization across different scalar families problematic. We adopt a "Scalar Diversity" approach in our experimental design in order to incorporate multiple scalar items from a range of grammatical classes (see Figure 5 for a complete set of scalar items used). In particular, we use a food review paradigm to quantify literal semantics and pragmatic judgments as distributions over star-ratings.

In the following paper we investigate the role of scale representation in scalar implicature. We use a food review paradigm in which 'meaning' for a given scalar item was quantified as a distribution over star-ratings. Quantifying both literal semantic content and pragmatic judgments as distributions over stars provides an important interface to



Figure 1: The left panel shows a trial from Experiment 1a with the target scalar ‘liked’. Participants responses to the binary dependent variable are used to quantify literal semantics. The right panel shows a trial from Experiment 1b with the target scalar ‘liked’. Participants were asked to generate the star rating they think the speaker likely intended, given the target scalar (in this case ‘liked’).

our model, especially for scalar items with ambiguous literal semantics (e.g. items such as “liked/loved” or “memorable/unforgettable”). In addition to quantifying literal semantics and pragmatic judgments in Experiments 1a,b and Experiments 3a,b, Experiment 2 used a modified cloze task, based on van Tiel (2013, Experiment 2??), to generate a set of plausible alternatives to the scalar items used in Experiments 1a,b. We extend the set of alternatives measured in Experiments 3a,b and 4 using data from Experiment 2. Additionally, by looking at the relative frequencies of alternatives generated in Experiment 2 we were able to compute a saliency measure for each alternative. Using the literal semantic data from Experiments 1a, 3a and 4, we modeled pragmatic enrichment using RSA and compared model predictions to human judgments obtained in Experiments 1b and 3b. Model performance was significantly improved with the incorporation of alternatives. This result suggests that more “fully specified” scale representations may be active for human participants while making pragmatic judgments for scalar items.

Experiment 1a,b: Entailment scales

Experiment 1a and 1b were conducted to approximate literal listener semantic distributions (1a) and pragmatic judgments (1b) for entailment only scales (e.g. “liked/loved”, “good/excellent”, etc).

Methods

Participants Participants for both tasks were recruited on Amazon Mechanical Turk and paid \$0.20 for their participation. Thirty participants were recruited for Experiment 1a. Data for N participants was thrown out after participants either failed to pass two training trials or were not native English speakers, leaving a total sample of N. Fifty participants were recruited for experiment 1b. Data for N participants was thrown out after participants either failed to pass two training trials or were not native English speakers, leaving a total sample of N.

Design and procedure The left panel of Figure 1 shows a screen shot of a trial from Experiment 1a. Participants were

presented with a target scalar item and a star-rating (between 1-5 stars) and asked to judge the compatibility of the scalar item and star-rating. Compatibility was assessed through a binary “yes/no” response to a question of the form “Do you think that the person ___ed the food?” where a target scalar was presented in the “___ed.” Each participant was presented all scalar item and star-rating combinations with randomization.

The right panel of Figure 1 shows a screen shot of a trial from Experiment 1b. On each trial, participants were presented with a one-sentence prompt containing a target scalar item such as “Someone said they ___ed the food.” They were then asked to generate a star-rating according to what they think the reviewer likely gave. Each participant was presented all scalar items with randomization.

Results and Discussion

The left panel of figure two shows literal semantics for each of the scalar terms “liked” and “loved” as distributions over star-ratings.

Chi-squared test? Null results, not really important here anyway...

Experiment 2: What are the alternatives?

In Experiment 2 we asked participants to generate plausible alternatives for the scalar items presented in Experiment 1.

Methods

Participants Using Amazon’s Mechanical Turk, N workers were paid \$0.20 to participate. n workers were dropped from the analysis because they either failed to pass training trials or were not native English speakers. The final sample was N workers, all of whom were naive to the purpose of the experiment.

Design and procedure Participants were presented a target scalar term embedded in a sentence such as, “In a recent restaurant review someone said they thought they ___ed the food.” Participants were then asked to generate plausible alternatives by responding to the question, “If they’d felt differently about the food, what other words could they have used instead of ___?” and asked to generate three alternatives.

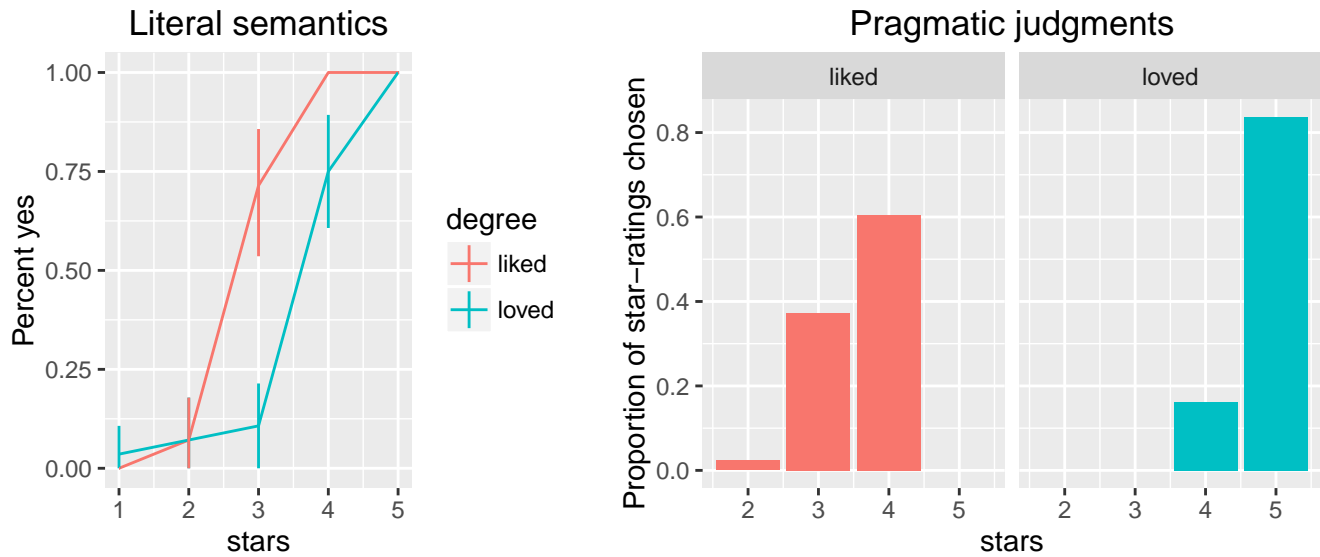


Figure 2: The left panel plots literal semantic distributions for the scalar pair 'liked/loved'. The y-axis is the percentage of 'yes' responses for a scalar term that is compatible with the number of stars on the x-axis (ie 100% of respondents beleived that 'liked' was compatible with both 4 and 5 stars). Erro bars are 95% confidence intervals. The right panel plots pragmatic judgments for the scalar pair 'liked/loved'. The y-axis is the proportion of selection of the star-rating on the x-axis (ie over 80% of judgments when prompted with 'loved' were 5-stars).

Results and Discussion

Figure 3 plots the combined counts of alternatives generated for the scalar items “liked” and “loved.” Alternative distributions for the other scalar pairs (e.g., “good/excellent”, “memorable/unforgettable”) were similarly long-tailed. In Experiments 3a,b we took the top two alternatives generated for each scalar item to enrich the entailment only scales from Experiments 1a,b. In experiment 4 we further enriched the scales from Experiments 3a,b with a neutral valence scalar chosen from the alternative sets.

Experiment 3a,b: Incorporating top alternatives

Experiment 3a,b are identical to Experiments 1a,b except the set of target scalar was expanded to include the top two alternatives generated for each scalar family. The orange colored items in table 1 denote additional scalar items added in Experiments 3a,b.

Participants

Participants for both studies were recruited on Amazon Mechanical Turk and paid \$0.20 for their participation. Thirty participants were recruited for Experiment 3a (Literal Listener task). Data for N participants was thrown out after participants either failed to pass two training trials or were not native English speakers, leaving a total sample of N.

Fifty participants were recruited for Experiment 3b (Pragmatic Listener task). Data for N participants was thrown out after participants either failed to pass two training trials or were not native English speakers, leaving a total sample of N.

Procedure

The procedure of Experiments 3a,b follow the same form as Experiments 1a,b with the expanded set of target scalar items.

Results and Discussion

Distributional differences.

Experiment 4: Symmetric scales

Experiment 4a is identical to Experiments 1a and 3a except the set of target scalar was expanded to include a neutral valence scalar item for each scalar family. The purple colored items in Table 1 denote additional scalar items added in Experiments 4.

Participants

Participants for both studies were recruited on Amazon Mechanical Turk and paid \$0.20 for their participation. Thirty participants were recruited for Experiment 4 (Literal Listener task). Data for N participants was thrown out after participants either failed to pass two training trials or were not native English speakers, leaving a total sample of N.

Procedure

The procedure of Experiment 4 follow the same form as Experiments 1a and 3a with the addition of a neutrally valence scalar item.

Results and Discussion

Distributional differences.

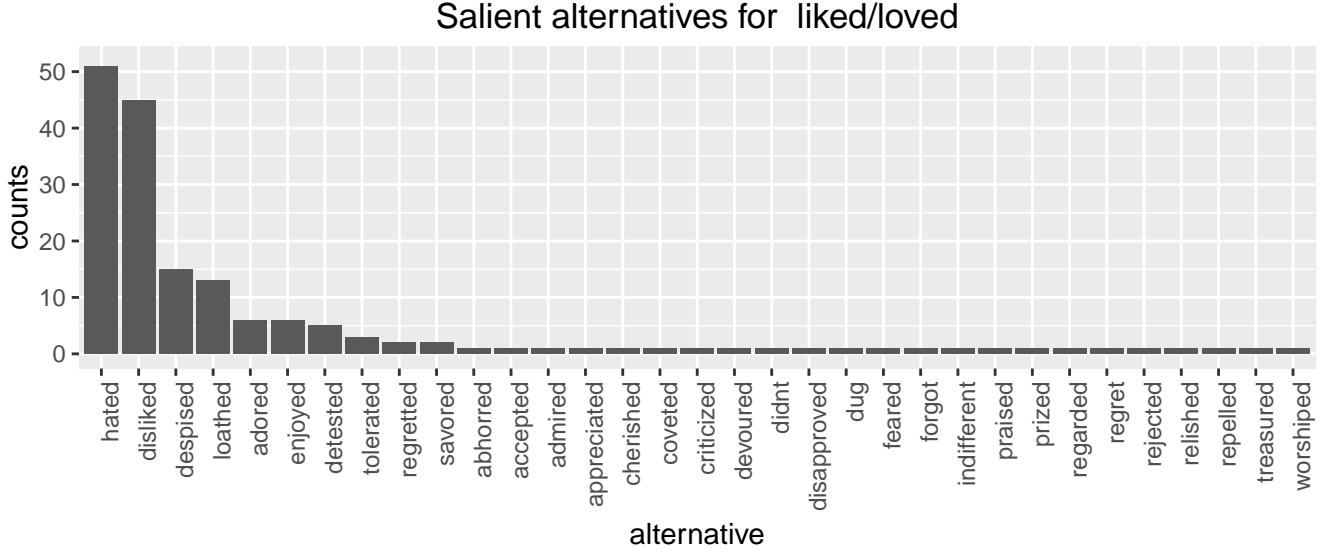


Figure 3: Caption goes here

Model

We adopt a Bayesian model of scalar implicature. In particular, RSA frames language understanding as a special case of social cognition (Stuhmüller & Goodman, 2014) in which listener and speaker agents reason about one another. We focus on the problem of a listener inferring the meaning of a speaker’s utterance. Let u be an observed utterance by our speaker with an intended meaning m . We assume that the listener has access to a space of possible meanings M and some model relating meanings $m \in M$ to u .

Upon hearing an utterance u the Listener agent evaluates all candidate word meanings, computing their posterior probabilities $p(m|u)$. This quantity is proportional to the product of the prior probability $p(m)$ of a particular meaning and the likelihood $p(u|m)$.

$$p(m|u) = \frac{p(u|m)p(m)}{p(u)} \\ = \frac{p(u|m)p(m)}{\sum_{m \in M} p(u|m)}$$

The prior $p(m)$ represents the listener’s expectations about plausible meanings (star ratings), *independent* of the utterance u . To avoid biasing the data we assume priors over meanings are uniform, however future investigations may want to experiment with different priors, either by inferring them or using empirical measures

The likelihood $p(u|m)$ assumes that speakers choose words to maximize informativity in context. We quantified the informativeness of a word by its surprisal minus a cost. Cost, in this case is operationalized as a salience measure - more salient items should have a larger impact on implicature generation than less salient alternatives.

$$p(u|m) = \frac{e^{-\alpha(-\log(p(m|u)) - \text{cost}(u))}}{\sum_{m \in M} e^{-\alpha(-\log(p(m|u)) - \text{cost}(u))}}$$

The posterior $p(m|u)$ reflects the listener’s degree of belief that m is the intended meaning of the speaker, given her utterance u .

Using literal semantics data from experiments 1a and 3a and 4, we can simulate the role alternatives play in scalar implicature. In particular, we can vary the number of scalar items present during implicature generation. In the following analyses we focus on three scenarios 1) Entailment only scales (good, excellent), 2) Negative Alternatives (horrible, bad, good, excellent) 3) Symmetric scale (horrible, bad, okay, good, excellent). Table 1 contains the set of scalar items used for each simulation.

Model fitting and comparison

Results Table here

Figure 4 plots model / human correlation for our three simulations. Alpha was tuned separately for each of these runs. . .

General Discussion

By varying the type of scale representations available to our Bayesian model we investigated the effects of alternatives on scalar implicature. Model fit with human judgement was significantly improved by the inclusion of alternatives beyond the typical “strong” and “weak” scalar items. In fact, we found that both neutral and negative valence scalar items contribute to human-like implicature generation within our framework.

\begin{CodeChunk} \begin{figure}[h]

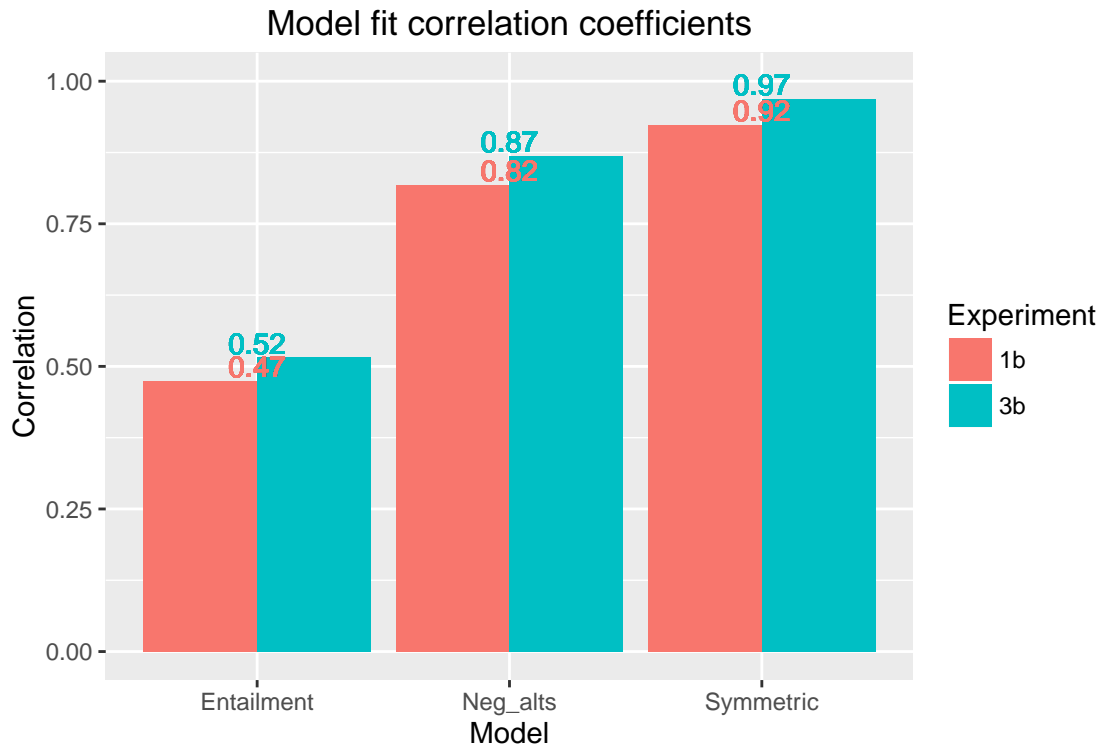


Figure 4: The left panel shows improved model fit as scale representations are enriched with more scalar items. Correlations are computed using pragmatic judgment data from Experiments 1b and 3b. The right panel plots model predictions using full symmetric scales versus human judgments from Experiment 3b.

Alternative sets				
good / excellent	liked / loved	memorable / unforgettable	palatable / delicious	some / all
excellent	loved	unforgettable	delicious	all
good	liked	memorable	palatable	most
okay	felt indifferent about	ordinary	mediocre	some
bad	disliked	bland	gross	little
horrible	hated	forgettable	disgusting	none

Entailment items used in Experiments 1a,b
 Top two alternatives added in Experiments 3a,b
 Neutral item added in Experiment 4

This table shows the stimuli used in Experiments 1a,b, 3a,b and 4. Colors denote the additions made in each experiment. For example, in Experiment 1a we measured literal semantics for ‘liked/loved’. In Experiment 3a,b we extended this group to ‘hated/disliked/liked/loved’. In Experiment 4 we extended this group to ‘hated/disliked/felt indifferent about/liked/loved’.

Acknowledgements

Thanks to NSF BCS XYZ. Thanks to Michael Franke, Judith Degen, and Noah Goodman.

References

the role various scale representations might play in scalar implicature. In Experiment 1a,b we gather data for Entailment only scalar items such as “some/all,” “good/excellent,” etc. In Experiment 1a we quantify literal semantics for scalar items as compatibility distributions over star-ratings (Literal listener task). In Experiment 1b, we ask participants to make judgments about speaker intentions by generating star-ratings in response to reviews (Pragmatic listener task). In Experiment 2, we ask participants to generate plausible alternatives to the scalar items used in Experiments 1a,b. Experiments 3a,b are identical to 1a,b with the addition of scalar items from Experiment 2. In this we explore fuller scale descriptions, gathering both literal semantics and pragmatic judgments for fuller scale descriptions such as “none/some/most/all” or “horrible/bad/good/excellent.”

In three experiments we Our “Scalar Diversity” approach involves a food review paradigm we employed throughout our three experiments. In both Literal Listener (1a, 3a) and Pragmatic Listener (1b, 3b) tasks, participants are shown one sentence prompts containing a target scalar term along with a star-rating. In individual trials, star-ratings were

either pre-populated (in literal listener tasks) or used as the dependent variable (pragmatic listener tasks). This design allowed us to quantify literal semantics for individual scalar terms as compatibility distributions over star-ratings (literal listener tasks 1a and 3a). Similarly, we quantified pragmatic judgments by prompting participants with a target scalar term and asking them to generate a star-rating.