# NLP Project Proposal: Definiteness

Christopher Brown

chrisbrown@utexas.edu

April 8, 2013

### Abstract

Noun phrases occurring in natural language can be described as either definite or indefinite. Definiteness is a semantic concept that relates to givenness, familiarity, topicality, and focus. It can be denoted by a variety of determiners, the absence of a determiner, or by context. This project focuses on the definite markers, *a* and *the*, and investigates what types of models are best at predicting the most suitable determiner when it is not available.

## 1 Semantics of definiteness

The theoretical semantics of definite noun phrases has an extensive history; we can start with Gottlob Frege and Bertrand Russell and names. Frege said that the meaning of a name (or definite noun phrase) was the object denoted by that description. Russell complicated matters with his "The King of France is bald" example, in which "the King of France" cannot refer to anything (at least not since 1848). Russell would call this sentence false, despite the description's failure to refer to anything; the failure of the noun phrase propagates out to the failure of the sentence.

Keith Donnellan and P. F. Strawson disagree; they say that referencing is something speakers, not words, do, and that "the King of France" is simply unsuccessful. Compare "a King of France is bald," which can refer to any (dead) past King of France, and thus have a truth value. Definite noun phrase usually denote a unique entity, and if this fails, the sentence is incoherent. But compare "After you enter the lobby, take the elevator to the 13th floor," which is acceptable even if there is more than one elevator.

This is just the beginning, but it's clear that the difference between definite and indefinite noun phrases involves a number of factors, such as context, real world knowledge, and other pragmatic phenomena. It is an active field of study in linguistics or philosophy of language, but has received very little attention in computational linguistics.

## 2 Applications

Here are two scenarios where we need to determine the appropriate definiteness of a noun phrase:

1. **Generation.** When generating language from a logical structure, we need to know what kind of referential noun phrase to use when representing a certain entity. This is relevant to anaphora generation, but it's simpler; disregarding pronouns for current purposes, if we use a full noun phrase more than once, we need to know when to use an indefinite description and when to use the definite form.

2. **Translation.** Some languages, most notable Russian, do not have determiners to denote definiteness. When translating between these languages and languages with overt definiteness markers, how do we determine where to insert determiners.

For example, Google Translate renders each of the four alternations of "A/The man bit a/the dog" as Человек укусил собаку. Translating back into English produces "A person bitten by a dog." I don't know Russian; it may be genuinely ambiguous, and only determinable via context; but I know that translating back into English accurately requires choosing the correct determiners, for which we need a model of definiteness (as well as anaphora / centering).

There are a few papers on the second of these applications, translation, but they are specific to only a few language pairs (Ishikawa, 1995, Siegel, 1996).

## 3   Experiment

The basic research question will be to determine what type of language models produce the best prediction of definiteness. Because English has determiners, any parsed document is "labeled" data. I will simply replace *a* and *the*'s with the placeholder «DET», retaining the original token as the label for that token. While this is an artificial evaluation metric, it could be useful in both of the application scenarios. In anaphora generation, each instance has to be filled by either a full noun phrase, either definite or indefinite, or a pronoun. In the noun phrase case, a placeholder would be inserted by a first pass of the language generation model, and then resolved by the anaphora resolution pass. In machine translation, noun phrases translated from a determiner-less language to one with determiners could be padded with a placeholder, which would then be resolved by some post-processing step performed on just the target language.

My project will consist of evaluating different models for predicting the deleted marker; it is not a typical sequence labeling problem, since relatively few of the tokens need to be resolved. But neither is each instance independent of the previous; in fact, a prior mention is presumably a strong indicator that an entity is now in the active context, and subsequent instances should use the definite determiner. The sequence is crucial in a sense of building and maintaining a 'center,' which is a common approach when resolving anaphora (Grosz et al., 1995, Beaver, 2000). But with an intelligent set of features, we might achieve sufficient accuracy with a simple logistic regression; given a set of features, evaluate whether the placeholder should be definite or indefinite?

My investigation will also compare feature sets, to determine which are the most efficient at predicting the deleted definiteness markers. Is the parse structure important, or just the surrounding tokens? Position in sentence / document? One presumably useful feature would be to use a cumulative index for each token type in the document—1 for the first instance, 2 for the second, and so on—which might handle the prior mention issue without requiring a sequence model.

## References

David Beaver. Centering and the optimization of discourse. http://www.blutner.de/Optimal/dat/beaver_centering.pdf or http://goo.gl/UgyKk, July 2000.

Barbara J. Grosz, Aravind K. Joshi, and Scott Weinstein. Centering: a framework for modeling the local coherence of discourse. *Comput. Linguist.*, 21:203–225, June 1995.

Kiyoshi Ishikawa. Crosslinguistic notions of (in)definiteness. In *Proceedings of the 10th Pacific Asia Conference on Language, Information and Computation*, 1995.

Melanie Siegel. Preferences and defaults for definiteness and number in japanese to german machine translation. In *Language, Information and Computation (PACLIC 11)*, pages 43–52, 1996.