
Semantic Indexing

Anant Bhardwaj
Ekaterina Ob'yedkova
Ravi Netravali

ANANTB@CSAIL.MIT.EDU
KATYA_@MIT.EDU
RAVINET@MIT.EDU

1. Introduction

Majority of information retrieval and web search use the inverted index as the backbone for their keyword based search. Unfortunately, inverted index fail to capture the semantic of the language and thus keyword based approach can only go as far as giving the relevant hits. We propose semantic indexing scheme which can encode natural languages without losing semantic relationships.

The proposed semantic encoding would involve parsing of each sentence, deriving its semantic model and storing the encoded representation. Matthew Fay and Jodyann Coley in his work SemanticSQL have attempted to encode natural language in SQL. While SemanticSQL can encode simple sentences, it fails to encode and accurately recover the semantic information from the stored representation for complex sentences. We aim to build a model that can derive semantic structure from complex sentences and a new schema which can accurately store the derived semantic model. In the following paragraph we illustrate different cases their system currently can not handle and propose a new semantic model which incorporates them.

1. Probabilistic Parsing: Our system would incorporate a probabilistic sentence parser which would allow the system to intelligently identify the part-of-speech of previously unknown words and appropriately add them to the index store.
2. Multiple adverbs and prepositional phrases: The current implementation only allows a single occurrence of location, adverb, and instrument phrases. Our system would allow sentences such as "John slept under the tree in the park" to be understood, and the question "Where did John sleep" would return "under the tree in the park".
3. Noun Phrases with relative clauses, and with time: Show me all the widgets that were selected last Tuesday.
4. Handling complex negation (negation of x AND Y kind of phrases) cd

5. Intersective vs. Non-intersective adjectives

intersective adjectives: The car is a red Volkswagen. Grace is a hairy brown dog. We drove up yet another serpentine mountain road.

non-intersective adjectives: Viktor is a former Catholic. The alleged thief has arrived in court.

6. Complex causality handling: The current implementation only allows a primitive causality of 'if then' kind. Expanding it to more complex causality mapping would allow us handle complex causality.

Anna told Alex to return the chair because it was broken. They moved here because of the work. Just because John is tired we decided to stay at home. Just because John is tired does not mean Mary and John will not go to the family dinner. We can leave after the rain stops. If Mary finishes her homework and the weather is good then Mary will go for a walk.

Relating events based on time: Did John come home after the fire alarm went off?

7. Handling constraints around conjunctions and disjunctions... Depending on the weather, John will go for a walk or watch a movie.
8. Handling quantifiers. For instance, we would be able to answer questions like:
How many times have John watched the movie? How much milk have John got in the fridge? Is there any milk in the fridge?

The indexer would do a semantic encoding of english text corpus to build its index store. As new data comes in, the index store would get updated. For querying, we would allow a web interface that can take natural language query and would search the index store to find an answer to the query. We would also provide APIs for building the index and a SQL/python interface for searching the index.