

Thesis

Ryan Tanner

January 29, 2012

Contents

1	Introduction	2
2	Problem Statement	2
2.1	Extracting properties defined in a massively large body of text	2
2.2	Doing so in a time-efficient manner	2
2.3	On a textual level, the problem of finding connections across concepts and entities in a massive corpora of text	2
3	Other Approaches	2
3.1	Training oriented approaches	2
3.1.1	Manually-annotated training input	2
3.1.2	More accurate than my proposed solution	2
3.2	Weakly-linked crowd sourcing	2
4	Importance of the Problem	2
4.1	The Problem of Big Data	2
4.2	Tracing Influence	3
5	Approach to Solving the Problem	3
5.1	Treating grammatical dependencies as functions	3
5.2	Mapping the governors and dependents of those dependencies to textual aliases and named entities	3
5.3	Reducing a set of input documents to find connections between those aliases and entities based on their common properties	3
5.4	Constructing a graph of these connections where the connections form weighted vertices and entities form nodes	3
5.5	Visualizing this graph	3
5.6	Why a functional language?	3

6	Algorithm in Detail	3
6.1	Dependency Functions	3
6.2	Properties	4
7	Results	4
8	Future Recommendations	4
A	Some Relevant NLP Concepts	4
A.1	Dependency Grammars	4
B	Tools Used	4
C	Code Highlights	4

1 Introduction

2 Problem Statement

This thesis is an attempt to tackle the problem of extracting facts and connections from written text. Massive quantities of text are produced daily and methods for quickly getting relevant information out of that text are needed. There are many

2.1 Extracting properties defined in a massively large body of text

2.2 Doing so in a time-efficient manner

2.3 On a textual level, the problem of finding connections across concepts and entities in a massive corpora of text

3 Other Approaches

3.1 Training oriented approaches

3.1.1 Manually-annotated training input

3.1.2 More accurate than my proposed solution

3.2 Weakly-linked crowd sourcing

4 Importance of the Problem

4.1 The Problem of Big Data

Google alone processes over twenty petabytes of data per day (Dean and Ghemawat, 2008).

4.2 Tracing Influence

5 Approach to Solving the Problem

Most approaches to this problem rely on extracting as much information as possible from a given input. This approach comes at the problem from the opposite direction and tries to extract a little bit of information very quickly but over an extremely large input set.

5.1 Treating grammatical dependencies as functions

This approach is based on the premise that dependency grammar relations can be treated as functions and modeled as such. Furthermore, I hypothesize that these functions can be curried, just as in a functional language. Every word in a sentence, save for the head, is dependent upon another word and each of these dependencies has a type. This structure forms a tree. By doing a depth-first traversal of this tree and recursively composing each individual dependency function into a curried function, we end with a function specific to that sentence.

In this approach, dependency functions are short operations which extract properties from the given relation. These functions take two nodes of a tree as input, the governor and the dependent. Based on the types of the tokens in each node a partial or full property is added to the accumulator map and returned up the tree. This map is comprised of entities mapped to properties representing pieces of information extracted from the relationship. More about properties can be found in section 6.2 on page 4.

5.2 Mapping the governors and dependents of those dependencies to textual aliases and named entities

5.3 Reducing a set of input documents to find connections between those aliases and entities based on their common properties

5.4 Constructing a graph of these connections where the connections form weighted vertices and entities form nodes

5.5 Visualizing this graph

5.6 Why a functional language?

6 Algorithm in Detail

6.1 Dependency Functions

The grammar dependencies used here are those described in the Stanford typed dependencies manual [?](#). Currently 53 grammatical relations are defined for the English language. Each of these has a corresponding function in this algorithm. Though the specifics of

each function differ, all follow the same simple pattern. Dependency functions take two parameters, a governor and a dependent, and return a map of tokens to a list of properties. Furthermore, these grammatical relations have a typed hierarchy where relations can inherit from other relations. Each function therefore can use its supertype's own function and only add the minimum processing necessary for its specific relationship.

6.2 Properties

7 Results

8 Future Recommendations

A Some Relevant NLP Concepts

A.1 Dependency Grammars

B Tools Used

C Code Highlights

References

Dean, Jeffrey and Sanjay Ghemawat. “MapReduce: simplified data processing on large clusters.” *Commun. ACM* 51 (January 2008): 107–113.