

B.Sc. in Computer Science and Engineering Thesis

Building A Knowledge Graph Using Twitter Data

Submitted by

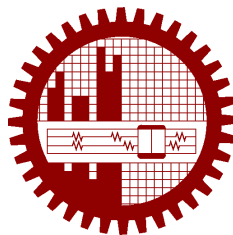
Sukanta Kundu
201205068

Anindya Biswas
201405006

Qazi Fahim Farhan
201405052

Supervised by

Dr. Muhammad Masroor Ali



Department of Computer Science and Engineering
Bangladesh University of Engineering and Technology

Dhaka, Bangladesh

February 2018

CANDIDATES' DECLARATION

This is to certify that the work presented in this thesis, titled, “Building A Knowledge Graph Using Twitter Data”, is the outcome of the investigation and research carried out by us under the supervision of Dr. Muhammad Masroor Ali.

It is also declared that neither this thesis nor any part thereof has been submitted anywhere else for the award of any degree, diploma or other qualifications.

Sukanta Kundu

201205068

Anindya Biswas

201405006

Qazi Fahim Farhan

201405052

CERTIFICATION

This thesis titled, “**Building A Knowledge Graph Using Twitter Data**”, submitted by the group as mentioned below has been accepted as satisfactory in partial fulfillment of the requirements for the degree B.Sc. in Computer Science and Engineering in February 2018.

Group Members:

Sukanta Kundu

Anindya Biswas

Qazi Fahim Farhan

Supervisor:

Dr. Muhammad Masroor Ali

Professor

Department of Computer Science and Engineering

Bangladesh University of Engineering and Technology

ACKNOWLEDGEMENT

We are indebted to our supervisor Professor Dr. Muhammad Mashroor Ali for introducing us to the Knowledge Graph. We've learned how to think critically like a researcher, carry on a research, read scientific papers, and how to write, speak and present properly. We show our gratitude for his dedication and patience in correcting our mistakes, reviewing our ideas, methodologies, proofs, etc. We couldn't have done it without him.

Dhaka
February 2018

Sukanta Kundu

Anindya Biswas

Qazi Fahim Farhan

Contents

List of Figures

List of Tables

List of Algorithms

ABSTRACT

In this thesis, we wish to work on creating a Knowledge Graph using Twitter micro-data. A Knowledge Graph is a representation of human knowledge in the form of a graph. Just like any graph, it has nodes (entities), attributes and edges (relations). Since tweets give us limited data (maximum length of 300 characters), it becomes very difficult to extract adequate information to create a knowledge graph.

Chapter 1

Introduction

The present world is a place with immeasurable data or informations due to world wide web. Each and every day this amount is increasing. It is getting harder day by day for people to extract exact information he/she needs and to analysis them to have a desired goal. However this big data has become a challenge for computer science to organize, store and most importantly efficient searching. (As far as the word efficient searching means searching a particular word with all its variations.)

A single word can represent a person, place, name of organization, restaurant . Again a word slightly different spelling can represent totally different thing. It is very hard for a single person or a particular software to do this. So, we need to make machine learn way of efficient searching.

It may be helpful for us to write grammars for machine to learn than writing all the rules of searching by ourselves. Now for machine to learn to work or search in order to search efficiently we need data or information to be organized in a way that machine can analyze those data. Here knowledge graph comes into play.

Chapter 2

Related Works

Previously there have been many researches conducted on knowledge graphs. Let us discuss them one by one.

In the paper ‘Knowledge Graph Identification’ [?], the authors have extracted uncertain entities and their relations to create an ‘extraction graph’. Next they removed noise, added missing information and determined candidate facts that resulted into a Knowledge Graph Identification. Reasoning jointly about candidate facts, extracting their confidence, identifying coreferences, imposing ontological constraints, removing duplicate entities, resolving ontological constraints were some of the other tasks. To disambiguate and relate a concept with similar concept, the researchers have used Probabilistic Soft Logic (PSL) [?] ontological constraints, identify coreferences etc.

PSL has some advantages. Models can be easily defined using declarative rules with first order logic syntax. We can have continuously valued random variables that is convenient with uncertainty or probability. Using weights, we can control the importance of model rules.

One common problem in building a Knowledge Graph is to find appropriate links between various concepts. Mihalcea and Csoma et al [?] have worked on link detection. They used machine learning on a wikipedia dump to link a concept(article) with its corresponding mention. There they have also calculated link probability of various phrases or mentions. Define link probability of a phrase as the number of wikipedia articles that use the phrase as anchor divided by number of articles that mention the phrase at all. Next they disambiguated links by taking help of the surrounding words (the words themselves and their parts of speech). Here a problem occurs on which meaning to select. The best method is to take the most common meaning first. Then calculate the relatedness with the context. Define $relatedness(a, b) = \frac{\log(\max(|A|, |B|)) - \log(|AB|)}{\log|W| - \log(\min(|A|, |B|))}$, where a, b are the Article of interest, A, B are sets of all articles linking to a, b respectively and W is the set of all articles.

After that, by comparing the common meaning with the relatedness, the correct link was found.

For example, the word ‘tree’ has a common meaning ‘plant’ and so 97% of the time, it will link to a page on plants. But if we find ‘binary tree’ in a document, the relatedness of ‘tree’ to ‘binary’ suggests that this phrase is related to data structure. That is how almost everyone disambiguates links. Followed by this, they did Topic indexing which aims to identify the most significant topics, summarize the document and organize it into various categories, and to answer queries fast. In topic indexing, the key problem is to find the correct significant terms and disambiguate them to appropriate topic.

The Next work deals with a semi-supervised graph regularization method on twitter data wikification [?]. There the authores described the problems as unlinkability (absense of a valid concept from knowledge base), ambiguity and prominence. Some more major problem for working with twitter data are informal writing style, shortness and noisiness. To deal with these problems, the authors developed a graph-based semi-supervised learning algorithm for wikification. Their model at the same time, detects mentions and disambiguates them at both local and global levels. They have also developed meta path-based unified framework to detect relevant mentions. Their method works in this way: given a set of tweets $\langle t_1, t_2, \dots, t_n \rangle$, first find candidate concept mentions $\langle m, c_1, c_2, \dots, c_k \rangle$. From there, we find the most probable concept c related to mention m and so we get $\langle m, c \rangle$. After that, they constructed a relational graph $G = \langle V, E \rangle$.

Here set of nodes $V = \langle v_1, \dots, v_n \rangle$, and set of edges $E = \langle e_1, \dots, e_e \rangle$. Here $v_i = \langle m_i, c_i \rangle$.

Finally, we have an important work of Larry Heck and Hongzhao Huang [?]. First they developed a robust way to represent concepts. Many NLP tools use string concatenation to represent ideas. But most of the time it fails to represent appropriate concepts. For example ‘Madrid’ is a place, ‘Real’ means something authentic or true, but ‘Real Madrid’ refers to a football club. This is a counter example where string concatenation fails. A possible alternative is to use word hashing with n -gram (tri-gram in this case) word representation. This way turns words into concept-vectors. This representation is robust because it can be used even for unseen words. After that they did neural embedding of knowledge graph. To do so, they tracked a concept and its corresponding subgraph, encode the knowledge as featured vector. then they trained Deep Neural Network to get semantic relationships among various concepts. After that, they took tweets and used their deep neural network on those tweets.

Appendix A

Algorithms

A.1 Sample Algorithm

In Algorithm ?? we show how to calculate $y = x^n$.

Algorithm 1 Calculate $y = x^n$

Require: $n \geq 0 \vee x \neq 0$

Ensure: $y = x^n$

$y \leftarrow 1$

if $n < 0$ **then**

$X \leftarrow 1/x$

$N \leftarrow -n$

else

$X \leftarrow x$

$N \leftarrow n$

end if

while $N \neq 0$ **do**

if N is even **then**

$X \leftarrow X \times X$

$N \leftarrow N/2$

else $\{N$ is odd $\}$

$y \leftarrow y \times X$

$N \leftarrow N - 1$

end if

end while

Appendix B

Codes

B.1 Sample Code

We use this code to find out...

```
1 #include <stdio.h>
2 int Fibonacci(int);
3
4 main()
5 {
6     int n, i = 0, c;
7
8     printf("Enter_the_value_of_n:_");
9     scanf("%d",&n);
10
11     printf("\nFibonacci_series\n");
12
13     for (c = 1 ; c <= n ; c++)
14     {
15         printf("%d\n", Fibonacci(i));
16         i++;
17     }
18
19     return 0;
20 }
21
22 int Fibonacci(int n)
23 {
```

```
24  if (n == 0)
25      return 0;
26  else if (n == 1)
27      return 1;
28  else
29      return (Fibonacci(n-1) + Fibonacci(n-2));
30 }
```

B.2 Another Sample Code

```
1 SELECT associations2.object_id, associations2.term_id,
2      associations2.cat_ID, associations2.term_taxonomy_id
3 FROM (SELECT objects_tags.object_id, objects_tags.term_id,
4      wp_cb_tags2cats.cat_ID, categories.term_taxonomy_id
5 FROM (SELECT wp_term_relationships.object_id,
6      wp_term_taxonomy.term_id, wp_term_taxonomy.term_taxonomy_id
7 FROM wp_term_relationships
8 LEFT JOIN wp_term_taxonomy ON
9      wp_term_relationships.term_taxonomy_id =
10     wp_term_taxonomy.term_taxonomy_id
11 ORDER BY object_id ASC, term_id ASC)
12 AS objects_tags
13 LEFT JOIN wp_cb_tags2cats ON objects_tags.term_id =
14     wp_cb_tags2cats.tag_ID
15 LEFT JOIN (SELECT wp_term_relationships.object_id,
16     wp_term_taxonomy.term_id as cat_ID,
17     wp_term_taxonomy.term_taxonomy_id
18 FROM wp_term_relationships
19 LEFT JOIN wp_term_taxonomy ON
20     wp_term_relationships.term_taxonomy_id =
21     wp_term_taxonomy.term_taxonomy_id
22 WHERE wp_term_taxonomy.taxonomy = 'category'
23 GROUP BY object_id, cat_ID, term_taxonomy_id
24 ORDER BY object_id, cat_ID, term_taxonomy_id)
25 AS categories on wp_cb_tags2cats.cat_ID = categories.term_id
26 WHERE objects_tags.term_id = wp_cb_tags2cats.tag_ID
27 GROUP BY object_id, term_id, cat_ID, term_taxonomy_id
28 ORDER BY object_id ASC, term_id ASC, cat_ID ASC)
29 AS associations2
30 LEFT JOIN categories ON associations2.object_id =
```

```
31         categories.object_id
32 WHERE associations2.cat_ID <> categories.cat_ID
33 GROUP BY object_id, term_id, cat_ID, term_taxonomy_id
34 ORDER BY object_id, term_id, cat_ID, term_taxonomy_id
```


Generated using Undergraduate Thesis L^AT_EX Template, Version 1.4. Department of
Computer Science and Engineering, Bangladesh University of Engineering and
Technology, Dhaka, Bangladesh.

This thesis was generated on Sunday 25th November, 2018 at 4:38am.