

sg

October 15, 2014

Abstract

1 Introduction

2 Bare Essentials

This report will demand that the reader is on the same page vis-a-vis a few terminologies. Trading off space and brevity for clarity: a)Entity b)Entity Pair c)Relation d)Mention e)Match f)Extraction

3 Problem

The next few sections lay the foundation for discussion of our solution.

4 Snowballing

Let us motivate the idea by considering the following related problem:

Suppose we want to populate the repository of founders of companies, and all that we know for a fact is that Elon Musk is the founder of SpaceX. The problem can be divided into two parts, each of them rely on an intuition about how human beings form sentences in general.

- The first of them is given an entity pair, and a corpus of documents, find out all the sentences that express a relation between the entity-pair.

Command line ninjas will quickly think of the following solution: `grep -i 'entity1' sentences—grep -i 'entity2'`

The intuition behind this perhaps the most obvious solution is that *a sentence that houses both the entities can be expected to express a relation between them*. A quick web search with the query “entity1” and “entity2” will show that this intuition is not out of the blue.

- Sentence structure depends on the relation being expressed. In verbose, if two sentences express the same relation, there will be (okay, there can be expected to be) many *features* that are similar in both of them. These include POS tags, words around the entities, dependency path between the entities to name a few.

Putting together the intuitions above, we can solve the problem as follows: Collect all the sentences which have SpaceX and Elon Musk in them, extract features from these sentences. Favor those features which repeat. Now use this set of features to extract similar pairs from other sentences. A fancier solution would be to re use the extractions to learn more features, and continuing the process till the point of diminishing returns.

This seemingly shaky method actually works [1] and is popular by the name of snowballing.

5 Distant Supervision: Snowball scaled up

5.1 Introduction

If the idea of Snowball looks convincing, Distant supervision should follow naturally. Earlier, we considered only one relation for one entity pair. Scale the amount of both of these up and we have distant supervision. The basic setup is as follows:

- a) KB : A knowledge base consisting of facts. The facts are 3-tuples; the entities and the corresponding relation. For example:

Entity	Entity	Relation
Donald Knuth	Wisconsin	Born In
Srinivasa Ramanujan	Erode	Born In
Alan Turing	London	Born In
Alon Musk	SpaceX	Founder Of

has 4 different facts

- b) Corpus The repository of text where we expect to find the sentences that express facts that we know. We need another repository, called the test set, where we will run our extractor to obtain new facts. These two can be the same.

5.2 Matching

We next need to align our knowledge base with the corpus. This process is also called matching.

Data: Corpus C, Knowledge Base KB

Result: Training data, D, A set of matches

Break C into a set of sentences, S;

for each sentence s in S **do**

 let E = all entity pairs in s ;

for each entity pair (e_1, e_2) in E **do**

if \exists relation r in KB with $r(e_1, e_2)$ **then**

 add s to D with label r

end

end

end

Algorithm 1: Distant Supervision

5.3 Training

Recall that obtaining the sentences which express a relation gives us training data, which we want to use to learn relation extractors, our goal. There are several ways to achieve this, starting from the naive ways of training sentence level classifier extending to fancier graphical model based learning.

We briefly discuss the different training methods as we look at a survey of works on Distant Supervision so far.

6 Survey

The first distant supervision paper came out in 1999. Since then, almost every knob that could be twisted in the ds machinery, has been twisted. Different types of relations and different types of datasets will pose different challenges of course, and this report deals with some of them.

This section tries to sketch the guideline.

- 1. Introduced DS
- 2. DS Assumption
- 3. Sentence level Naive Bayes Classifier

Craven and Kumlen, 1999

Distance supervision for the web

Mintz. et al 2009

Increase the number of relations

Learning 5000 relation extractors [LUCHS]

- 1. 102 Freebase relations + Wikipedia
- 2. Distance Supervision Assumption
- 3. Entity Pair Focussed Approach, not sentence level extraction
- 4. Multiclass LR

Relax the DS Assumption

Riedel et. al. 2010

More than one relation for an entity pair

Surdeanu et. al 2012 [Stanford MIML]

Handling large number of False negatives in automatic labeling

Min et. al 2013 [DS with incomplete KB]

- 1. Similar motivation and model as MultiR, differs in training

- 1. Builds up on the graphical model of Riedel et. al
- 2. Let an entity pair participate in more than 1 relation (Multi R)
- 3. Online training, approximate expectation by max, Reduce inference to known problems.

Hoffman et al 2011 [MultiR]

More than one relation for an entity pair

Model missing data

Ritter et. al. 2013 [DNMAR]

Smarter entity detection, type constraints while extracting

Koch et al 2014

- 1. Hitherto, the named entity matching has been adhoc. The authors explore coreference resolution and named entity disambiguation.
- 2. During relation extraction, coarse type constraints are imposed to improve precision

The following 2 assumptions
a. If a relation does not exist in knowledge base, then there won't be any sentence which expresses it.
b. If a relation exists in knowledge base, there will be atleast one sentence that expresses it
Lead to false negatives and false positives respectively, the work relaxes these two assumptions

- 1. Focus is on extracting a large number of relations
- 2. Handles the problem of sparsity by learning lexicons for each relation from weblists
- 3. Hints at hierarchial extraction, first a classifier is trained that prunes out pages which are likely to contain a relation,
- 4. For each attribute a linear chain CRF is trained, relation extraction is treated as a sequence labelling problem.
- 5. For numerical relations, include a feature that models closeness.
- 6. We can exploit their idea of first finding out the high level classes.

- 1. Relation holds in atleast one of the sentences that contain the entity pair.
- 2. Also allow sentence level extraction

- 1. Authors argue that a large number of labeled examples are false negatives
- 2. Algorithm that learns from only positive and unlabeled labels at the entity-pair level

7 MultiR

For our experiments

8 Numerical Relation extraction using MultiR

9 L’homme propose, et Dieu dispose

Our life would have been simpler if MultiR would do the desired, extract the numerical relations from the corpus provided information about some of the numerical relations. But as it happens, it was not meant to be. The model not doing well may mean several things. Checking sanity of the training data was the first among them. Turns out that vanilla distant supervision will lead to an **unprecedented** amount of false positives, which are the root cause of everything that went wrong in the distant supervision pipeline.

9.1 Numbers and False positives

Why would numbers lead to false positives in the first place? The problem stems from the fact that numbers have no identity of their own; they represent count of some real entity or phenomenon.

- **Numbers can appear in many more contexts with an entity** The number of ways in which any two entities can appear together in a sentence is far less than the number of ways in which a number and a quantity can appear together. For example, Consider the entity pair “Bill Gates” and “Microsoft” and the entity-number pair “Bill Gates” and “3” (say). While former will usually co-occur in finite contexts (Founder, CEO, Evangelist etc.), the latter may co-occur anywhere Bill Gates happen to be around something which is 3, the number of cars, billion dollars donated, number of units headed, position in the company, number of business units shutdown by Microsoft and so on.
- The situation is worse for smaller whole numbers, which are more frequent. This intuitively makes sense as we are more often see 2,3 or 11 than 111212233 or 11.42143.
- **The match mines**
During the initial phases of our experiments, we stumbled upon the *match mines* These were basically huge tables, world cup scores of all the matches played and so on. A couple of such sentences were responsible for 21% of the matches! It is easy (and very important) to get rid of such sentences. For subsequent runs, we first sort the sentences by length and then remove top 1000 of them.

The first point of debugging a model which is not performing as well as it should is doing sanity check on the training data.

10 Fighting false Positives with units

11 Numbers are weak entities: a case for keywords

12 Results

13 Possibilities

References

- [1] <http://dl.acm.org/citation.cfm?id=336644>Agichtein, Eugene, and Luis Gravano. "Snowball: Extracting relations from large plain-text collections." Proceedings of the fifth ACM conference on Digital libraries. ACM, 2000.
- [2] Mintz, Mike, et al. "Distant supervision for relation extraction without labeled data." Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 2-Volume 2. Association for Computational Linguistics, 2009.