# A Detailed Analysis of Core NLP for Information Extraction

Simran Kaur[1] and Rashmi Agrawal[2]

[1,2]Faculty of Computer Applications,

Manav Rachna International University, Faridabad, Haryana

[1]simran.fca@mriu.edu.in [2]rashmi.fca@mriu.edu.in

*ABSTRACT*

*The amount of unstructured text present in all electronic media is increasing periodically day after day. In order to extract relevant and succinct information, extraction algorithms are limited to entity relationships. This paper is compendium of different bootstrapping approaches which have their own subtask of extracting dependencies like who did, what, whom, from natural language sentence. This can be extremely helpful in both feature design and error analysis in application of machine learning to natural language processing.*

## 1. Introduction

Bootstrapping is a statistical approach which is used for extracting large amount of information from un-annotated seeds and corpus (text data). For bootstrapping we need large amount of training data which is annotated in a general fashion to extract relevant features of data. Bootstrapping builds effective entities from domain corpus by applying different learning rules. Several algorithms like Basilisk algorithm have been developed to acquire semantic lexicons automatically or semi-automatic tasks, including information extraction, anaphora resolution, question answering, and prepositional phrase attachment. Although some semantic dictionaries like word net exists, these resources do not contain the specialized vocabulary for specific domains. The architecture for information extraction in a pipelined fashion organization is shown below:

Surface or raw text

```
┌─────────────────────────┐
│ TEXT SEGMENTATION       │
│ AND TOKENIZATION        │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│ POS (PARTS OF SPEECH    │
│ TAGGER) NLTK            │
└─────────────────────────┘
```

List of strings

```
┌─────────────────────────┐
│ ENTITY RECOGNITION      │
│                         │
└─────────────────────────┘
```

List of tuples

```
┌─────────────────────────┐
│ RELATION DETECTION      │
│                         │
└─────────────────────────┘
```
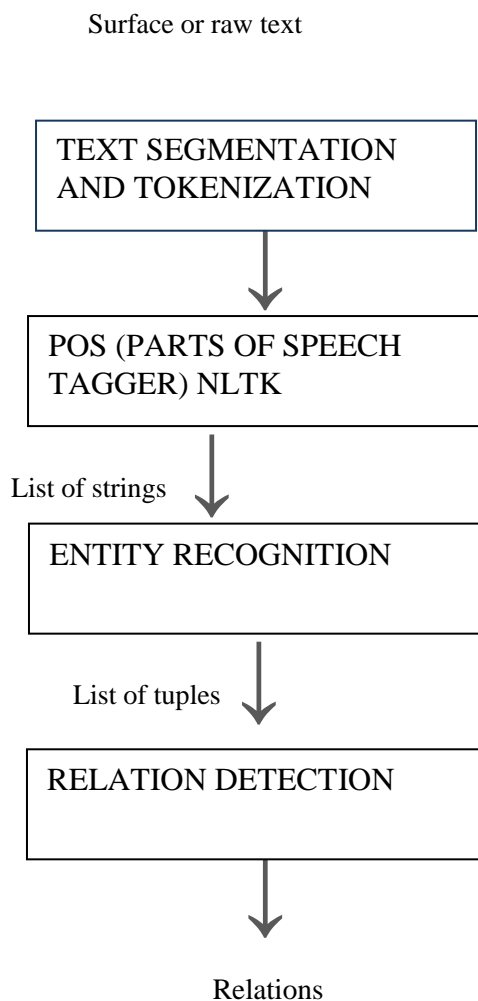
Relations

Figure 1: Pattern or information extraction

To perform the function of information extraction, need a natural language tokenizer (NLTK) which has three parts: sentence tokenizer, word tokenizer and part of speech tagging. Entity detection or extraction is done using this tokenizer but it can be done with bootstrapping techniques as well. e.g.: Infosys is located in Bengaluru. The tuples generated in this case is: ([org: 'Infosys'] 'in' [loc: 'Bengaluru']).

This paper reviewed wide variety of techniques involved in information extraction and problems they solve.

## 2. Confidence Scoring Algorithms

The machine learning algorithms are reasonably successful when enough annotated and labeled data is given. E.g.: clustering algorithm works great on tagged data and labeled seeds but tagging is a tedious job. So, for this purpose to learn from large pool of un-annotated data, bootstrapping learning algorithms are proposed where they usually proceed with seeds to perform labelling of data.eg: milk and juice come under the semantic category of beverages and can be used with verb phrases drank or imbibed. The important characteristic used here is that every time the output of previous iteration is fed as input to the next iteration which grow more seeds.

This algorithm works as follows:

i.    Start with null list of seeds.

ii.   Initialize your list with selected seeds,

iii.  Leverage things in the list and find more things from the corpus

iv.   Score newly found words and add them to the main list

v.    Go to step 2

vi.   Stop after sum fixed epoch values after some 100 iterations or some other stop condition.

In order to extract facts which are relevant to corpus, Riloff et.al [8] considered KNOWITALL system by Etzioni et.al. [5] where an alternative approach is applied for extracting entities i.e. extracting patterns from rules and vice versa. Major weakness of the simple algorithm was tendency to decrease in precision over time as no labelling was present initially. The simple bootstrapping just extracted as many entities without any labelling. Here confidence is score of a particular pattern in a sentence.

The next section reviewed various approaches to bootstrapping algorithms for semantic categorization as well as subjective extraction of noun phrases by using the weighted score of the words or tokens generated.

2.1. Nomen Algorithm

Lin et.al [9], (2003) investigated bootstrapping on un supervised learning of semantic classification. It basically checks whether two different algorithms have same result on the problem or not. It works on biomedical terms like 'mad cow disease' or 'Ebola' virus which often lack capitalization. It is very often that both the diseases and symptoms refer to same thing leading to ambiguity in categorization.

The first step of the algorithm is sentence preprocessing i.e. tokenizer or a lemmatizer is needed to extract tokens and then part of speech tagging for which they tag tokens according to their use in the sentence.

Window size of a noun phrase or a tag is taken ranging from 3 to 8 i.e. the number of words in a noun phrase. For each of the tag generate patterns like this :< disease>mad cow disease </disease> and the window size is three here.

For every learner generate here pos (p), neg (p) and unk (p) where pos (p) are the accepted words present in the pool and neg (p) are the new words accepted from other pools and unk (p) are the unknown words in the pool. The generalized pattern which is followed is (adj+noun) to determine the boundary conditions.

Here accuracy is calculated by given formula:

Confidence = pos (p)-neg (p) /pos (p) +neg (p) +unk (p)

If accuracy<precision, then the words are removed and only top n patterns are accepted.

## 2.2. Basilisk Algorithm

Basilisk is a supervised bootstrapping algorithm proposed by Thelen and Riloff et.al [8], (2002) where input to the system is un-annotated text that automatically generates semantic lexicons The input to Basilisk is an unannotated text corpus and a few manually defined seed words for each semantic category. As the process begins an extraction pattern learner is carried out over the corpus to extract semantic words or noun phrases. Before bootstrapping begins, an extraction pattern learner over the corpus is reviewed which generates patterns to extract every noun phrase in the corpus. The process begins by selecting a subset of the extraction patterns that tend to extract the seed words which is called the pattern pool. The nouns extracted by these patterns become candidates for the lexicon and are placed in a candidate word pool. Basilisk scores each candidate word by gathering all patterns that extract it and measuring how strongly those contexts are associated with words that belong to that pattern of that semantic category. The five best candidate words are added to the lexicon, and the process starts over again.

The input to Basilisk is a text corpus and a set of seed words. The seed words are generated by sorting the words in the corpus by frequency and manually identifying the 10 most frequent nouns that belong to each category. These seed words form the initial semantic lexicon. This section we describe the learning process for a single semantic category.

To identify new lexicon entries, Basilisk relies on extraction patterns to provide contextual evidence that a word belongs to a semantic class. As our representation for extraction patterns,

used the Auto Slog System (Riloff, 1996). Auto Slog'-TS is an extension of auto slog system used earlier in bootstrapping algorithms which generates extraction pattern and then for evaluation using statistical formula which gives primary extraction patterns that represent linguistic expressions that extract a noun phrase in one of three syntactic roles: subject, direct object, or prepositional phrase object. For example, three patterns that would extract people are: "was arrested", "murdered", and "collaborated with". Extraction patterns represent linguistic contexts that often reveal the meaning of a word by virtue of syntax and lexical semantics. Extraction patterns are typically designed to capture role relationships. For example, consider the verb "robbed" when it occurs in the active voice. The subject of "robbed" identifies the perpetrator, while the direct object of "robbed" identifies the victim or target. Before bootstrapping begins, they run Auto Slog exhaustively over the corpus to generate an extraction pattern for every noun phrase that appears. The patterns are then applied to the corpus and all of their extracted noun phrases are recorded. For example, three patterns that would extract people are: "was arrested", "murdered", and "collaborated with". Extraction patterns represent linguistic contexts that often reveal the meaning of a word by virtue of syntax and lexical semantics. Extraction patterns are typically designed to capture role relationships. For example, consider the verb "robbed" when it occurs in the active voice. The subject of "robbed" identifies the perpetrator, while the direct object of "robbed" identifies the victim or target. Before bootstrapping begins, Auto Slog runs exhaustively over the corpus to generate an extraction pattern for every noun phrase that appears. The patterns are then applied to the corpus and all of their extracted noun phrases are recorded.

PR (REL TEXT/TEXT CONTAINING PATTERN (i)) =rel-freq (i)/total-freq (i)

- Where rel-freq(i) = frequency of the relation pattern extracted
- Total-frq(i) = frequency of total patterns present

The algorithm works as follows:

1. Reproduce all extraction patterns as well as pattern evaluation; lexicon = {seed words} i: = 0
2. Score all extraction patterns
3. Pattern pool = top ranked 20+i patterns
4. Candidate word pool = extractions of patterns in pattern     pool
5. Score candidate words in candidate word pool
6. Add top 5 candidate words to lexicon
7. i: = i + 1
8. Again, Go to Step 1.

Table1 : patterns  extracted by autoslog-ts

| Best pattern | "downed in <x>"   (F=3, N=6) |
|---|---|
| Known locations | nicaragua, san miguel*, city |
| New locations | area, usulutan region, soyapango |
| Best pattern | "to occupy <x>"   (F=4, N=6) |
| Known locations | nicaragua, town |
| New locations | small country, this northern area,san sebastian neighborhood, private property |

## 2.3. Snowball Algorithm

Snowball algorithm by Agichtein and Gravano[3] , (2000) tell entities and its relationships from unstructured text. In this algorithm relation is a table which maps different entities.eg: university and location are two entities and their relation can be depicted in table given below.

Table 2: Seed tuples provided to Snowball.

| University | Location |
|---|---|
| Rochester university | Rochester,ny |
| Manav rachna university | Surajkund delhi |
| Delhi university | Delhi |

Snowball which is built on top of the DIPRE (Dual Iterative Pattern Relation Expansion) method by Brin [2], (1999) works on the ides of pattern relation duality, that a good pattern will have good tuples present in it and vice versa following an alternative approach. The only user-provided input to the Snowball system is the training corpus and a small handful of manually compiled relations (the seeds).  Then it searches for that particular tuple or the tuples in similar proximity as well as analyzes the text around it.

The algorithm generates the  pattern on basis of a 5-tuple system <left, tag1, middle, tag2, right> where left, right and middle words are the weights of the words. tag1 and tag2 are the tuple names university and location as mentioned above.

## 3. Data Collection

Many NLP related tasks like named entity recognition, subjectivity classification, text classification and word sense disambiguation rely on un-annotated data which is not tagged. Worldwide web has become a popular choice of un-annotated data which is not tagged. Worldwide web has become a popular choice as a source of un-annotated data. All the NLP tasks are performed using unstructured data on worldwide web for supervised and weakly supervised algorithms. The use of un-annotated text for information extraction uses supervised algorithms and in bootstrapping methods (Riloff,1996) and bootstrapping algorithm begins with seed words.

Previously existing bootstrapping algorithm's uses domain oriented corpus for information extraction but now they are also using MUC-4 training text. Meta bootstrapping method was also trained on web pages where domain specific web pages like corporate relationships exist. The know it all system also uses un annotated web pages for extraction of information using domain independent relationships extracting ontology.

## 3. Methodology

NLP known as Natural Language Processing by Bebo White[11],(2014) is a technique of Artificial Intelligence for extracting gist out of the corpus of data and converting words in data.

| PROPERTY NAME | DESCRIPTION |
|---|---|
| TOKENIZE | It Tokenizes the text as well as handle noisy data. It saves character offsets. By adding Character Offset Begin Annotation and Character Offset End Annotation. |
| CLEANXML | It removes all the extensible mark-up language tokens. |
| SSPLIT | It splits your sequence of words to sentences |
| POS | It applies part of speech tagging |
| LEMMA | It lemmatizes all the tokens in the corpus |
| NER | It Recognizes named entities from the Corpus such as name, temporal and numerals |
| REGEXNER | Implements a simple, rule-based NER over token sequences using Java regular expressions. The goal of this Annotator is to provide a simple framework to incorporate NE labels that are not annotated in traditional NL corpora. For example, the default list of regular expressions that we distribute in the models file recognizes ideologies (IDEOLOGY), nationalities (NATIONALITY), religions (RELIGION), and titles (TITLE). Here is a simple example of how to use RegexNER. For more complex applications, you might consider TokensRegex |
| SENTIMENT | It implements Socher et al sentiment model Where it uses a binarized tree. |
| TRUECASE | It recognizes upper case of the sentences using crf sequence tagger. |
| PARSE | Provides full syntactic analysis, using both the constituent and the dependency representations. The constituent-based output is saved in Tree Annotation. it generates three dependency-based outputs, as follows: basic, un collapsed dependencies, saved in Basic Dependencies Annotation; collapsed dependencies saved in Collapsed Dependencies Annotation; and collapsed dependencies with processed coordination's, in Collapsed CC Processed Dependencies Annotation |
| DEPARSE | Provides a fast-syntactic dependency parser. it generates three dependency-based outputs, as follows: basic, un collapsed dependencies, saved in Basic Dependencies Annotation; collapsed dependencies saved in Collapsed Dependencies Annotation; and |

| | |
|---|---|
| | collapsed dependencies with processed coordination's, in Collapsed CC Processed Dependencies Annotation. Most users of our parser will prefer the latter representation |
| DECOREF | It implements nominal coreference resolution where entire paragraph is stored in coref chain annotation. |
| RELATION | It finds relation between two entities which is trained on relation types. Eg live in, located in, org based in, work for and none. |
| NATLOG | It marks quantifier scope and token polarity with respect to natural language semantics. e.g. logic operators. |
| QUOTE | Deterministically picks out quotes delimited by "or 'from a text. All top-level quotes, are supplied by the top-level annotation for a text. If a Quotation Annotation corresponds to a quote that contains embedded quotes, these quotes will appear as embedded Quotation Annotations that can be accessed from the Quotation Annotation that they are embedded in. The Quote Annotator can handle multi-line and cross-paragraph quotes, but any embedded quotes must be delimited by a different kind of quotation mark than its parents |

The important packages of the system are described below in detail:

i.Tokenization: A Tokenizer by Bebo White[11], 2014  follows tokenization approach in natural language understanding by splitting a string into its sub-classes. It provides a lexical scanner that handles all the operators and delimiters.
Eg: A tokenizer tokenizes sentence into its morphological forms as given: "my name is Simran" is tokenized to 'my','name','is','Simran' which are the individual tokens generated.

ii.Stemming: The idea of stemming by Bebo White[11],2014 in natural language processing is a sort of normalizing method. A word can vary and have different variations.
 Eg:  I was studying my chapter.
    I studied my chapter.
    I study.
So verb study has three variations above which can be stemmed to its root or canonical form. Hence studying is brought to its root form study which is shown below in the illustration.
Studying---Study  It can be either prefix stemming or suffix stemming which brings word to its root form. So it is a processing interface that removes suffixes and prefixes from the word and bring it down to its affix form.

1. Porter Stemmer Algorithm: This algorithm removes suffixes from the word forms.
Eg: studying--study
2. Lancaster Stemmer: A word stemmer that is based on lancaster stemmer algorithm.
3. Regex Stemmer: This word stemmer used regular expressions for converting into morphological affixes.
Eg: bars--bar
iii. Segmentation:
Sentence tokenizer or Sentence disambiguation by Bebo White[11],2014 is also known as sentence breaking or segmentation into its constituent words or tokens.
Eg: "My name is Simran. I live in India." is segmented to two sentences : My name is Simran, I live in India.

iv. Collocation:

Collocations are the words that occur commonly in same context and frequently. For example, the top ten bi-gram collocations in Genesis are listed below, as measured using Point-wise Mutual Information.

Collocation by Bebo White[11],2014 or lexical collocation means two or more words co-occur in a sentence more frequently than by chance. A collocation is an expression that forms a specific meaning. It may be noun phrase like large villa, verbal phrase like go down, idioms, or technical terms. Collocations are defined by constricted compositionality, that is, it is difficult to predict the meaning of collocation from the meaning of its parts. For example, He is known for his fair and square dealings and everybody trusts his work.Here fair and square means honest but if we take the individual words though the word fair gives somewhat closer meaning as it means just the word square confuses us. So instead of taking individual words one should take the collocation fair and square and find meaning. It shows that collocations play a key role in understanding sentences. Collocations are recursive  in nature  so they may contain more than two words in a sentence.

 v. Tagging: Parts of Speech Tagging by Bebo White[11],2014 means assigning lexical categories to words whether it is a noun phrase, verb phrase or prepositional phrase.

Eg: Cat chases a rat. Here Cat is assigned NN lexical category(noun);Chases is assigned VB lexical category(verb); NN is assigned lexical category(noun).

vi. Parsing: Parsing by Bebo White[11],2014 or Analyzing the Syntactic Structure of the sentence using Context Free Grammar.

A grammar is said to be recursive in nature only when a lexical category that is present on the left handside of the production rules appear on the right hand side of the production rules.The production Nom -> NP Nom (where Nom is the category of nominals) involves direct recursion on the category Nom, whereas indirect recursion on S arises from the combination of two productions, namely S -> NP DET and DET-> N S.
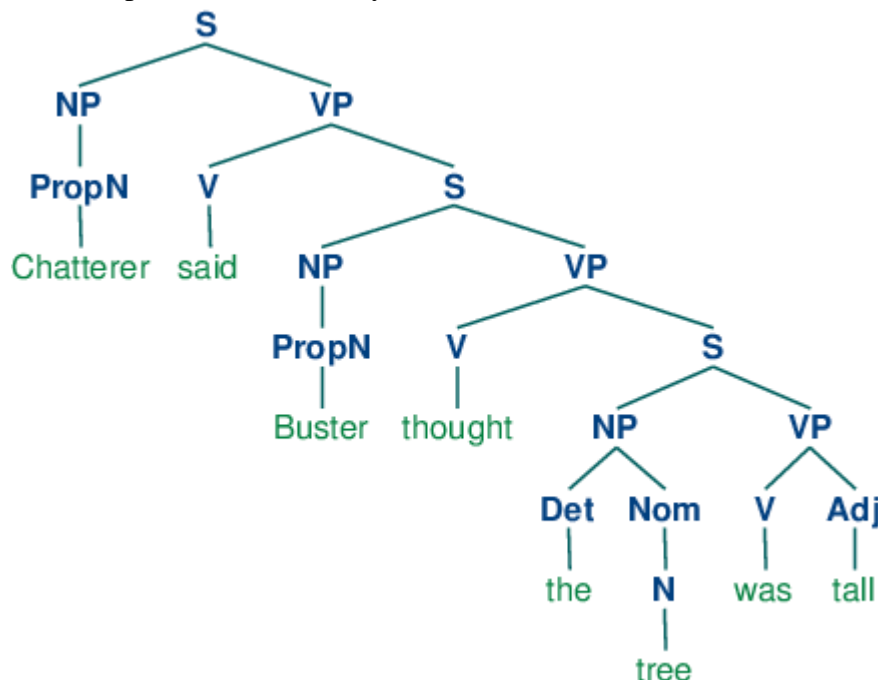


Figure 2: Tree Structure

vii. Word Sense Disambiguation: Word Sense Disambiguation in NLP by Bebo White[11],2014 is a grammatical dilemma that handles ubiquitous ambiguity .

Eg: in the sentence :"While hunting in India, I shot a deer in my pajamas. How he got into my pajamas, I don't know."
Here we can see ambiguity in the phrase "shot a deer " that whether the person shot with an gun or camera
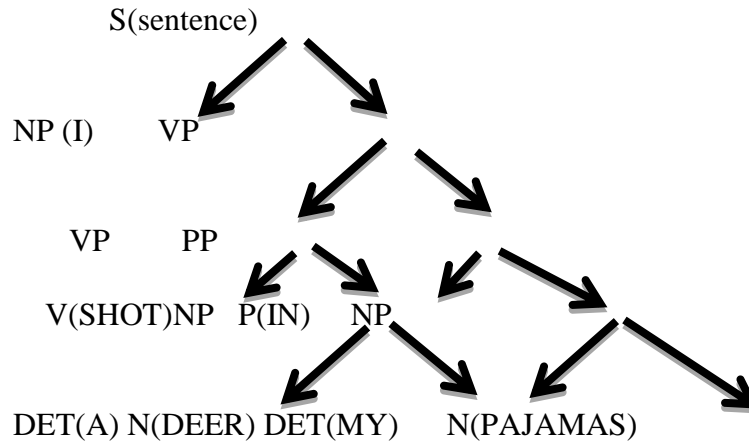


Figure 3: Tree Structure

In the above syntax the parse tree for the sentence is generated which signifies whether prepositional phrase point towards shooting with camera or shooting event hence removing ambiguity.

Henceforth all the techniques when are combined together and applied on big data-set or corpus of data for machine translation, speech recognition gives us logjam of fruitful results.

NLP is importantnt in all the field i.e. scientific, educational, medicinal, corporate and cultural fields. NLP is experiencing fast maturation as its belief and method acting are deployed in a variety of new language technologies. For this reason it is important for a wide range of people to have a working cognition of NLP. Within industry, this includes people in human-computer fundamental interaction, business information reasoning, and web software alteration.  So basic goals of NLP are language analysis and language technology (Erik Cambria,2014).

In language analysis data modelling, text mining and knowledge representation techniques are applied while in language technology statistical algorithms are used deploying data structures.

In order to generate enhanced dependencies, syntactic structure of the sentences, named entities and relations between different entities and test this natural language analysis tools on a corpus of data given below and get the following results or our expected outcomes:

## 5. Pipelined Working Of The Tool

The annotator pipeline of the tool works as follows:

i.Parts of speech tagging:

It is a software that is applied for pos tagging. This tool is implemented in java platform. It generates fine grain level tags in noun plural form.
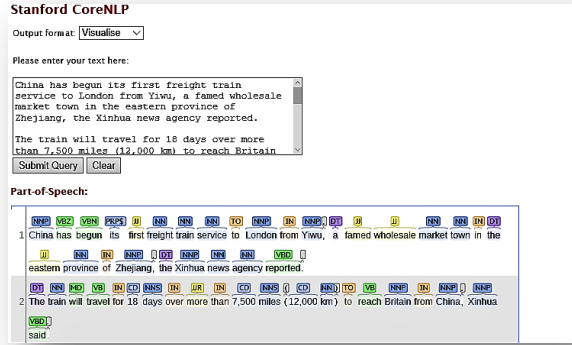


Fig 2: part of speech tagger

ii. Named entity recognition:

It recognizes or extracts basically named, numerical and temporal entities and normalize them.
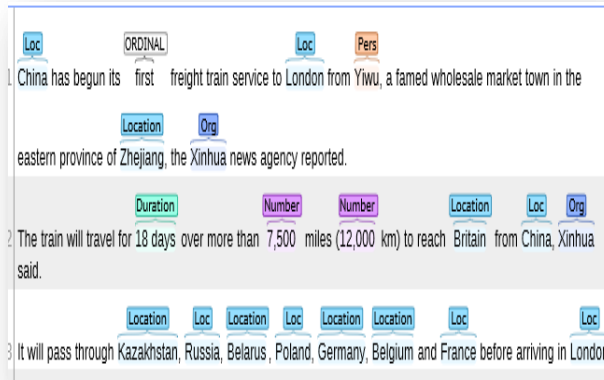


Fig 3: named entity recognition

iii. Co-Reference resolution

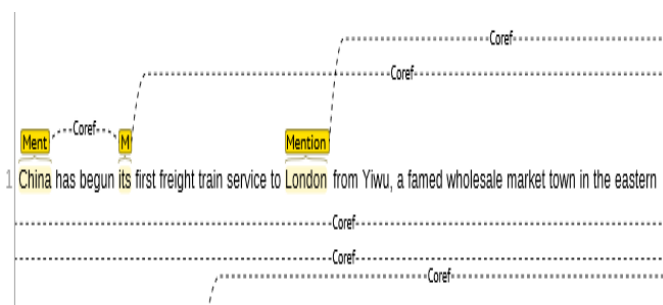It generated headwords and their f1 measures by using statistical methods.



Fig 4: co-reference resolution

iv. Dependencies:

It generates relation between entities using relationship annotators in java e.g.: live-in, work for, located etc.
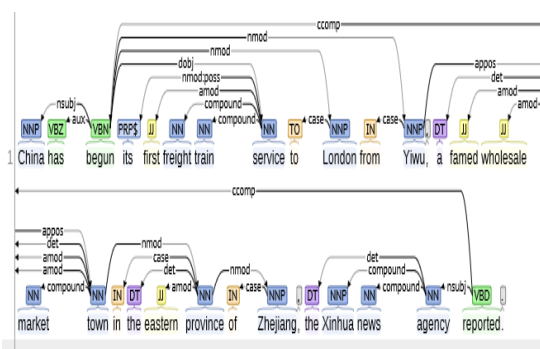


Fig 6: dependencies

Corpus: Stanford Core NLP provides a set of natural language analysis tools. It can give the base forms of words, their parts of speech, whether they are names of companies, people, etc., normalize dates, times, and numeric quantities, mark up the structure of sentences in terms of phrases and word dependencies, indicate which noun phrases refer to the same entities, indicate sentiment, extract particular or open-class relations between entity mentions, get quotes people said, etc.

Stanford Core NLP's goal is to make it very easy to apply a bunch of linguistic analysis tools to a piece of text. A tool pipeline can be run on a piece of plain text with just two lines of code. Core NLP is designed to be highly flexible and extensible. With a single option you can change which tools should be enabled and which should be disabled. Stanford Core NLP integrates many of Stanford's NLP tools, includes, parts of speech tagger(POS) the named entity recognizer (NER), the parser, the co-reference resolution system, sentiment analysis, bootstrapped pattern learning, and the open information extraction tools. Moreover, an annotator pipeline can include additional custom or third-party annotators. Core NLP's analyses provide the foundational building blocks for higher-level and domain-specific text understanding applications.

Sentence wise tokens generated and dependencies between the words in the sentence using tool.

1.Sentence no1:

44

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 1 | Stanford | Stanford | 0 | 8 | NNP | ORGANIZATION | | PER0 | |
| 2 | CoreNLP | CoreNLP | 9 | 16 | NNP | O | | PER0 | |
| 3 | provides | provide | 17 | 25 | VBZ | O | | PER0 | |
| 4 | a | a | 26 | 27 | DT | O | | PER0 | |
| 5 | set | set | 28 | 31 | NN | O | | PER0 | |
| 6 | of | of | 32 | 34 | IN | O | | PER0 | |
| 7 | natural | natural | 35 | 42 | JJ | O | | PER0 | |
| 8 | language | language | 43 | 51 | NN | O | | PER0 | |
| 9 | analysis | analysis | 52 | 60 | NN | O | | PER0 | |
| 10 | tools | tool | 61 | 66 | NNS | O | | PER0 | |
| 11 | . | . | 66 | 67 | . | O | | PER0 | |

*Parse tree*
(ROOT (S (NP (NNP Stanford) (NNP CoreNLP)) (VP (VBZ provides) (NP (NP (DT a) (NN set)) (PP (IN of) (NP (JJ natural) (NN language) (NN analysis) (NNS tools))))) (. .)))

*Uncollapsed dependencies*

- root ( ROOT-0 , provides-3 )
- compound ( CoreNLP-2 , Stanford-1 )
- nsubj ( provides-3 , CoreNLP-2 )
- det ( set-5 , a-4 )
- dobj ( provides-3 , set-5 )
- case ( tools-10 , of-6 )
- amod ( tools-10 , natural-7 )
- compound ( tools-10 , language-8 )
- compound ( tools-10 , analysis-9 )
- nmod ( set-5 , tools-10 )

*Enhanced dependencies*

- root ( ROOT-0 , provides-3 )
- compound ( CoreNLP-2 , Stanford-1 )
- nsubj ( provides-3 , CoreNLP-2 )
- det ( set-5 , a-4 )
- dobj ( provides-3 , set-5 )
- case ( tools-10 , of-6 )
- amod ( tools-10 , natural-7 )
- compound ( tools-10 , language-8 )
- compound ( tools-10 , analysis-9 )
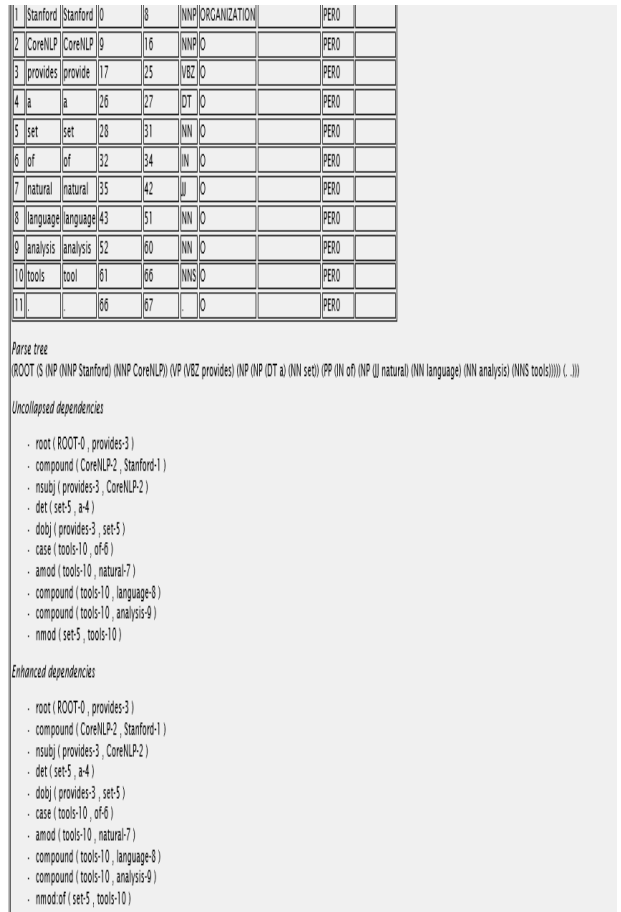- nmod:of ( set-5 , tools-10 )

Figure 7: sentence structure

Here in the above screenshot it inserts corpus of data in Stanford NLP tool and select the option as pretty print. The tokens generated by tool are:

Table 3: Token Generation

| Tokens | Part of speech tagger |
|---|---|
| Stanford | Noun phrase |
| Core | Noun phrase |
| Nap | Noun phrase |
| A | Determiner |
| Set | Noun |
| Of | In(connectors) |
| Natural | Adjective |
| Language | Noun phrase |
| Analysis | Noun phrase |

The root of this corpus is Stanford and core-NLP tokens which are having highest frequency.

# 6. Role of Bootstrapping Algorithms

These algorithms maps specific words and relations by applying different pattern based approaches as they could see in snowball algorithm. Similarly, in basilisk algorithm they used subjective nouns and tag them to classes for semantic categorization as well as subject phrases so basilisk is just not limited to semantic categorization, it is also useful for learning large number of typological syn-sets. The basilisk algorithm can be employed with synergic combination of other bootstrapping algorithms. It is higher level of basilisk approach where semantic categorization of seeds is done. These seeds are feed into the system to extract patterns. The best patterns are selected from which they select centroid of each pattern to form a good lexicon and a fair gold standard for evaluation.

An Exploratory research was carried by Riloff et.al to examine multiple approaches to bootstrapping in natural language processing. Bootstrapping is an extremely powerful approach which extracted good patterns from unstructured language. Different algorithms followed different approaches giving different results and are amenable to to a wide variety of natural language processing tasks, such as semantic categorization, relation extraction, and subjectivity analysis. Bootstrapping algorithms must be provided with proper seeds and just not take improper seeded words imbibing ir-relevant data. The choice of seeds is pivotal to the success of the bootstrapping process and it is not at all clear how to determine what the "best" seed might be. Bootstrapping systems are best suitable to natural language processing tasks due to their ability to learn and navigate the syntactically rich, unstructured, and extremely complex nature of loosely structured natural languages.

# 7.Conclusion

The bootstrapping algorithm exploit two-fold ideas: (1) evidence extraction from patterns used for inferencing different semantic categories and (2) learning of multiple bootstrapping algorithms to extract the semantic categories. In this paper, the result generated by bootstrapping algorithm are tokens and enhanced relations between them using Basilisk algorithm which provides succinct and fair gold summary which has high confidence and accuracy. Ahead multiple semantic categories can be extracted by working on basilisk algorithm and combining it with other bootstrapping algorithms (meta-bootstrapping), in order to improve the results of tagged patterns which would further help in other application areas like biomedical natural language processing, maintaining clinical inventories, providing speech aid to challenged children and machine translation where we convert semantic features of one language to another language.

# 6.References

[1] Agichtein, E., E. Eskin, and L. Gravano (2000). Combining strategies for extracting relations from text collections. In Proceedings of the 2000 ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery (DMKD 2000).

[2] Agichtein, E. and L. Gravano (2000, June). Snowball: Extracting relations from large plain-text collections. In Proceedings of the 5th ACM International Conference on Digital Libraries (DL'00), San Antonio, TX.

[3] Agichtein, E., P. Ipeirotis, and L. Gravano (2003). Modeling query-based access to text databases. Proceedings of the Sixth International Workshop on the Web and Databases, WebDB, 87–92.

[4] Brin, Sergey. 1999. Extracting patterns and relations from the World Wide Web. In The World Wide Weband Databases, pages 172–183. Springer Berlin Heidelberg.

[5] O. Etzioni, M. Cafarella, D. Downey, S. Kok, A. Popescu, T. Shaked, S. Soderland, D. Weld, and A. Yates. Web-Scale Information Extraction in KnowItAll. In Proceedings of the 13th International World Wide Web Conference (WWW-04), pages 100–110, New York City, New York, 2004.

[6] Thelen, Michael and Riloff, Ellen. 2002. A bootstrapping method for learning semantic lexicons using extraction pattern contexts. In Proceedings of the ACL-02 conference on Empirical methods in natural language processing, Volume 10, pages 214– 221. Association for Computational Linguistics.

[7] Yangarber, Roman, Lin, Winston, and Grishman, Ralph. 2002. Unsupervised learning of generalized names. In Proceedings of the 19th international conference on Computational linguistics, Volume 1. Association for Computational Linguistics.

[8] Lin, Winston, Yangarber, Roman, and Grishman, Ralph. 2003. Bootstrapped learning of semantic classes from positive and negative examples. In Proceedings of ICML-2003 Workshop on The Continuum from Labeled to Unlabeled Data, Volume 4, No. 4.

[9] Riloff, Ellen, Wiebe, Janyce and Wilson, Theresa. 2003. Learning subjective nouns using extraction pattern bootstrapping. In Proceedings of the seventh conference on Natural language learning at HLTNAACL, Volume 4, pages 25–32. Association for Computational Linguistics.

[10] Riloff, Ellen and Wiebe, Janyce. 2003. Learning extraction patterns for subjective expressions. In Proceedings of the 2003 conference on Empirical methods.