

# Coupling Semi-Supervised Learning of Categories and Relations

Harsh Patel, Uday Ramesh and Jashandeep Singh Sran

School of Computer Science

Illinois Institute of Technology

Chicago, IL 60616

{hpatel100, uramesh1, jsran}@hawk.iit.edu

## Abstract

Consider how to extract semi-supervised learning information, especially to extract instances of noun categories (such as "athletes" and "teams") and relationships (such as "playsForTeam"). A semi-supervised approach that uses a few tagged and many untagged examples are often unreliable because it often produces a set of internally consistent but incorrect extracts. I have. We propose that this problem can be overcome by simultaneously learning many different categories and related classifiers in the presence of an ontology that defines the constraints that combine the training of these classifiers. Experimental results show that simultaneously learning a combined collection of 3 categories and related classifiers leads to a much more accurate extraction than training the classifiers individually.

## 1 Introduction

A lot of knowledge is expressed in natural language on the web. Converting to a structured knowledge base that contains facts about entities (such as "Disney") and the relationships between them (Company Industry ("Disney", "Entertainment", etc.)) can be very useful in many applications. A fully supervised method for learning to extract such facts from text works well, but the cost of collecting many labeled examples of each type of knowledge extracted is impractical... Researchers also considered semi-supervised learning methods that rely primarily on unlabeled data,

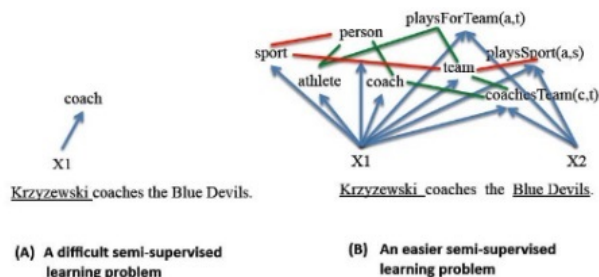


Figure 1: Combined training of information extractors with many related categories and relationships (B) is significantly more accurate than the simple but much more difficult task (A) of learning a single information extractor. It shows that it improves too.

However, these approaches tend to be plagued by the fact that they face less constrained learning tasks, which often results in inaccurate extraction.

By combining the training of many information extractors, we present a semi-supervised learning approach that yields more accurate results. The intuition behind our approach (summarized in Figure 1) is that single extractor-type semi-supervised training such as "coach" covers many interconnected entities and relationship types. It's much more difficult than training the extractors at the same time. In particular, with prior knowledge of the relationships between these various entities and relationships (eg, "coach (x)" means "person (x)", not "sports (x)"). For training, unlabeled data will be much more useful during constraints.

In our previous work, we combined learning from multiple categories or used the static category recognition feature to test the arguments of the learned relationship extraction feature, but our work is to have multiple semi-supervised learnings. This is the first task of recognizing that we will combine category and relationship training at the same time. Our experiments show that this binding leads to a more accurate extraction. Based on the results reported here, it is assumed that the accuracy of information extraction can be significantly improved by combining training with hundreds or thousands of extractors.

## 2 Problem Statement

It is helpful to first explain the use of common terms. An ontology is a collection of unary and binary predicates, also known as categories and relations, respectively. An instance of one category, or a category instance, is a noun phrase. An instance of a relation or relational instance is a pair of noun phrases. Instances can be positive or negative with respect to a particular predicate. That is, the predicate may or may not apply to that particular instance. Promoted instances are the instances that the algorithm considers to be positive instances of the predicate. Also associated with both categories and relationships are patterns. Wildcard strings (eg "Match with arg1" and "Head coach of arg1, arg2"). Promoted patterns are patterns that are supposed to be high-probability indicators of predicates.

The task of this task is reliability, starting with a large corpus of sentences annotated with some positive initial instances and patterns of each predicate, and part of speech, in relation to a particular ontology category. Is to learn an extractor that automatically fills in high instances. Part of speech (POS). keyword. Focus on extracting the facts mentioned several times in the corpus, which can be probabilistically evaluated with the help of corpus statistics. It does not resolve the string to the actual entity. The issue of resolving synonyms and disambiguating strings that can refer to multiple entities remains for future work.

## 3 Related Work

Work on perform various tasks learning has shown that directed learning of different "related" works together can yield higher exactness than learning the capacities independently (Thrun, 1996; Caruana, 1997). Semi-regulated perform various tasks learning has been displayed to increment precision when errands are connected, allowing one to utilize an earlier that empowers comparative parameters (Liu et al., 2008). Our work likewise includes semi-directed preparing of numerous coupled functions, however varies in that we expect express earlier information on the exact manner by which our multiple capacities are connected (e.g., that the upsides of the capacities applied to a similar info are commonly exclusive, or that one suggests the other).

In this paper, we center around a 'bootstrapping' technique for semi-managed learning. Bootstrapping approaches start with a few labeled 'seed' models, utilize those seed guides to prepare an underlying model, then, at that point, utilize this model to label a portion of the unlabeled information. The model is then retrained, utilizing the first seed models in addition to oneself named models. This cycle emphasizes, slowly growing how much-marked information. Such methodologies have shown guarantee in applications, for example, site page arrangement (Blum and Mitchell, 1998), named element characterization (Collins and Artist, 1999), parsing (McClosky et al., 2006), and machine interpretation (Ueffing, 2006).

Bootstrapping ways to deal with data extraction can yield noteworthy outcomes with minimal introductory human exertion (Brin, 1998; Agichtein and Gravano, 2000; Ravichandran and Hovy, 2002; Pasca et al., 2006). In any case, after numerous cycles, they usually experience the ill effects of semantic float, where mistakes in labeling aggregate and the gained idea 'floats' based on what was

expected (Curran et al., 2007). Coupling the learning of predicates by utilizing positive examples of one predicate as bad models for others has been displayed to assist with restricting this float (Riloff and Jones, 1999; Yangarber, 2003). Moreover, ensuring that connection contentions are of sure, expected types can assist with alleviating the advancement of wrong occasions (Paca et al., 2006; Rosenfeld and Feldman, 2007). Our work expands on these plans to couple the synchronous bootstrapped preparing of multiple classifications and different relations.

Our way to deal with data extraction depends on utilizing high accuracy context oriented designs (e.g., 'is the chairman of argl' proposes that argl is a city). An early example-based way to deal with data extraction acquired 'is a' relations from text utilizing conventional contextual designs (Hearst, 1992). This approach was subsequently increased to the web by Etzioni et al. (2005).

Another exploration investigates the undertaking of 'open information extraction', where the predicates to be learned are not indicated ahead of time (Shinyama and Sekine, 2006; Banko et al., 2007), however, arise rather from the examination of the information. Interestingly, our methodology relies emphatically on information in the cosmology about the predicates to be learned, and the connections among them, to accomplish high precision.

Chang et al. (2007) present a structure for discovering that upgrades the information probability in addition to imperative-based punishment terms then catches earlier information and exhibits it with semi-regulated learning of division models. Limitations that catch space information guide bootstrap learning of an organized model by punishing or disallowing infringement of those imperatives. While comparative in the soul, our work varies in that we think about learning many models, instead of one organized model, and that we think about a lot bigger scope application in an alternate area.

## 4 Approach

### 4.1 Coupling of Predicates

As referenced over, our methodology depends on the notion of coupling the learning of numerous capacities to compel the semi-managed learning issue we face. Our framework learns four distinct sorts of capacities. For every class  $c$ :

1.  $f_{c,inst} : NP(C) \rightarrow [0, 1]$
  2.  $f_{c,patt} : Patt_C(C) \rightarrow [0, 1]$
- and for each relation  $r$ :
1.  $f_{r,inst} : NP(C) \times NP(C) \rightarrow [0, 1]$
  2.  $f_{r,patt} : Patt_R(C) \rightarrow [0, 1]$

where  $C$  is the info corpus,  $NP(C)$  is the arrangement of legitimate thing phrases in  $C$ ,  $Patt_C(C)$  is the arrangement of substantial class designs in  $C$ , and  $Patt_R(C)$  is the arrangement of legitimate connection designs in  $C$ . "Legitimate" thing phrases, class examples, and connection designs are characterized in Segment 4.2.2. The learning of these capacities is coupled in two ways:

1. Sharing among same-arity predicates as per legitimate relations.
2. Relation contention type-checking.

These strategies for coupling are made conceivable by earlier information in the info cosmology, past the arrangements of classes and relations referenced previously. We give general depictions of these strategies for coupling in the following areas, while the subtleties are given in segment 4.2.

### 4.2 Algorithm Description

In this segment, we portray our calculation, CBL (Coupled Bootstrap Student), exhaustively. The contributions to CBL are a huge corpus of POS-tagged sentences and an underlying cosmology with pre-characterized classes, relations, fundamentally unrelated relationships between same-arity predicates, subset relationships between certain classifications, seed occurrences for all predicates, and seed designs for the categories. Classes in the information philosophy additionally have a banner demonstrating whether occurrences should be formal people, places or things, normal things, or whether they can be either (e.g., cases of 'city' are formal people, places or things).

<p><b>Algorithm 1: CBL Algorithm</b></p> <p><b>Input:</b> An ontology <math>\mathcal{O}</math>, and text corpus <math>C</math></p> <p><b>Output:</b> Trusted instances/patterns for each predicate</p> <p>SHARE initial instances/patterns among predicates;</p> <p><b>for</b> <math>i = 1, 2, \dots, \infty</math> <b>do</b></p> <p>    <b>foreach</b> predicate <math>p \in \mathcal{O}</math> <b>do</b></p> <p>        EXTRACT candidate instances/patterns;</p> <p>        FILTER candidates;</p> <p>        TRAIN instance/pattern classifiers;</p> <p>        ASSESS candidates using classifiers;</p> <p>        PROMOTE highest-confidence candidates;</p> <p>    <b>end</b></p> <p>    SHARE promoted items among predicates;</p> <p><b>end</b></p>
--

#### 4.2.1 Configuration for the Algorithm

Full (Full Algorithm): - The Algorithm described in the proposed solution like below.

#### 4.2.2 Input

Larger corpus (collection of Texts) of part-of-speech (POS) tagged sentences and an Initial ontology with pre-defined categories, relations, mutually exclusive relationships between same-arity predicates, subsets relationships between same categories, seed instances for all predicates, and seed patterns of categories.

#### 4.2.3 Sharing

Seed instances, and patterns are shared among predicates according to mutual exclusions, subset, and type-checking constraints.

### 4.3 Candidate Extraction

Finds new candidate instances by using newly promoted patterns to bring out the noun phrases that co-occurs with those patterns in text corpus by using a Map-reduce framework.

#### 4.3.1 Category Instances

A noun phrase will be searched for by CBL. The common goal is here to find noun phrase coming from text corpus. To segment for a noun phrase, use POS tags and ignore delimiters. Only if the category's common noun specification is met. We are looking for stop words and capital letters if present in the text corpus. The goal is here to looking for stop words from nltk package and capital letters. If it doesn't contain any of the above mentioned things, will get the common noun from `extract_common_text_phrase()` in `instance_extract_common_text_phrase.py` [in CBLAlgorithm/instance/ package].

#### 4.3.2 Relation Instances

If promoted pattern was found, a candidate relation instances extracted if both placeholders are valid noun phrases.

```

5 def extract_common_text_phrase(roots):
6     if roots is not None:
7         similar_phrase = ''
8
9         # check any stop words are there.
10        obj1 = False
11        # check any capital letters are there.
12        obj2 = False
13
14        for child in roots.leaves():
15            # word is child[0] then,
16            word = child[0]
17
18            # Add word into empty similar_phrase
19            similar_phrase = similar_phrase + word
20
21            # Add "." after the similar_phrase
22            similar_phrase = similar_phrase + ' '
23
24            # If there are no stop words there, obj1 will be false.
25            if not utils.check_stop_word_map.check_stop_word_map(word):
26                obj1 = False
27            else:
28                obj1 = True
29
30            # If there are any capital letters there, obj2 will be true.
31            if utils.has_capital_letters.has_capital_letters(word):
32                obj2 = True
33            else:
34                obj2 = False
35
36        # Strip out similar_phrase if any unnecessary thing.
37        similar_phrase = similar_phrase.strip()
38
39        # If similar_phrase is there and obj1 and obj2 is not, means no stop words are there and no capital letters, then return it.
40        if similar_phrase and not obj1 and not obj2:
41            return similar_phrase

```

To extract noun phrase

```

6 def extract_noun_phrase(slab, grammar, typeCheck):
7     if slab is not None:
8         if grammar is not None:
9             if typeCheck is not None:
10                similar_phrase = ''
11
12                # Create a new chunk parser, from the given start state and set of chunk patterns.
13                regexParser = RegexParser(grammar['EXPRESSION'])
14
15                # words_list is list
16                regexParser_parser = regexParser.parse([str2tuple(obj) for obj in slab.split()])
17
18                for child in regexParser_parser.subtrees():
19                    if (typeCheck == 'proper' or typeCheck == 'all'):
20                        # PN means proper nouns.
21                        if child.Label() == 'PN':
22                            # Extract proper noun.
23                            similar_phrase = instance_extract_proper_noun.extract_proper_text_phrase(child)
24                        if (typeCheck == 'common' or typeCheck == 'all'):
25                            # CN means common nouns.
26                            if child.Label() == 'CN':
27                                # Extract common text text.
28                                similar_phrase = instance_extract_common_text_phrase.extract_common_text_phrase(child)
29
30                return similar_phrase

```

To extract instance candidate instances according to the position, we do have CBLAlgorithm/instance/instance\_extract\_corpus\_type.py file.

### 4.3.3 Category Patterns

If the preceding words are verbs followed by a series of adjectives, prepositions, or delimiters, CBL highlights them as a candidate pattern(e.g. 1, 'being acquired by arg1') or nouns and adjectives followed by a sequence of adjectives, prepositions, or determiners (e.g., 'former CEO of arg1'). (e.g. 2, 'arg1 broke the home run record'), or verbs followed by a preposition (e.g., 'arg1 said that'). To extract instance candidate patterns according to the instances, we do have CBLAlgorithm/pattern/patterns\_get\_extracted\_category\_findings.py file.

```

10 def get_extracted_category_findings(text_corpus, collections):
11     supplement_findings = list()
12
13     # Find regression findings.
14     regex_findings = utils_init_regex_expression.init_regex_expression(collections)
15
16     # Get findings list.
17     findings_list = utils_init_halfway_pattern.init_halfway_pattern(regex_findings)
18
19     # Iterate over all corpus.
20     for idx in text_corpus:
21         # For each corpus, look in candidate obj or not.
22         if utils_check_candidate_text_corpus.check_candidate_text_corpus(idx, regex_findings):
23             # Split the corpus string to the occurrences of the findings, returning a list containing the resulting substrings
24             slabs = re.split(findings_list, idx)
25
26             # Get the length of slabs and find left and right markings.
27             for slab_idx in range(len(slabs) - 1):
28                 LP = ''
29                 RP = ''
30
31                 # Return a copy of the string with leading and trailing whitespace removed.
32                 LS = slabs[slab_idx].strip()
33                 if LS is not None:
34                     LP = pattern_get_extract_terms.get_extract_terms(LS, PATTERN_GRAMMAR_LEFT)
35
36                 # Return a copy of the string with leading and trailing whitespace removed.
37                 RS = slabs[slab_idx + 1].strip()
38                 if RS is not None:
39                     RP = pattern_get_extract_terms.get_extract_terms(RS, PATTERN_GRAMMAR_RIGHT)
40
41                 if LP is not None:
42                     supplement_findings.append(utils_filling_supplement.filling_supplement(LP, 'LEFT'))
43                 elif RP is not None:
44                     supplement_findings.append(utils_filling_supplement.filling_supplement(RP, 'RIGHT'))
45
46     return supplement_findings

```

### 4.3.4 Relation Patterns

If both arguments from a promoted relation instance are found in a sentence, then the intervening sequence of words gets brought out. To extract instance relation patterns according to the instances, we do have CBLAlgorithm/pattern/pattern\_extract\_relation.py file.

```
8 # Extract the relation markings according to instances.
9 def extract_relation_patterns(text_corpus, collection_1, collection_2):
10     supplement_findings = []
11
12     # Find regression findings for l_regex_findings
13     l_regex_findings = utils_init_reg_expression.init_reg_expression(collection_1)
14     # Find regression findings for r_regex_findings
15     r_regex_findings = utils_init_reg_expression.init_reg_expression(collection_2)
16
17     init = '('
18     mid = '(.+?)'
19     end = ')'
20
21     regex_pattern = init + l_regex_findings + end
22     regex_pattern = regex_pattern + mid
23     regex_pattern = regex_pattern + init + r_regex_findings + end
24
25     for idx in text_corpus:
26         # Function to check whether corpus_text is candidate obj or not.
27         if utils_check_candidate_text_corpus.check_candidate_text_corpus(idx, regex_pattern):
28             # Return a list of all non-overlapping matches in the string.
29             POS_findings = re.findall(regex_pattern, idx)
30
31             for findings in POS_findings:
32                 new_findings = ''.join([str2tuple(tag)[0] for tag in findings[1].split()])
33                 obj1 = new_findings.strip()
34                 if obj1 is not None:
35                     # If any coupling relation
36                     if pattern_check_coupling_relation.check_coupling_relation(obj1):
37                         # Getting supplement of the object.
38                         findings = utils_filling_supplement.filling_supplement(obj1, 'BOTH')
39                         supplement_findings.append(findings)
40     return supplement_findings
```

### 4.3.5 Candidate Filtering

Filtering to maintain high precision (consistency measure). Mainly to avoid extremely specific patterns. An instance can be only considered for assessment if it co-occurs with at least two promoted patterns in the text corpus, and if the co-occurrence count with all promoted patterns is at least three times its co-occurrence counts with negative patterns.

### 4.3.6 Candidate Assessment

For each predicate, CBL algorithm trains a discretized Native Bayes classifier to classify candidates' instances. The current sets of promoted and negative instances are used as training example for the classifier. Patterns are mapped using estimate of precision of each pattern  $p$ :

$$Precision(p) = \frac{\sum_{i \in I} count(i, p)}{count(p)}$$

- 1) I: set of promoted instances for predicate currently being considered.
- 2) Count (i, p): co-occurrence count of instance i with pattern p.
- 3) Count(p): Hit count of pattern p. Because p believes that the rest of the occurrences or pattern p are not positive examples of the predicate, this is a negative estimate.
- 4) All the co-occurrence counts needed for the assessment step are collected in the same MapReduce pass as those required for filtering candidates.

```
1 # Count it occurs from mutually exclusive p_list.
2 def count_get_N(inst, error_list, premises_list):
3
4     # Initiate count as 0.
5     N = 0
6
7     # Traverse premises_list
8     for idx in premises_list:
9         if idx.name in error_list:
10             continue
11
12         # Get the candidate_instances
13         inst_candidate = inst.candidate_instances
14
15         # if markings is in inst_candidate
16         if inst in inst_candidate:
17             # marking_list = inst_candidate[markings]
18             marking_list = inst_candidate[inst]
19             # count N
20             N = N + sum(marking_list.values())
21     return N
```

Filter the candidate instances. Input categories or relations will have promoted instances and their co-occurrence count with patterns.

```

4 # Input: predicate list
5 def filter_candidate_Mc(list):
6     advanced_F = list()
7     # For each predicate in the list
8     for p in p_list:
9         # Get the all the filtered_candidate_instances as a dictionary.
10        p.filtered_candidate_instances = dict()
11        # Get the all the last_promoted_instances as a list.
12        p.last_promoted_instances = list()
13        # Get the word_exceptions and record predicate's name into it.
14        exclusion_list = p.word_exceptions
15        # Append p.name into exclusion_list
16        exclusion_list.append(p.name)
17
18        # Traverse the p.candidate_instances.items()
19        for inst, markings in p.candidate_instances.items():
20
21            # Count it occurs from actually excludes p_list.
22            count_markings = filter_count_get_M_count_get_M(inst, exclusion_list, p_list)
23
24            if count_markings == 0:
25                # Decrease the markings by 1
26                count_markings = 1
27            else:
28                # Increase the markings by 1
29                count_markings = count_markings + 1
30
31            # Traverse markings.values()
32            for n in markings.values():
33                if n == CBL_FILTER_CONSTANT * count_markings:
34                    # Add inst into filtered_candidate_patterns
35                    p.filtered_candidate_instances[inst] = sum(markings.values())
36
37                    if type(inst) == tuple:
38                        inst = list(inst)
39
40                    if inst not in p.last_promoted_instances:
41                        if inst not in p.instances:
42                            # Append inst into each predicate's last_promoted_patterns
43                            p.last_promoted_instances.append(inst)
44                            # Append inst into each predicate's patterns
45                            p.instances.append(inst)
46                            if p not in advanced_F:

```

### 4.3.7 Candidate Promotion

Ranks the candidate according after each iteration and shared among predicates. This is output of CBL algorithm.

```

def cblAlgorithmPromotedPatterns(text_collections, cblPattern):
    new_promoted_categories = dict()
    new_promoted_relations = dict()

    # preProcess text_collection.
    texts = utils_get_initialisation_text_corpus.get_initialisation_text_corpus(text_collections)

    # Learn category patterns and instances and added into the dictionary.
    learner_learn_category_instances.learn_category_instances(texts, cblPattern)
    learner_learn_category_patterns.learn_category_patterns(texts, cblPattern)
    new_promoted_categories.update(filterer_filter_candidate_M.filter_candidate_M(cblPattern.last_promoted_categories))
    new_promoted_categories.update(filterer_filter_candidate_marking.filter_candidate_marking(cblPattern.last_promoted_categories))

    # Learn relation patterns and instances and added them into dictionary.
    learner_learn_relation_instances.learn_relation_instances(texts, cblPattern)
    learner_learn_relation_patterns.learn_relation_patterns(texts, cblPattern)
    new_promoted_relations.update(filterer_filter_candidate_M.filter_candidate_M(cblPattern.last_promoted_relations))
    new_promoted_relations.update(filterer_filter_candidate_marking.filter_candidate_marking(cblPattern.last_promoted_relations))

    # update the pattern's last_promoted_categories and last_promoted_relations
    cblPattern.last_promoted_categories = new_promoted_categories
    cblPattern.last_promoted_relations = new_promoted_relations

    return new_promoted_categories, new_promoted_relations

```

## 5 Experimental Evaluation

We planned our exploratory assessment to attempt to respond to the accompanying inquiries: Might CBL at any point repeat ordinarily regardless accomplish high accuracy? How supportive are the kinds of coupling that we utilize? Could we at any point broaden existing semantic assets?

### 5.1 Configurations of the Algorithm

We ran our algorithm in three configurations:

- 1) Full: The calculation as depicted in Segment 4.2.

## 5.2 Experimental Procedure - We tried to implement our own, by reading online resources

- 1) We ran every design for 3 cycles. To evaluate the accuracy of advanced occasions, we sampled approximately 7297 examples from the advanced set for each predicate in every design after 1, 2, and 3 iterations, pooled together the examples for each predicate, and afterward passed judgment on their rightness.
- 2) The appointed authority didn't realize which run a case was tested from. We assessed the accuracy of the advanced examples from each pursue 1, 2, and 3 iterations as the quantity of right advanced occasions isolated by the number inspected.
- 3) While tests of 7297 cases don't deliver tight certainty intervals around individual evaluations, they are adequate for testing for the impacts in which we are intrigued.
- 4) It's a basically, extract the word from sentence, map POS tagged words and then, check its category and relation from pattern.
- 5) We try to check result match percentage with each iteration. Finally, we get able to map relation between categories and relation.
- 6) To run app, python app.py.

## 5.3 Results

**Data:** To download stop-words:

```
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords
stop_words = set(stopwords.words('english'))
```

- 1) For POS (Part-of-Speech tagger) downloaded tagger:  
english-left3words-distsim.tagger
- 2) To tokenize the input corpus (we used):  
stanford-postagger.jar
- 3) United Nations General Debate Corpus Data from Harvard Data verse: (Test-Data)  
<https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/OTJX8Y%27>
- 4) We're all aware that paragraphs are made up of words from distinct parts of speech (POS). In the English language, there are eight different POS: noun, pronoun, verb, adjective, adverb, preposition, conjunction, and intersection. The POS regulates how a word in a phrase function in terms of meaning. This demonstrates that the POS tagging of a term has a significant impact on comprehending the content of a phrase. Without a doubt, we can use it to extract useful data from our data.
- 5) Split the word and check with POS in the English language: noun, pronoun, verb, adjective, adverb, preposition, conjunction, and intersection.



```

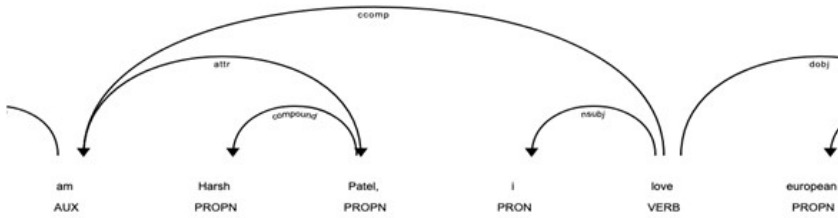
This -> PRON
is -> AUX
a -> DET
sample -> NOUN
sentence -> NOUN
... -> PUNCT

Word in the 'text', which is noun: sample
Word in the 'text', which is noun: sentence

The -> DET
children -> NOUN
love -> VERB
green -> NOUN
classmate -> NOUN
Dependency label: 'nsubj' -> children
Dependency label: 'advmod' -> classmate

TXT/Session 73 - 2018
TXT/Session 29 - 1974
TXT/Session 50 - 1995
TXT/Session 64 - 2009
TXT/Session 51 - 1996
TXT/Session 27 - 1972
TXT/Session 70 - 2015
TXT/Session 34 - 1979
TXT/Session 53 - 1998
TXT/Session 38 - 1983
TXT/Session 35 - 1980
TXT/Session 44 - 1989
TXT/Session 26 - 1971
TXT/Session 71 - 2016
TXT/Session 48 - 1993
TXT/Session 62 - 2007
TXT/Session 75 - 2020
TXT/Session 45 - 1990
TXT/Session 59 - 2004
TXT/Session 48 - 2013
TXT/Session 42 - 1987
TXT/Session 56 - 2001
TXT/Session 65 - 2018
TXT/Session 32 - 1977
TXT/Session 67 - 2002
TXT/Session 72 - 2017
TXT/Session 28 - 1970

```



Wordlist:= I am Harsh Patel, i love european food

6) We need to import data from TXT/Session\* /IND\*.txt.

```

TXT/Session 73 - 2018
TXT/Session 29 - 1974
TXT/Session 50 - 1995
TXT/Session 64 - 2009
TXT/Session 51 - 1996
TXT/Session 27 - 1972
TXT/Session 70 - 2015
TXT/Session 34 - 1979
TXT/Session 53 - 1998
TXT/Session 38 - 1983
TXT/Session 35 - 1980
TXT/Session 44 - 1989
TXT/Session 26 - 1971
TXT/Session 71 - 2016
TXT/Session 48 - 1993
TXT/Session 62 - 2007
TXT/Session 75 - 2020
TXT/Session 45 - 1990
TXT/Session 59 - 2004
TXT/Session 48 - 2013
TXT/Session 42 - 1987
TXT/Session 56 - 2001
TXT/Session 65 - 2018
TXT/Session 32 - 1977
TXT/Session 67 - 2002
TXT/Session 72 - 2017
TXT/Session 28 - 1970

```

7) We need to gather all the data and clean the data first.

Year	Speech	Country	Session	Speech_clean
0	2018 On my own behalf and on behalf of my country, ...	TXT/Session 73	2018/IND	73 On my own behalf and on behalf of my country, ...
1	1974 Mr. President, I have already had occasion to ...	TXT/Session 29	1974/IND	29 Mr President, I have already had occasion to c...
2	1995 It gives me great pleasure to\congratulate Mr...	TXT/Session 50	1995/IND	50 It gives me great pleasure to congratulate Mr ...
3	2009 I offer my congratulations \nto Mr. Treki on h...	TXT/Session 64	2009/IND	64 I offer my congratulations to Mr Treki on his...
4	1996 It gives me great pleasure to\congratulate A...	TXT/Session 51	1996/IND	51 It gives me great pleasure to congratulate Am...

8) We need to extract information, we need to better grasp the structure of a text, I'll display the dependencies of an example instance in a tree format, which provides a better knowledge of the structure.



```

df_rule1:=
      Year      Sent \
0      2018 On my own behalf and on behalf of my country, ...
1      2018 As a woman, I feel doubly proud that this hon...
2      2018 I also recall, with equal pride, that the fir...
3      2018 I also thank former President Miroslav Lajčák...
4      2018 We received the tragic news this mor...
...
7292 2014 In preparing and implementing the post 2015 ...
7293 2014 For 2015, let us come together to give a new ...
7294 2014 The year 2015 should be a banner year in his...
7295 2014 I hope that we will all live up to its promise
7296 2014

      Output
0      []
1      []
2      []
3      [I thank Lajčák]
4      [We receive news]
...
7292  []
7293  []
7294  []
7295  []
7296  []

[7297 rows x 3 columns]
37.11114156502672 % pattern match
37.11114156502672

```

We then get counting the number of sentences containing the verb and nouns.

```

      Sent      Year      Noun1 \
0      I also thank former President Miroslav Lajčák... 2018      [I]
1      We received the tragic news this mor... 2018      [We]
2      The United Nations is the world's principal m... 2018 [nations]
3      The United Nations is the world's principal m... 2018 [that]
4      A common refrain since 2015 has been that we ... 2018      [we]
...
3070 If we can keep all of that in mind, through ... 2014 [who]
3071 If we can keep all of that in mind, through ... 2014 [we]
3072 By taking full advantage of this moment, we ... 2014 [we]
3073 We could give the journey of the United Nati... 2014 [We]
3074 We could give the journey of the United Nati... 2014 [We]

      Verb      Noun2
0      thank      [Lajčák]
1      receive     [news]
2      seek        [balm]
3      correct     [imbalances]
4      reach       [horizon]
...
3070 contribute   [ideas]
3071 find          [ways]
3072 achieve       [consciousness]
3073 give          [journey]
3074 give          [journey, form]

[3075 rows x 5 columns]

```

Let's look at the most often used verbs in the sentences.

```

[('have', 269), ('take', 122), ('give', 73), ('provide', 64), ('welcome', 63), ('support', 57), ('make', 54), ('see', 48), ('need', 43), ('face', 40), ('bring', 39), ('require', 37), ('reflect', 34), ('urge', 30), ('represent', 29), ('show', 28), ('express', 26), ('receive', 25), ('offer', 23), ('constitute', 23)]

```

Get the list of the verb, for 'have' as a verb.

	Sent	Year	Noun1	Verb	Noun2
10	Through the Pradhan Mantri Jan Dhan Yojana, t...	2018	[Yojana]	have	[accounts]
11	Through the Pradhan Mantri Jan Dhan Yojana, t...	2018	[Yojana]	have	[accountsinside, have, accounts]
14	We have a prayer in India Sarve Santu Niramay...	2018	[We]	have	[prayer]
16	Similarly, we have launched the largest scale...	2018	[everyone]	have	[roof]
17	Similarly, we have launched the largest scale...	2018	[everyone]	have	[roof, head]
...	...	...	...	...	...
3011	India's ancient wisdom sees the world as one ...	2014	[country]	have	[philosophy]
3012	India is a country where, beyond nature, we ...	2014	[we]	have	[communication]
3014	Owing to our ideology, we have a firm belief...	2014	[we]	have	[belief]
3022	I have the same policy towards Pakistan	2014	[I]	have	[policy]
3034	Why, when we have a good forum like the Unit...	2014	[we]	have	[forum]

269 rows x 5 columns

Get the list of the verb, for ‘constitute’ as a verb.

	Sent	Year	Noun1	Verb	Noun2
362	They constitute a blueprint that is more comp...	2015	[They]	constitute	[blueprint]
436	This constitutes a guarantee against the misu...	1979	[This]	constitute	[guarantee]
457	If I have spoken at length on nuclear disarm...	1979	[weapons]	constitute	[danger]
796	Collective self reliance through south south ...	1989	[reliance]	constitute	[plank]
1277	Even though these forces constitute only a ve...	1987	[forces]	constitute	[fraction]
1447	Fourthly, despite their division into nation ...	1977	[people]	constitute	[family]
1465	We are told that nuclear weapons are necessar...	1977	[that]	constitute	[core]
1704	A successful launching of the global negotiat...	1981	[launching]	constitute	[success]
1953	Reports indicate that the armed forces in tha...	1984	[Tamils]	constitute	[majority]
2015	Youth, which constitutes a crucial segment of...	1984	[which]	constitute	[segment]
2223	It constitutes a unique international forum w...	1985	[it]	constitute	[forum]
2250	The policies of apartheid of the racist regim...	1985	[policies]	constitute	[source]
2260	The United Nations Security Council, convened...	1985	[which]	constitute	[basis]
2261	The pursuit of apartheid, the occupation of N...	1985	[occupation]	constitute	[threats]
2323	It constitutes a crime against humanity	2000	[it]	constitute	[crime]
2369	The very fact that these global problems are ...	1975	[Nations]	constitute	[forum]
2532	Israel arrogant defiance of the will of the I...	1986	[all]	constitute	[chapters]
2534	We would like to underscore once again the im...	1986	[that]	constitute	[contribution]
2586	When a developing country is able to successf...	2019	[achievements]	constitute	[message]
2675	At the same time, we believe that any outside...	1991	[intervention]	constitute	[abridgement]
2717	Sponsorship of terrorism in another country c...	1991	[Sponsorship]	constitute	[violation]
2770	The policies of the major developed countries...	1988	[policies]	constitute	[determinants]
2858	Suffice it to say that Security Council resol...	1978	[resolutions]	constitute	[basis]

11) **Information Extraction of Adjective Noun Structure:** I extracted the Adjective Noun Structure in the preceding rule, but the data did not feel full. This is since many nouns have had an adjective or a compound dependent word that adds to the meaning of the noun. We can learn more about subject and object if we extract these including the noun.

We have a pattern match of more than 51% for Adjective Noun Structure phrase, and we can test it for all the words in the corpus.

```
52.15827338129496 % pattern match
```

```
52.15827338129496
```

For We then get counting the number of sentences containing the Adjective Noun Structure.

```
76.66164177059065 % pattern match
```

```
76.66164177059065
```

Out of 7297, 1668 sentences matched our pattern rule.

```
df_rule:=
  Year      Sent \
0    2018  On my own behalf and on behalf of my country, ...
1    2018  As a woman, I feel doubly proud that this hon...
2    2018  I also recall, with equal pride, that the fir...
3    2018  I also thank former President Miroslav Lajčák...
4    2018  We received the tragic news this mor...
...    ...
7292 2014  In preparing and implementing the post 2015 ...
7293 2014  For 2015, let us come together to give a new ...
7294 2014  The year 2015 should be a banner year in his...
7295 2014  I hope that we will all live up to its promise
7296 2014

                                Output
0                                [own behalf, third session]
1                                []
2  [equal pride, first woman, eminent position]
3                                [second session]
4                                []
...
7292                                [development agenda]
7293  [new direction, new lease]
7294                                [turning point]
7295                                []
7296                                []

[7297 rows x 3 columns]
```

Now, if we combine two rules mentioned ‘Noun-Verb-Noun Phrases’ and ‘Noun-Verb-Noun Phrases’, we match.

	Year	Sent	Output
0	2018	On my own behalf and on behalf of my country, ...	[]
1	2018	As a woman, I feel doubly proud that this hon...	[]
2	2018	I also recall, with equal pride, that the fir...	[]
3	2018	I also thank former President Miroslav Lajčák...	[ I thank Lajčák]
4	2018	We received the tragic news this mor...	[ We receive news]
5	2018	From this rostrum, on behalf of India, I wish...	[]
6	2018	I would also like to assure them that India w...	[]
7	2018	The United Nations is the world's principal m...	[ nations seek balm, that correct skewed econ...
8	2018	In 2015, we established 2030 as a critically ...	[]
9	2018	A common refrain since 2015 has been that we ...	[ we reach horizon, India find way]
10	2018	I assure the General Assembly through you, Ma...	[ I assure Assembly]
11	2018	We are totally committed to achieving those o...	[]
12	2018	Under the leadership of Prime Minister Narend...	[ India initiate unprecedented economic transf...
13	2018	I will provide an overview to illustrate the ...	[ I provide overview]
14	2018	Through the Pradhan Mantri Jan Dhan Yojana, t...	[ Yojana have accounts, Yojana have accounts ...
15	2018	The programme has enabled the poor to receiv...	[]
16	2018	Similarly, Ayushman Bharat Yojana, the world'...	[]
17	2018	That revolutionary scheme will benefit 500 mi...	[ scheme benefit Indians, who receive coverage]
18	2018	We have a prayer in India Sarve Santu Niramay...	[ We have prayer]
19	2018	The Ayushman Bharat Yojana, or National Healt...	[]
20	2018	Similarly, we have launched the largest scale...	[ we launch largest programme, everyone have ...
21	2018	Under the programme, we have set for oursel...	[ we set target]
22	2018	So far, more than 5 million homes for the p...	[]
23	2018	Two extremely effective programmes have also ...	[]
24	2018	I stress that more than 140 million Indians h...	[ real Indians take loans]
25	2018	The most significant aspect of the MUDRA sche...	[]
26	2018	At the heart of Prime Minister Modi's transfo...	[]
27	2018	All the programmes that I have just mentioned...	[ programmes have welfare]

12) **Information Extraction of Prepositions Phrases:** When we come across a preposition, we look to see if it has a noun as a head word. We search through all the tokens for prepositions. Proposition is it indicates where or when something is connected to something else.

For Short sentence, we get 48.50119904076738 % pattern match.

48.50119904076738 % pattern match  
48.50119904076738

For large corpus, we get 48.50119904076738 % pattern match.

48.50119904076738 % pattern match  
48.50119904076738

The data frame shows the result of the frame of entire corpus.

	Sent	Year	Noun1	Preposition	Noun2
0	On my own behalf and on behalf of my country, ...	2018	behalf	of	[country]
1	On my own behalf and on behalf of my country, ...	2018	election	as	[President]
2	On my own behalf and on behalf of my country, ...	2018	election	at	[session]
3	I also thank former President Miroslav Lajčák...	2018	session	of	[Assembly]
4	We received the tragic news this mor...	2018	news	of	[tsunami]
5	From this rostrum, on behalf of India, I wish...	2018	behalf	of	[India]
6	From this rostrum, on behalf of India, I wish...	2018	people	of	[Indonesia]
7	The United Nations is the world's principal m...	2018	wounds	of	[history]
8	The United Nations is the world's principal m...	2018	platform	for	[solutions]
9	In 2015, we established 2030 as a critically ...	2018	horizon	for	[Goals]

The top-promoted category

[('of', 7167), ('in', 1393), ('for', 982), ('to', 612), ('on', 466), ('with', 282), ('as', 173), ('between', 162), ('by', 142), ('from', 127), ('against', 99), ('at', 83), ('towards', 81), ('among', 73), ('over', 43), ('within', 42), ('under', 41), ('into', 37), ('about', 24), ('than', 23)]

13) The next iteration, decides, the better results.

	Sent	Year	Noun1	Preposition	Noun2
0	On my own behalf and on behalf of my country, ...	2018	behalf	of	[country]
1	On my own behalf and on behalf of my country, ...	2018	election	as	[President]
2	On my own behalf and on behalf of my country, ...	2018	election	at	[session]
3	I also thank former President Miroslav Lajčák...	2018	session	of	[Assembly]
4	We received the tragic news this mor...	2018	news	of	[tsunami]
...	...	...	...	...	...
12170	We could give the journey of the United Nati...	2014	journey	of	[Nations]
12171	I therefore feel that 70 years is a great op...	2014	opportunity	for	[]
12172	Let us come together and fulfill our promise t...	2014	improvements	to	[Council]
12173	For 2015, let us come together to give a new ...	2014	lease	on	[development]
12174	The year 2015 should be a banner year in his...	2014	year	in	[history]

12175 rows x 5 columns

With experiment, up to 5 iterations, we are getting better results. We present all the scenarios. It's somewhat we are expected, but with more POS, we are getting accurate results.

## 6 Conclusion

We have introduced a technique for coupling the semi-supervised learning of classes and relations and showed exactly that the coupling hinders the issue of semantic float related to bootstrap learning strategies. We suspect that learning extra predicates at the same time will yield considerably more exact learning. A surmised comparison with a current archive of semantic knowledge.

## 7 Project Reference URL

We read online articles from medium.com, open-source libraries, and some YouTube implementation videos (All I've mentioned below)

1. <https://www.analyticsvidhya.com/blog/2020/06/nlp-project-information-extraction/>
2. <https://towardsdatascience.com/from-text-to-knowledge-the-information-extraction-pipeline>
3. <https://medium.com/swlh/python-nlp-tutorial-information-extraction-and-knowledge-graphs>
4. <https://github.com/midhun-pk/cpl>
5. <https://www.youtube.com/watch?v=Tj3Dkiw-iZg>
6. <https://nanonets.com/blog/information-extraction/>
7. <https://www.findaphd.com/phds/project/deep-learning-for-semantic-based-information-extraction>

## 8 Team Member's Contribution

- 1) **Harsh Patel** : Worked on CBLAlgorithm - extraction, Filter, Promoted patterns, Dictionary part in category and relation, worked on architectural design part, understanding some basic examples on web and go through articles, working on final iteration part and report part.
- 2) **Uday Ramesh** : Worked on CBLAlgorithm - extraction, Filter, Promoted patterns, Dictionary part in category and relation, go through basic examples, articles related to it and report.
- 3) **Jashandeep Singh Sran**: Worked on MongoDB part, read examples through net,report, Asses part in CBL algorithm.

## References

- [1] Eugene Agichtein and Luis Gravano. 2000. Snowball: Extracting relations from large plain-text collections. In JCDL.
- [2] Michele Banko, Michael J. Cafarella, Stephen Soderland, Matt Broadhead, and Oren Etzioni. 2007. Open information extraction from the web. In /JCA/.

- [3] Avrim Blum and Tom Mitchell. 1998. Combining labeled and unlabeled data with co-training. In COLT.
- [4] Sergey Brin. 1998. Extracting patterns and relations from the world wide web. In WebDB Workshop at 6th International Conference on Extending Database Technology.
- [5] Rich Caruana. 1997. Multitask learning. *Machine Learning*, 28:41-75.
- [6] Ming-Wei Chang, Lev-Arie Ratinov, and Dan Roth. 2007. Guiding semi-supervision with constraint- driven learning. In ACL.
- [7] Michael Collins and Yoram Singer. 1999. Unsupervised models for named entity classification. In EMNLP.
- [8] James R. Curran, Tara Murphy, and Bernhard Scholz. 2007. Minimising semantic drift with mutual exclusion bootstrapping. In PACLING.
- [9] Jeffrey Dean and Sanjay Ghemawat. 2008. Mapreduce: simplified data processing on large clusters. *Commun. ACM*, 51(1):107-113.
- [10] Doug Downey, Matthew Broadhead, and Oren Etzioni. 2007. Locating complex named entities in web text. In JCAI/.
- [11] Oren Etzioni, Michael Cafarella, Doug Downey, Ana Maria Popescu, Tal Shaked, Stephen Soderland, Daniel S. Weld, and Alexander Yates. 2005. Unsupervised named-entity extraction from the web: an experimental study. *Artif. Intell.*, 165(1):91-134.
- [12] Usama M. Fayyad and Keki B. Irani. 1993. Multi interval discretization of continuous-valued attributes for classification learning. In UAI.
- [13] Marti A. Hearst. 1992. Automatic acquisition of hyponyms from large text corpora. In COLING.
- [14] Qiuhua Liu, Xuejun Liao, and Lawrence Carin. 2008. Semi-supervised multitask learning. In NIPS.
- [15] David McClosky, Eugene Charniak, and Mark Johnson. 2006. Effective self-training for parsing. In NAACL.
- [16] Luke K. McDowell and Michael Cafarella. 2006. Ontology-driven information extraction with ontosyphon. In ISWC.
- [17] Metaweb Technologies. 2009. Freebase data dumps. <http://download.freebase.com/datadumps/>.
- [18] Marius Paşca, Dekang Lin, Jeffrey Bigham, Andrei Lifchits, and Alpa Jain. 2006. Names and similarities on the web: fact extraction in the fast lane. In ACL.
- [19] Marius Pasca, Dekang Lin, Jeffrey Bigham, Andrei Lifchits, and Alpa Jain. 2006. Organizing and searching the world wide web of facts - step one: The onemillion fact extraction challenge. In AAAI/.
- [20] Deepak Ravichandran and Eduard Hovy. 2002. Learning surface text patterns for a question answering system. In ACL.
- [21] Ellen Riloff and Rosie Jones. 1999. Learning dictionaries for information extraction by multi-level bootstrapping. In AAAI/.
- [22] Benjamin Rosenfeld and Ronen Feldman. 2007. Using corpus statistics on entities to improve semi-supervised relation extraction from the web. In ACL.
- [23] Yusuke Shinyama and Satoshi Sekine. 2006. Preemptive information extraction using unrestricted relation discovery. In HLT-NAACL.
- [24] Mark D. Smucker, James Allan, and Ben Carterette. 2007. A comparison of statistical significance tests for information retrieval evaluation. In CIKM.
- [25] Sebastian Thrun. 1996. Is learning the n-th thing any easier than learning the first? In NIPS.
- [26] Peter D. Turney. 2001. Mining the web for synonyms: Pmi-ir versus Isat on toefl. In EMCL.
- [27] Nicola Ueffing. 2006. Self-training for machine translation. In NIPS workshop on Machine Learning for Multilingual Information Access.
- [28] Roman Yangarber. 2003. Counter-training in discovery of semantic patterns. In ACL.