

The Title

By

Alice Irankunda (alice.irankunda@aims.ac.rw)
Supervised by Dr.Yabebal FANTAYE

June 2017

*AN ESSAY PRESENTED TO AIMS RWANDA IN PARTIAL FULFILMENT OF THE REQUIREMENTS FOR THE AWARD OF
MASTER OF SCIENCE IN MATHEMATICAL SCIENCES*

⁷ **DECLARATION**

8 **ACKNOWLEDGEMENTS**

9 DEDICATION

¹⁰ **Abstract**

11 **Contents**

12 **Dedication** **iii**

13 **1 Introduction** **1**

14 **2 Literature review** **3**

15 2.1 Natural Language Processing (NLP) 3

16 2.2 Parse Tree 3

17 2.3 Named Entity Recognition and Classification NERC 4

18 **3 Third Chapter** **10**

19 **References** **11**

1. Introduction

1.1. Introduction

In every organization there is a way to communicate ,one of the most popular way to transmit the information is to produce a written report which explains how different activities of the organization are going. For the large organizations there a huge number of reports and it is so challenging to go through each and every report manually. This research has an aim of providing an easy way of visualizing and extracting the important information locked in reports from NGO and large organisations.

In 1919, The International Federation of Red Cross and Red Crescent societies (IFRC) has been founded, it has some millions of reports related to humanitarian support,How to know automatically the number of people who suffered from a disease? How to know the fraction of fund spent on shelter ?

In this research, there are some solutions to those questions by using combination of statistics and Natural Language Processing (NLP)techniques. Big data and Machine learning is for analysing the huge data by using statistical and computing algorithms. Document modelling by extracting entities is one of the way to deal with natural big data linguistic problems where entity is defined as a single unit of data like location , people, organization and so one. entities can be classified based on their relationship.

For example, consider MDRAF003 to be IFRC report "Afghanistan MDRAF003 26May2016.pdf", it is composed by 12 pages of texts where the entities will be extracted from.

what are the key procedures to be done in order to perform entities extraction ?

- The sentences which compose a report must be parsed.
- Entities also must be identified in the report and classified.
- Relationship between entities must be modelled.

In this research, there is a clear discussion about powerful techniques to answer the previous questions. Natural Language Processing techniques used to sentence level and content based analysis,Natural Language ToolKit (NLTK) for splitting the sentences into tokens and remove the common words and how to work with corpus.The used reports for the implementation of different language algorithms are from IFRC .

1.2. Motivation

Big data and Machine learning have recently become one of the major and strong solution finder to most difficult problems in heath, statistical prediction, company development and linguistics.

52 within huge reports,journals or articles ,this work will return significant classified entities which
53 will help the user to not struggle opening the report and get like amount spent in a given activity,
54 the sum of people who participated in an event etc.

2. Literature review

In today's life, many organizations are generating unstructured data while they are communicating, the entities to be extracted from those reports are plenty, In this research, all reports we considered are written in English.

2.1 Natural Language Processing (NLP)

To label the boundaries of sentences is one of the important prerequisite steps in NLP but the punctuation marks cause some ambiguity (Palmer and Hearst, 1994) for example it is challenging to differentiate the the point in abbreviations and a full stop. To handle this ambiguity some systems use the special purpose-regular expression grammar ,exception rule method etc.

David D.Palmer and Marti A.Hearst worked on the problem of punctuations (Palmer and Hearst, 1994), They developed an efficient system with high accuracy in automatic labelling the boundaries of the sentence by using the feed forwarding neural - networks where the input was the POS probabilities of all tokens which are surrounding the punctuation and output was found as the label to be assigned to the token.This work was able to correct up to 98.5% for punctuation of sentence- boundaries.A proposed new approach was how to represent the context of punctuation marks without ambiguities.

This research will also look at how neural networks can be used to label different tokens.

Capitalization can be used in different ways such as the beginning of the proper noun, the abbreviation, the post of high level profile people etc. Considering the English language text , if we are given a particular token it is not by chance to determine whether it is a name or not. some of the approaches to indicate a name are to use capitalization ,detection of sentence boundaries and dictionaries (Baluja et al., 2000).

2.2 Parse Tree

One of the sentences that compose our sample report says : " Assessment reports indicated 117 deaths, 544 people injured, 12,794 homes damaged and 7,384 houses destroyed", Suppose that this sentence is called "S"

There are two mains steps which can be performed to get the entities from this sentence :

- **Tokenizing:** This is a procedure of taking a sentence and extract the composing atomic linguistic elements means words,verbs,punctuations, adjectives etc . S has the following tokens: ['Assessment', 'reports', 'indicated', '117', 'deaths', ',', '544', 'people', 'injured', ',', '12,794', 'homes', 'damaged', 'and', '7,384', 'houses', 'destroyed']
- **POS:** part-of-speech is a process of attaching to every linguistic element of the sentence

a corresponding tagg based on grammar rules. The POS of S are: [('Assessment', 'JJ'), ('reports', 'NNS'), ('indicated', 'VBD'), ('117', 'CD'), ('deaths', 'NNS'), (',', ','), ('544', 'CD'), ('people', 'NNS'), ('injured', 'VBN'), (',', ','), ('12,794', 'CD'), ('homes', 'NNS'), ('damaged', 'VBN'), ('and', 'CC'), ('7,384', 'CD'), ('houses', 'NNS'), ('destroyed', 'VBD')]

The meanings of the used tags for S :

- JJ : **Adjective** : 'Assessment'
- NNS : **Noun,plural**: 'reports', 'deaths', 'people', 'houses'
- VBD : **Verbs, past tense**: 'indicated', 'injured', 'damaged', 'destroyed'
- CD : **Cardinal Number**: '117', '544', '12,794', '7,384',
- CC : **Coordinate Conjugation**: 'and'

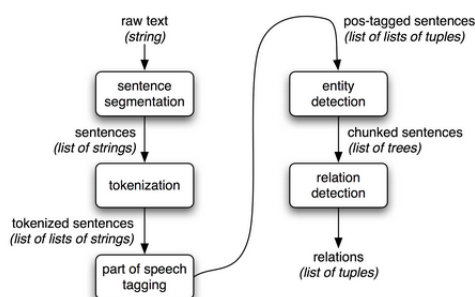
The parse tree is formed based on the POS, the classification of word and the way words are arranged in a sentence show a kind of relationship between words.

Figure 2.1: The parse of the above sentence



The process of classifying entities can be more explained in the following picture

Figure 2.2: information extraction process



2.3 Named Entity Recognition and Classification NERC

The term "Named entity" has been coined in 1996 in "sixth Message understanding Conference" (MUC-6 R. Grishman and Sundheim 1996). Entity can be referred as a task, the entity is "named" when it is restricted to one or many rigid designators (Sharnagat, 2014), example: persons, location, product are the named entities.

Based on the classification of Standard Generalizes Markup Language(SGML) a task can be divided into three subtasks:

• ENAMEX: location,product,country ,organization

• NUMEX : percentage,quantity

• TIMEX : time, date

The entities from different reports. For extracting entities in a report there are different models which can be used:

2.1.1. Hidden Markov Model

This model is based on Bayesian probability inference which has been initiated in 18th century. HMM is the earliest applied model for Natural Entities Recognition for English language. The way needed task to be performed is to find the most likely sequence of tagged names(TN) given a sequence of words(SW).

$$P(TN|SW) = \frac{P(SW|TN)P(TN)}{P(SW)} \quad (2.3.1)$$

The equation (2.3.1) is conditional probability, $P(TN|SW)$ can be called posterior and it is the probability of an event Sequence of word occurring given Tagged names has observed. $P(SW|TN)$ is also called likelihood means it is the probability of observing the sequence of words(SW) when the given hypothesis tagged name(TN) is true. on another hand $P(TN)$ doesn't depend on the evidences, $P(TN)$ is called prior means that it is true even if there is no given evidence at all(masters thesis). We can be ignored $P(SW)$ and the remaining objective is to maximise the probability of getting the sequence of tagged names when sequence of words is given.

$$Max [P(TN|SW)] \quad (2.3.2)$$

From the equation (2.3.2) of the maximization , the following estimation can be made

$$P(TN) \approx \prod_{i=1}^n P(TN_i|TN_{i-1}) \quad (2.3.3)$$

Where TN_i is a tag in the sequence of names (TN) , for the likelihood probability can be estimated as

$$P(SW|TN) \approx \prod_{i=1}^n P(SW_i|TN_i) \quad (2.3.4)$$

The above estimations was for a small sequence where TN_i is a tag in the sequence of names (TN) and SW_i is a tag at index i in a sequence words (SW). For the large training corpus , the needed step is estimate based on the number of times the tag occurs and the position of the tag in a given corpus.

$$P(T_i|T_{i-1}) = \frac{K(T_{i-1}, T_i)}{K(T_{i-1})} \quad (2.3.5)$$

Based on the training corpus, $K(T_{i-1}, T_i)$ is referred as a how many times the tag T_i occurs after the tag T_{i-1} . in the corpus, $K(T_{i-1})$ is considered as the number of occurrences for the tag T_{i-1} .

Therefore the estimation can be performed as follow:

$$P(C_i|T_i) = \frac{K(T_i, C_i)}{K(T_i)} \quad (2.3.6)$$

From the equation (2.3.6), the term $K(T_i, C_i)$ is referred as the sum of the times that a word " C_i " has a tag T_i in the training corpus. The process of computing the posterior using the above steps is called Markov model.

2.1.1.1. Advantages of Hidden Markov Model

It is one of the most powerful statistical and machine learning (ML) techniques in modelling and high qualified in entities extraction. When the researcher is willing to train new data, HMM is very robust and efficient in computations.

2.1.1.2. Disadvantages of Hidden Markov Model

One of the limitations of HMM is that the researcher must have the notion of model topology and statistical techniques on how to deal with large amount of training data.

2.1.2. Supporting Vector Machine based model (To be edited for non linear data)

This model has an aims of classifying the named entities by using the linear support vector machine which separate input train documents into two categories, a document must be categorized as either positive or negative and be represented in two dimensional graph. Hyperplane is for separating train documents based on their categories and " w " is a weight vector which is perpendicular to hyperplane is represented by the following equation:

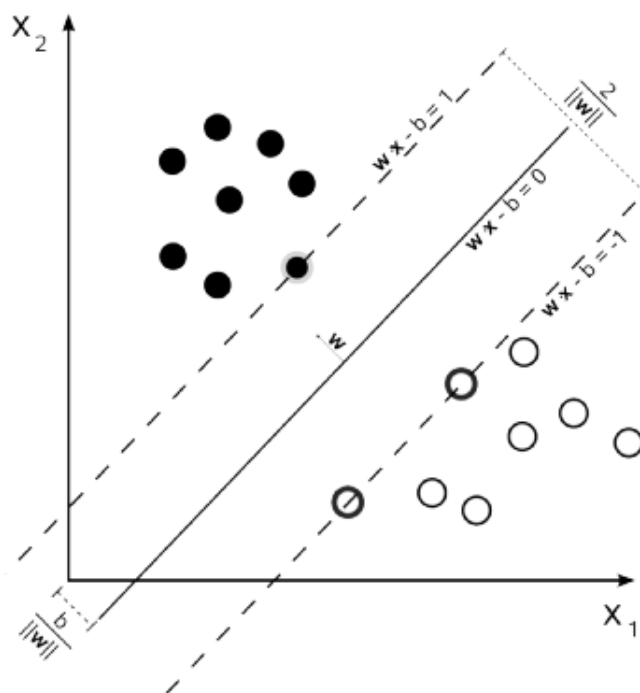
$$w \cdot x - b = 0 \quad (2.3.7)$$

128

From the (2.3.7), the offset of the hyperplane is $\frac{b}{\|w\|}$

The target is to maximize the the margin between the the points which represent two categories. remember that the vectors which pass through each of the point representative is perpendicular to the w , suppose that there will be an imaginary line which join two borders points h_-

Figure 2.3: SVM in hyperplane representation



and h_+ . Supporting vectors which are demonstrated by the dashed lines on the figure above are formed by :

$$w \cdot x - b = 1 \quad \text{and also} \quad (2.3.8)$$

$$w \cdot x - b = -1 \quad (2.3.9)$$

There are many algorithms with different approaches to optimization problems but all tends to the same solution says that minimize nwn automatically maximize the margin between h_- and h_+ where the boundary is a half way. Now, add another constraint for each document category from the equations (2.3.10) and (2.3.11), in order to hit the target

$$w \cdot x - b \geq 1 \quad \text{and also} \quad (2.3.10)$$

$$w \cdot x - b \leq -1 \quad (2.3.11)$$

2.1.2.1 Disadvantages of SVM

The classification of particular documents is not easy to be performed by SVM without destroying the constructed weights but with hand-written rule model. the machine learning prefers to use the decision tree procedure than SVM. in addition the decision tree has a detailed boolean-like model which is more popular to user.

Overview of rule/pattern based systems hand-written rule, decision tree, bootstrapping and

Hand-written rule

It is one of the standard approaches of NER and IE, it has been used for extracting the patterns from automated pages such as Amazon, NLP is so useful for unstructured human-written text by delivering part-of-speech (POS), syntactic parsing and categories of semantic words.

Rule /pattern based extraction

Many IE systems use rule/pattern to extract words and also phrases by looking to the context of those words or based on their surroundings. (Califf and Mooney, 2003). Some systems decided if the procedure of extracting the words should rely on the meaning of each word independently or on the context of their surroundings in a phrase. The limitation of this method is that some words do not have a closer mining to their surroundings that is why Patwardhan Siddharth with help of Ellen Riloff in workshop called "ACL 2006" presented another approach which was generating an automated IE system to learn patterns from a large fixed data set within a specific domain (Patwardhan and Riloff, 2007).

Our research deals with reports generated through a template, compared to the work of (Patwardhan and Riloff, 2007) templates usage is a limitation.

2.1.3. Text classification and Naive Bayes

It is one of the most important algorithms in text classification by using base rule and bag of words to classify the entities (Manning, 2012). The user instead of going through the report and start posing many queries, text classification algorithm transient the need information. Its aim is to build a function θ which takes the bag of words and returns the class of sentiment C either positive or negative.

$$\theta$$

$$\Updownarrow$$

ARCS initiated its response immediately after the earthquake struck to address the immediate needs. The National Society (NS) regional branches were at the forefront of the response and worked with Disaster Response Units (DRU). ARCS staff and volunteers were deployed promptly to support rescue efforts, provide first aid to the injured and distribute immediate relief supplies to affected people alongside undertaking initial assessments. A total of 900 volunteers were mobilised to support this response operation. ARCS also supported to transport critically injured people to hospital and mobilized community members for voluntary non-remunerated blood donations.

$$\Updownarrow$$

$$C$$

The procedure is to look for all words and retrieve those which form the subsets. Bag of words are formed after throwing away all words except the subsets. The use of the function θ is for

¹⁶⁴ attributing to each item of the bag of words a sentiment.

¹⁶⁵ Information extraction is a combination of segmentation, classification and clustering

3. Third Chapter

References

- Shumeet Baluja, Vibhu O Mittal, and Rahul Sukthankar. Applying machine learning for high-performance named-entity extraction. *Computational Intelligence*, 16(4):586–595, 2000.
- Mary Elaine Califf and Raymond J Mooney. Bottom-up relational learning of pattern matching rules for information extraction. *Journal of Machine Learning Research*, 4(Jun):177–210, 2003.
- Christopher Manning. Information extraction and named entity recognition, 2012.
- Christopher D Manning, Prabhakar Raghavan, and Hinrich Schütze. Text classification and naive bayes. *Introduction to information retrieval*, 1:6, 2008.
- David D Palmer and Marti A Hearst. Adaptive sentence boundary disambiguation. In *Proceedings of the fourth conference on Applied natural language processing*, pages 78–83. Association for Computational Linguistics, 1994.
- Siddharth Patwardhan and Ellen Riloff. Effective information extraction with semantic affinity patterns and relevant regions. In *EMNLP-CoNLL*, volume 7, pages 717–727, 2007.
- Rahul Sharnagat. Named entity recognition: A literature survey. *Center For Indian Language Technology*, 2014.