

A HYBRID APPROACH WITH FUZZY CLUSTERING FOR NER

Sobhana N V*

ABSTRACT

Named entity recognition is an extremely important and fundamental task of text mining. Some of the applications of text mining are Search or Information Access, Social media monitoring, E-Discovery, Records Management, National Security or Intelligence, Enterprise Business Intelligence/Data Mining, Competitive Intelligence etc. Machine learning methods like CRF, MEMM and SVM have been widely used for learning to recognize such entities from an annotated corpus. In this paper, we propose a novel kernel function for SVM and a combined approach which includes SVM and Conditional Random Fields (CRF) for named entity recognition (NER). The proposed kernel is based on calculating a novel distance function between the string based features and different contextual information of the words along with the variety of features that are helpful in predicting the various named entity (NE) classes. The proposed distance function makes use of certain statistics primarily derived from the training data and fuzzy clustering information. The training set consists of more than 2 lakh words and has been manually annotated with a NE tag set of seventeen tags. The system is able to recognize 17 classes of NEs with 81.99% Precision and 78.36 Recall.

Keywords: Named Entity Recognition, Precision, Recall

* Associate Professor, Rajiv Gandhi Institute of Technology, Kottayam

1. Introduction

Named Entity Recognition (NER) is a key part of information extraction system. NER involves identification of proper names in texts and their classification into a set of predefined categories of interest. Different categories are usually person names, location names, organization names, date & time expressions etc. A variety of techniques has been used for NER. The different approaches to NER include a. linguistic approaches b. machine Learning (ML) based approaches c. hybrid systems. The linguistic methods usually use rules manually written by linguists. There are several rule based NER systems, containing mostly lexicalized grammar, gazetteer lists, and list of trigger words, which are capable of providing up to 92% F-measure accuracy for English [1]. Linguistic approach uses hand crafted rules which require skilled linguistics. The main disadvantage of these rule based method is that they need vast experience and grammatical knowledge of the particular language or domain and these systems are not easily adaptable to other domains or languages [2]. Machine learning approaches are trainable and are thus much cheaper than that of rule-based ones. Some of the machine learning techniques used for the NER tasks are hidden markov model [3], Maximum Entropy Markov Model (MEMM) [4], Conditional Random Fields [5],[6]. Hybrid systems have been generally more effective for NER. The combination of MaxEnt, hidden markov model (HMM) and handcrafted rules for making NER is explored in [7].

Section 2 gives Characteristics of geological text. Section 3 discusses features used for Geological NER. Section 4 gives brief introduction to Conditional Random Fields, a machine learning approach to sequence labelling task. Section 5 describes the details of Geological Corpus. Section 6 explain the experiments and Results. The paper is concluded in section 7.

2. Geological named entity recognition

Geology is the study of origin, history and structure of the earth. Text mining on geological documents is an important area in scientific data mining. These documents contain spatial references and geo references in the form of spatial coordinates stored in database. They contain geospatial and temporal information. This spatial and temporal information is very important but normal text mining algorithms will fail to extract such information.

Named Entity Recognition (NER) is an important tool in almost all Natural Language

Processing (NLP) application areas. Proper identification and classification of named entities (NEs) are very big challenge to the NLP researchers. Geological NER has applications in several domains including information extraction, information retrieval, question answering [8], automatic summarization, machine translation [9] etc from Geological text.

Named entity recognition (NER) (also known as entity identification and entity extraction) is a subtask of information extraction that seeks to locate and classify atomic elements in text into predefined categories such as the names of persons, organizations, locations, expressions of times, quantities, monetary values, percentages, etc. The different categories of Geological Named Entities considered here are Country, State, City, Region, Mountain, Island, Water bodies, River, Village, Mineral, Year, Organization, Measures, Person, Time, Fault and Rock. We have used corpus based machine learning technique to recognize, classify, and identify these geological entities.

In general NER is a hard problem. Words have different applications and there is an infinite number of proper names. In English, the problem of identifying NEs is solved to some extent due to capitalization feature. Most of the named entities begin with a capital letter which is a discriminating feature for classifying a token as named entity. Geological documents contain spatial references and geo references in text. Major task in handling geographical references in text is Name Resolving.

3. Characteristics of geological text

Geological documents contain textual description of geological phenomena, images and maps of geographic space in the form of spatial references, geo references and temporal information. Geographic references can be defined spatially using a point (ex. longitude and latitude) or a set of points. The information in the textual document such as place name and the corresponding linked geographic location is called geographic footprint. Geographic footprint is represented by coordinates (longitude, latitude)

3.1 Features used for Geological NER

Different features may be used for identifying NE's. The features aids in deciding to which class a named entity belongs. The main features for the NER task have been identified based on the different possible combination of available word and tag context. The features also include prefix

and suffix for all words. The term prefix/suffix is a sequence of first or last few characters of a word, which may not be a linguistically meaningful prefix.

$$F = \{ W-2, W-1, W_i, W+1, W+2, |\text{prefix}| \leq 3, |\text{suffix}| \leq 3, \text{POS tag, Digit information, NE tag} \}$$

Context word feature: Previous and next words of a particular word can be as a feature.

Word prefix: A fixed length prefix of the current and/or the surrounding word(s) can be used as features.

Word suffix: Word suffix information assists in identifying NEs. This feature can be used in two different ways. The fixed length word suffix of the current and/or the surrounding word(s) can be used as a feature. For example, suffixes like -pur, -bad, etc are indicators of a name of a location.

Part of Speech (POS) Information: The POS of the current and/or the surrounding word(s) can be used as features.

Digit features: Several binary digit features have been considered depending upon the presence and/or the number of digits in a token (e.g., ContainsDigit [token contains digits], FourDigit [token consists of four digits], TwoDigit [token consists of two digits]), combination of digits and punctuation symbols (e.g., ContainsDigitAndComma [token consists of digits and comma], ContainsDigitAndPeriod [token consists of digits and periods]), combination of digits and symbols (e.g., ContainsDigitAndSlash [token consists of digit and slash], ContainsDigitAndHyphen [token consists of digits and hyphen], ContainsDigitAndPercentage [token consists of digits and percentages]). These binary valued features aids in recognizing miscellaneous NEs such as time expressions, date expressions, percentages, numerical numbers etc.

Named Entity Information: The NE tag of the current or previous word can be considered as the feature.

4. The proposed kernel for SVM

SVM uses a line or surface to separate the data. Thus, SVM is suitable for binary classification problems but not multiple-class problems where there are more than two candidate objective classes. In most cases, name entity recognition is a multiple-class task. As a result, the initial binary SVM is not fit for most name entity recognition tasks. We can use two main types of approaches to solve multiple-class problems. One is to update an SVM kernel function that can merge the multiple classification surface problems into an optimization so as to solve multiple class classification in one pass. The alternative is to apply multiple binary classifiers until they finish the job

The linear SVM computes the dot product between instances

$$K(X, Y) = \phi(X, Y)$$

If x_1, x_2, \dots, x_n are features of X, and y_1, y_2, \dots, y_n are features of Y

$x_i \cdot y_i = 1$ if x and y are same

$x_i \cdot y_i = 0$ otherwise

However when we are dealing with word features and other string based features, such dot product based similarity computation is not able to capture the NER task specific similarity between the strings. For example, the words 'Prof.' and 'Chairman' have some similarity in the context of the NER task as both occurs frequently at the preceding position of the person names; 'small' and 'large' are related words, both being adjectives used in similar contexts; 'town' and 'district' have similarity as both of these are common location terms and occur frequently at the surrounding positions of the location names. Such task specific similarity is important not only in word features but also in other string features like suffix, prefix and n-grams.

We have attempted to capture this semantic similarity with the distance between the instances. It is based on fuzzy clustering information. We have used these similarity functions as kernel in SVM. These individual functions are combined with a suitable weight and the combined function is also used as a composite kernel.

5. Fuzzy clustering based kernel for NER

In the clustering based kernel we use cluster information of the feature values (e.g., words) as a measure of similarity between them. Cluster information has been used in different NLP tasks in the past. Several types of clustering techniques (e.g., Brown et al., 1992; Pereira et al., 1993;

Ushioda, 1996; Biemann, 2006) have been proposed and used in various NLP tasks. Miller et al. (2004) used the hierarchical word clustering algorithm proposed by Brown et al. to extract binary string representation of the words which were encoded in features that are incorporated in a discriminatively trained name tagging model.

5.1 Fuzzy Logic

In *fuzzy logic* based systems, the collection of data are considered as *fuzzy sets*.

Traditional crisp sets include or do not include an individual element; there is no other case than true or false. But in the case of fuzzy sets, it allows partial membership. Fuzzy logic thermostats to control the heating and cooling based on the linguistic terms like cold, moderate, warm and hot

$$\mu = f(s, x)$$

μ : is the fuzzy membership value for the element

x : is the value from the underlying domain.

5.2 Fuzzy c-means clustering

Non-unique partitioning of the data in a collection of clusters. Crisp clustering techniques have difficulties in handling extreme outliers. Fuzzy clustering, each data point will have an associated degree of membership for each cluster in the range [0-1], indicating the strength of its association in that cluster

Membership function is

$$\mu_j(x_i) = \frac{\left(\frac{1}{d_{ji}}\right)^{\frac{1}{m-1}}}{\sum_{k=1}^p \left(\frac{1}{d_{ki}}\right)^{\frac{1}{m-1}}}$$

New cluster centres are computed as

$$c_j = \frac{\sum_i \mu_j(x_i)^{\frac{1}{m-1}} x_i}{\sum_i \mu_j(x_i)^{\frac{1}{m-1}}}$$

6. Clustering of words

Here we have used the Fuzzy clustering algorithm. The input to the algorithm is a list of words to be clustered and a large raw corpus (we have used a raw corpus containing 2 lakh words). The output from the clustering algorithm provides the average distance from cluster members to the center of each cluster.

During kernel computation to obtain the similarity between the words. The similarity is obtained from the distance from the center of the cluster

6.1 Kernel computation

The distance between two string values (e.g., words) is also can be computed from the distance. Finally the individual distances are combined in a weighted fashion to obtain the kernel value of a feature group.

$$K(X,Y)=\lambda \sum_i \phi(X,Y)$$

7. Conditional Random Fields

CRFs are often used for the labeling or parsing of sequential data, such as natural language text or biological sequences. CRFs work well in named entity recognition tasks. Many features can be used in CRFs. For example, term appearance (e.g., capitalization, affixes, etc.) and orthographic features (e.g., alphanumeric characters, dashes, Roman numeral characters, etc.) are used frequently.

Conditional Random Fields (CRFs) are undirected graphical models used to calculate the conditional probability of values on designated output nodes given values assigned to other designated input nodes. A conditional random field (CRF) is a type of discriminative probabilistic model used for the labelling sequential data such as natural language text. Conditionally trained CRFs can easily include large number of arbitrary non independent features. The expressive power of models increased by adding new features that are conjunctions to the original features. When applying CRFs to the named entity recognition problem an observation sequence is the token sequence of a sentence or document of text and state sequence is its corresponding label sequence.

In the special case in which the output nodes of the graphical model are linked by edges in a linear chain, CRFs make first order Markov assumption and can be viewed as conditionally trained probabilistic finite automata (FSMs)

The conditional probability of a state sequence $s = \langle s_1, s_2, \dots, s_T \rangle$ given an observation sequence $o = \langle o_1, o_2, \dots, o_T \rangle$ is

$$P(s/o) = \frac{1}{Z} \exp \sum_{t=1}^T \sum_k \lambda_k f_k(s_{t-1}, s_t, o_t)$$

Where $f_k(s_{t-1}, s_t, o_t)$ is a feature function whose weight λ_k is to be learned via learning. CRFs define the conditional probability of a label sequence based on total probability over the state sequences

where $l(s)$ is

$$P(l/o) = \sum_{s: l(s)=l} P(s/o)$$

the sequence of labels corresponding to the labels of the states in sequences. Z_o is a normalization factor over all state sequences. To make all conditional probabilities sum up to 1, we must calculate the normalisation factor

$$Z_o = \sum_s \exp \sum_{t=1}^T \sum_k \lambda_k f_k(s_{t-1}, s_t, o_t)$$

The feature functions could ask arbitrary questions about two consecutive states, any part of the observation sequence and the current position. For example a feature function may be defined to

have a value 0 in most cases and have value 1 when s_{t-1} , s_t are certain states and the observation has certain properties.

However, CRFs have many drawbacks. First, CRFs use a limited size of context rather than the whole text because of computational limitation, thereby limiting the contextual information. Second, splitting the context of the whole text into small pieces of context will generally separate inherent relationships among them, and simply combining these pieces of context again cannot reproduce the original context due to the loss of relationships during splitting. For example, a CRF geological term identifier uses a two-word context. The whole text could be split into many pieces of two-word contexts. As a result, the same term in the different places of the text could be tagged with different results due to the variation in the context. However, SVM deals with the whole text so it does not have such restrictions. Third, CRFs are affected by the data distribution. If we want to achieve better results, the data should have an exponential distribution.

8. SVM-CRFs Combined Geological Name Entity Recognition

One of the new research areas in machine learning is combining useful algorithms together to provide better performance or for achieving smooth and stable performance. SVM and CRFs are two conventional algorithms that can deal with named entity recognition tasks well. As stated earlier, the feature context used by SVM is global and it does not have the same constraints as CRFs. SVM is initially the best fit for binary-class tasks and it does not perform well on multiple-class tasks. CRFs generally require more computational time and space than SVMs. Thus, although CRFs have many drawbacks, they are very good at sequential data tagging tasks, which is a typical problem in name entity recognition. Thus, we combined Similarity Kernel based SVM and CRFs because they can complement and facilitate each other.

In our approach, Geological named entity recognition was regarded as a two-step task. The first step was to determine whether a candidate term was a Geological one. If it was a Geological, we determine its class of entity. The first step was a binary classification task where the result was either yes or no. We then used CRFs to infer the type of Geological term. Finally, we merged the results returned by SVM and CRFs, before performing an amendment process.

9. Experimental results and analysis

A NE tagged Geological corpus has been used for NER experiment and it contains geology related information in India. This corpus is split into two sets. One forms the training data and the other forms the test data

They consist of 90% and 10% of the total data respectively. CRF is trained with training data and test data is tagged using CRF model.

More than 2 lakh words have been used as training set for the CRF based NER system. The size of the test file is 23K words and the data is labeled with 17 labels. We have used different standard measures such as Precision, Recall and F-measure for evaluation.

Recall is the ratio of number of NE words retrieved to the total number of NE words actually present in the file (gold standard). Precision is the ratio of number of correctly retrieved NE words to the total number of NE words retrieved by the system. These two measures of performance combine to form one measure of performance, the F-measure, which is computed by the weighted harmonic mean of precision and recall.

	Class	Precision	Recall	F-measure
1	Country	98.75	90.29	94.33
2	State	95.83	94.52	95.17
3	City	81.11	73	76.84
4	Region	82.54	71.23	76.47
5	Mountain	92.31	85.71	88.89
6	Water bodies	84	72.41	77.78
7	Island	78.71	78.57	78.64
8	River	88.89	88.89	88.89
9	Village	70.59	50	58.54
10	Mineral	94.26	80.42	86.79
11	Organization	46.79	91.82	61.99
12	Measures	91.53	92.05	91.79
13	Year	96	97.27	96.63

14	Person	68.38	80.69	74.03
15	Fault	33.33	14.29	20.00
16	Rock	71.15	75.51	73.27
17	Time	35.71	76.92	48.78
	Overall	77.05	77.27	75.81

Experimental Results of NER using CRF(baseline)

We have got Precision of 77.05%, Recall of 77.27% and F-measure of 75.81% by the combination of features (prefix and suffix of length up to three of the current word, information about the surrounding words, POS information, digit features, and NE tag) for identifying named entities.

	Class	Precision	Recall
1	Country	99	91
2	State	96.3	95
3	City	96	82
4	Region	98	93
5	Mountain	96	92.4
6	Water Bodies	86	86
7	Island	85	85
8	River	90.11	90.11
9	Village	85	85
10	Mineral	94.3	77
11	Organisation	47.5	47.5
12	Measures	93.2	93.2
13	Year	97	97
14	Person	69.11	69.11
15	Fault	34.53	20.1
16	Rock	90.2	92.1

17	Time	36.7	36.7
		81.99706	78.36588

Experimental Results of hybrid approach with K-Means Clustering

We have got Precision of 81.99%, Recall of 78.36% for the combined approach. NE's such as Country, State and Year have high F-measure values because of their higher appearance in the corpus

	Class	Precision	Recall
1	Country	99.5	92.5
2	State	98.3	97
3	City	97	85
4	Region	99	94.5
5	Mountain	97	94.4
6	Water Bodies	88	87
7	Island	86	87
8	River	91	91.11
9	Village	87	86
10	Mineral	96.4	78
11	Organisation	48	48.5
12	Measures	95	95
13	Year	98	98
14	Person	71.33	71.75
15	Fault	35.66	23
16	Rock	92	94.2
17	Time	37.3	37.3
		83.32	80.02

Experimental Results of hybrid approach with Fuzzy Clustering

Realworld data is uncertain and vague. Crisp clustering techniques have difficulties in handling extreme outliers. Fuzzy logic is capable of supporting, to a reasonable extent, human type reasoning in natural form by allowing partial membership for data items in fuzzy subsets. Due to this we have got Precision of 83.32%, Recall of 80.02% for the combined approach with fuzzy clustering. NE's such as Country, State and Year have high F-measure values because of their higher appearance in the corpus.

10. Conclusion

In this paper, we have developed a NER system using SVM and CRF with the help of a NE tagged Geological Corpus. We also presented a new named entity tagset that was developed for annotation of this corpus. We have considered features such as prefix and suffix of length up to three of the current word, POS information, digit features, information about the surrounding words and their tags. This proposed method of hybrid approach with fuzzy clustering has obtained better accuracy than CRF based method (baseline) and hybrid approach with K-Means clustering .

References

- [1] Wakao, T., Gaizauskas, V. and Wilks, Y. 1996. Evaluation of an algorithm for the recognition and classification of proper names, In Proceedings of COLING-96.
- [2] Singh, A. K. and Surana, H. 2007. Can Corpus Based Measures be Used for Comparative Study of Languages, In Proceedings of Ninth Meeting of the ACL Special Interest Group in Computational Morphology and Phonology, ACL 2007.
- [3] Bikel, D. M. , Schwartz, R. L. and Weischedel R. M. 1999. An Algorithm that Learns What's in a Name. Machine Learning, pp. 211-231.
- [4] Borthwick, 1999. Maximum Entropy Approach to Named Entity Recognition, Ph.D. thesis, New York University.
- [5] Lafferty, J. D., McCallum, A. and Perera, F. C. N. 2001. Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data, pp.282-289, ICML 2001.
- [6] Li W and McCallum A. 2003. Rapid development of Hindi named entity recognition using conditional random fields and feature induction, ACM Transactions on Asian Language Information Processing (TALIP), pp.290–294.

- [7] Srihari, R., Niu, C. and Li, W. 2000. A Hybrid Approach for Named Entity and Sub-Type Tagging, In Proceedings of the sixth conference on Applied natural language processing.
- [8] Toral, A., Noguera, E. Llopis, F. and Munoz, R. 2005. Improving question answering using named entity recognition, In Proceedings of the 10th NLDB congress, Lecture notes in Computer Science.
- [9] Spain, A., Babych B. and Hartley. 2003. A. Improving machine translation quality with automatic named entity recognition, Springer-Verlag 2003.
- [10] Wallach, H. M. 2004. Conditional random fields: An introduction, Technical Report MS-CIS-04-21, University of Pennsylvania, Department of Computer and Information Science, University of Pennsylvania.
- [11] Zhu, Fei, and BairongShen. "Combined SVM-CRFs for Biological Named Entity Recognition with Maximal Bidirectional Squeezing", PLoS ONE, 2012
- [12] Saha, S.K.. "A composite kernel for named entity recognition", Pattern Recognition Letters, 20100901
- [13] Sobhana, N. V, S.K Ghosh, and PabitraMitra. "Entity Relation Extraction from geological text using Conditional Random Fields and subsequence kernels", 2012 Annual IEEE India Conference (INDICON), 2012
- [14] Takukudo. 2005. CRF++, an open source toolkit for CRF, <http://crfpp.sourceforge.net>.