# Domain Analysis of Information Extraction Techniques

**Talha Mahboob Alam[1] and Mazhar Javed Awan[2]**

[1]Computer Science and Engineering Department, University of Engineering and Technology, Lahore, Pakistan
[2]Department of Computer Science, University of Management and Technology, Lahore, Pakistan
[1]talhamahboob95@gmail.com, [2]mazhar.awan@umt.edu.pk

*Abstract*— In this research, we extant a short outline of Information Extraction, which is also a natural language processing domain that tries to find required information in structured, semi structured and unstructured Data. We draw a taxonomy of information extraction tasks and techniques. The other important thing is that we also extract learning methods like supervised, semi supervised and unsupervised learning and which methods are used in these types of learning. Our domain analysis consists on social media, Biomedical, chemical and unstructured data. There are different tasks included in information extraction which makes this activity more manageable as well as to easy to work in specific domain. We also detect weakness of existing techniques.

*Keywords*— Biomedical, Information Extraction, Supervised, Semi Supervised, Social Media and Unsupervised

## I. INTRODUCTION

Information extraction (IE) is the task to consequently separate structured data from unstructured and semi-structured data. There were two tasks in IE .The primary task is Name Entity Recognition, which is characterized as "the procedure of find certain terms in the processing of content like person, location, organization etc". The second task is Small expression co-reference determination, which is "the procedure whether two expressions in the characteristic Language refer to a similar substance and the arrangement of references by pronouns and relation between entities is known as relation extraction" [1].
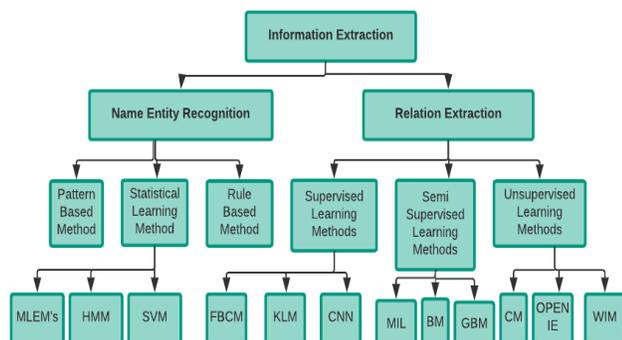


Fig. 1: Taxonomy of IE

## II. DOMAIN ANALYSIS

There are four domains in which we have done our research. Social media, Biomedical, Chemical as well as unstructured data sets are used to done our domain analysis. We also analyze that which method is not used in specific domain shows in tabular form. The data in these specific domains which we chose are developing violently and helpful outcomes are appearing every day and increases in future. A huge amount of datasets of these domains are accessible on the web. If we take the example of social media and unstructured data, online social networks bleed information as well as vast range of data is on internet in the form of unstructured data.

Each time somebody posts something on Twitter and Facebook and as well as a large amount of raw data is uploaded in the form of blogs, newspapers and e-mails. If we take the example of bio medical and chemical data, the new generation of sequencing technologies enables the processing of billions of DNA sequence data per day, and the application of electronic health records (EHRs) is documenting large amounts of patient data as well as number of compounds are increased in the shape of industrial and synthetic reports.

| Information Extraction | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mame Entity Recognition | | | | | Relation Extraction | | | | | | | | |
| Pattern Based | Rule based | Statistical Learning | | | Semi Supervised | | | Supervised | | | Unsupervised | | |
| | | MLE'm | HMM | SVM | MIL | BM | GBM | KLM | CNN | FBCM | Open IE | CM | WIM |
| Social Media Data | | | ✓ | ✓ | ✓ | | ✓ | ✓ | | | ✓ | ✓ | ✓ | ✓ |
| BioMedical Data | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | |
| chemical Data | | ✓ | | ✓ | | ✓ | ✓ | ✓ | | | ✓ | ✓ | ✓ | |
| Unstructured Data | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ |

Fig. 2: Domain analysis of different techniques

### A) Named Entity Recognition (NER)

NER is a subtask of IE that tries to find and characterize named entities in content into pre-characterized classes, for example, the names of people, organization, location, date, time and amounts. For instance, "UMT" mention to the individual "University of management and technology" location "Lahore", or some other Entity that has a similar shortened form. To decide the way of the organization for

"UMT", which happens in a specific report, must be its location Considered.

### 1) Pattern Based Method (PBM)

Old-fashioned IE frameworks first develop a pattern and after that utilized the word reference to extract the important data from the new unlabeled data [2].

Thorsten et al. introduce a PBM for named entity recognition using unstructured data. Utilizing this technique, they developed a few classes of seed entity records into significantly high accuracy records and enhanced the precision of a depending irregular field-based named entity tagger [3]. Valery et al. purposed a pattern based recognition technique to find organization names in the text of news which is in Russian language. Using business news, the recognition method depends on use of local context signs of information and inside name words [4].

PBM in social media for the sake of entity recognition is not easy to use. There are many reasons like the data is not normalize like if we take the example of Twitter "*Talhaaa@Umtlahore*" the data will be normalized before recognition. Real time systems are very difficult to build like in the field of medical. The life of patients is related to results of tests if results are not accurate then system will be collapsed soon but Amita et al. fundamentally, the suitability of the use of pattern-based recognition in the analysis of salivary gland (SG) injuries. Optionally, an attempt was made to concentrate the cytomorphology of the different sores in detail and talk about the drawbacks and arrangements required in the testing conditions at cytology [5].

Pattern based technique is not used for chemical text especially compound and element detection. One possible issue is that the dictionary of is too large because the number of molecules and elements is huge that's why PBM is not used in substance or element recognition [6]. Other methods like condition random field (CRF), rule-based method and several other techniques are used for chemical name entity recognition [7], [8].

### 2) Rule Based Method (RBM)

The RBM is closely related to decision-making trees, except that they trees do not create a strict hierarchical division of the training data. Rather, overlaps are allowed to create more robust model for training.

Daniel et al. purposed a RBM for protein and gene entity recognition. The ProMiner framework utilizes a pre-handled equivalent word lexicon to distinguish potential names in the biomedical content and the expanded ProMiner framework has been connected to the experiments of the BioCreAtIvE challenge with exceptionally supportive outcomes [9]. It is very difficult to label or recognize entities from unstructured data But Rayner Alfred et al. suggests a RBM of Named-Entity Recognition for Malay articles. The proposed Malay NER is planned in view of a Malay part of speech (POS) labeling highlights and relevant components. A few physically developed word references will be utilized to deal with three named entities; Person, Location and Organization [10]. Baichuan et al. used the twitter to extract tweets because it is very necessary to examine the novel issue of automated question detection proof in the microblog condition. To extract tweets, setting highlights like short url's and Tweet specific highlights like Retweets are extravagantly chosen for characterization [11]. Masaharu YOSHIOKA Thaer M. DIEB exhibit an outfit way of Chemical Named Entity Recognizer (CNER) instruments with various attributes like RBM and machine learning method. What's more, he utilized text segmentation to normalize the text as well as detect the boundaries of sentence [12].

### 3) Statistical Learning Method

Information extracting using Statistical learning model, which emerged in recent years as new probabilistic framework of the process of extracting information. Generally modeling statistical language or simple modeling language (SML), refers to the task of a probability distribution of the estimation that captures the statistical rules of natural language training [13].

***Maximum Likelihood Estimation Method (MLEM's):*** NER in English is easy but it is very difficult in Chinese language because words are joined in sentences but Jian et al. used MLEM to find the probabilities in unstructured data. They consider the issue of Chinese NER distinguishing Statistical model (SM) and figure the probabilities through MLEM. Word segmentation and NER have been incorporated and got a structure together that comprises of few class-based models [14]. Zhuo et al. worked on Biomedical NER using Maximum likelihood estimation to calculate probabilities. For enhancing the performance and decreasing the training time, to find ordered pairs use maximum likelihood in the training of process model. The parameter estimation process is slow but using Conditional random fields with MLEM it gives better results [15]. MLEM was not wildly used for Chemical name entity recognition because the detection of entity is not easy task due to the inclusion of different formulas as well as same terms which gives different meanings and vice versa. Kamal Sarkar purposed an approach, utilizing MLEM for Entity recognition from Social Media Text. Created trigram HMM tagger requires to handle label trigram likelihood, which is figured by the most extreme likelihood estimation from name trigram entities. To conquer the information unavailability issue, label trigram likelihood is smoothed utilizing extraction method which utilizes the most extreme probability measures from means for tag trigram, tag bigram and tag unigram [16].

***Hidden Markov Models (HMM):*** NER in Natural Languages like Kannada is a fundamental and exciting task particularly for unstructured information. S Amarappa and S. V. Sathyanarayana purposed a method using HMM for unstructured data in Kannada language for NER. The point of their work is to build up a Hybrid model using HMM and rule based model. The outcomes are talked about utilizing 100's of tests and results are impressive [17]. In spite of the fact that there exists an enormous number of biomedical data on the web, there is an absence of tools adequate to help individuals to get data or learning from them. Shaojun Zhao present a HMM for NER, with a word comparability based smoothing. Their test demonstrates that the word similarity based smoothing can enhance the execution by utilizing immense unlabeled information [18].

online social networking perform ineffectively because of an absence of dependable capitalization, irregular sentence structure and an extensive variety of vocabulary but Chun-Kai et al. propose a novel generative model that joins the concepts from HMM and n-gram language models into what we call an N-gram Language Markov Model (NLMM). even with this straightforward model, way to deal with NER on casual content beats existing frameworks prepared on formal English and matches best in class NER frameworks prepared close by hand-labeled Twitter messages [19]. Scientific data is exceptionally rich and may require chemical entities. A standout amongst the most vital difficulties in NER is the names of chemical formulas like *Lead* is another symbol *PB*. Ms. Snehal and Dr. Neeta Extract Chemical Names utilizing HMM. The forward-backward algorithm is used by HMM model, Viterbi Algorithm and Estimation-Modification technique for validation. To represent joint probability over observation and tagging order HMM needs number of all likely observation sequence [20].

*Support Vector Machine (SVM):* SVM methods use linear conditions to separate some other classes. The idea is to use a direct state that divides the two classes as much as possible from each other.

Joel mickelin purposed a framework for NER in Swedish text and recognize names of individuals, location of companies, and specific time. This framework was developed from the POS tagger and the SVM framework. The framework was prepared to perceive Named Entities by breaking down examples in training corpora comprising of arrangements of words having a place with every class [21]. Zhenfei Ju et al. utilized SVM for NER in biomedical data. By utilizing SVM, require a few components to direct recognizable contents. There are many components can be utilized as a part of SVM including two sections word nature and grammatical feature [22].

Existing NER methods, commonly intended for formal data written in standard language (e.g. articles, newspapers and statistics) don't perform well on user oriented data like tweets, pins, Facebook status as well as comments. Kresimir et al. used the technique of SVM for NER by using Croatian Tweets. The organized SVM reliably blows the baseline and the HMM Model but is likewise reliably outperforms by the CRF. The most well-known reason for mistakes for the organized SVM model are tokens marked as within a named entity (e.g., I-PER) anyhow when the previous token was not the start of a named entity (e.g., B-PER) [23]. The major disadvantage of SVM methods is that they are slow. However, they are very popular and tend to have high accuracy in many practical domains such as text.

### B) Relation Extraction

A relationship is different as a predicate that extends over two arguments, one argument representing ideas, entities, or persons in the world of nature, and the predicate relationship defines the union or collaboration among the entities signified by the arguments. Another imperative task of IE is the extraction of relationships. Extracting relationships is the task of representation Semantic connections between entities in the content [13].

### 1) Supervised Learning Methods

If data is supervised, there is a result for all observations. The learning goal is to build a prediction-based model on a series of observations with known results to predict the results of the new entities, when presented to the system. The input data are also called learning or training data, while new observations are called test data. To evaluate the validity of the model test data can be used since it was not used in the learning phase.

***Feature Based Classification Method (FBCM):*** RBM focuses on a less scope of data to explain the relation between different entities in the dataset as standards and formats. The FBCM utilizes a wide range of contextual information [24]. Learning of relations between drugs is significant for medical experts to keep away from unfriendly impacts when co-controlling medications to patients. Quoc-Chinh Bui et al. purposed a FBCM to extract relations between drugs from biomedical text. Their Method comprises of three stages. To start with, apply data preprocessing to change over sentences from a given dataset into structured representations. Second, map the drug-drug relation and then combine from that dataset into an appropriate syntactic structure. In light of that, a novel set of features is utilized to produce feature vectors for these drug-drug relation sets. Third, the acquired feature vectors are utilized to prepare a SVM classifier [25]. Conventional FBCM are unfitted for data of social media as far as exploiting the extra data in linked information. So, make the unique relations that can be extracted from linked data is known as relation extraction. Jiliang Tang and Huan Liu used FBCM for relation extraction using social media. Outline and direct analyses on datasets from social media and the exact outcomes show that the proposed system can altogether enhance the execution of feature selection [26].

The extraction of chemicals as well as diseases and their relations from unstructured scientific data is important for some regions of research. The manual extraction of these elements and relations, and their storing in organized databases is unwieldy and costly. E. Pons et al. extract the relations between chemical and disease using FBCM. The task of Relation extraction as a similar choice issue on every conceivable match of chemicals and infections found in each report. For each occurrence, three sorts of features were produced, based on prior knowledge, and on statistical and semantic data from the records [27]. Various NER frameworks were proposed in a few languages, for example, English, Chinese and Arabic. To utilize the extracted name entities, it is important to distinguish relations among them. Named entities extraction in Thai is not quite the same as those in different languages. Nattapong Tongtep and Thanaruk Theeramunkongp present a model, utilizing FBCM for separating relations among named entities from Thai news documents. Four supervised learning methods are combined, on the other hand to explore the execution of relation extraction utilizing distinctive feature sets [28].

In the relationship extraction task, the problem with FBCM is that sometimes data with explicit feature vectors cannot be represented easily. In these cases, the extraction of characteristics is a very complex task and leads to vectors of very high dimension, which in turn leads to problems of

calculation. Kernel-based methods attempt to solve this problem implicitly rely on scalar vector products are calculated at very high Dimensional spaces without any kind of vector must be explicit [29].

***Kernel Learning Methods (KLM):*** KLM provides a powerful and consistent framework for all these disciplines, motivating algorithms that can act on general types of data such as text and for the general types of relationships e.g., classification, relation extraction, and pattern recognition. Dmitry Zelenko et al. purpose a method for relation extraction through KLM from unstructured sources. The nodes of the parsed tree have properties and need to utilize the qualities in characterization. First, define a pairing function and a function of similarity to the node. The defined matching node determines whether the node or not. On account of relationship extraction, nodes are just Parse trees if their types and functions coordinate. The function of similarity at the node is processed as far as the properties of the nodes [30]. It is in sensible to assume that kernel based particularly tree-based methodologies are not reasonable for Chinese language, in any occurrence at current stage. Most existing biomedical relation extraction techniques require manual formation of biomedical dictionaries or parsing layouts in view of learning area. Jiexun Li et al. intend KLM to consequently extract biomedical relations from medical content. To build up a system of KLM for biomedical relation extraction, this framework can be divided into four subtasks: entity recognition, relation annotation, kernel construction, learning and evaluation [31].

***Convolutional Neural Network (CNN):*** Up to now, relation extraction frameworks have made broad utilization of components created by analysis of data. The solution to the problem of relation extraction can be obtained by the use of CNN. Relation extraction from biomedical data, for example, examines articles, magazines, or database records have been a subject of many research activities and shared difficulties. Relation extraction is the way toward identifying and ordering the semantic relation among entities in a given biomedical data. Sunil Kumar et al. purposed a Relation extraction from biomedical data using CNN. The proposed model, which brings a comprehensive sentence with specified entities as information and produces a probability vector comparing to all conceivable relation types. Each feature is having vector illustration which is established randomly with the exception of word embedding feature. The proposed model has demonstrated better execution by SVM based baseline model [32]. Ji Young Lee et al. purposed a technique for Relation Extraction through CNN for unstructured text. Their model consists of three parts: preprocessing, convolution neural network (CNN) and post-processing [33].

## 2) Semi Supervised Learning Methods

Semi-supervised learning states to the practice of labeled and unlabeled data for training. The semi-supervised learning is middleware between supervised and unsupervised learning. In adding to untagged data, the algorithm is provided with certain supervised information, but not essentially for all cases. KLM's, as well as Feature-based relation Extraction in

a portion of training data, which is expensively obtained. One answer for this issue is softly focused learning Approaches that work with considerably fewer Training information.

***Multi Instance Learning (MIL):*** In the MIL, the training set is composed of many sections, each with many cases or bags. MIL emerged as an extension of semi supervised learning. Wei Wang et al. depicts MIL model, its ability to detect news articles about civil unrest events (from Spanish text) across ten Latin American countries and identify the key sentences relating to these events. The model, trained without annotated sentence labels [34]. Yu-Ju Chen and Jane Yung-jen Hsu purposed a model of Relation Extraction by MIL for Chinese language when data is unlabeled. MIL algorithms find occurrences in the sections then occurrences used for training are characterized by the sections to which they belong. In their work, stMIL is used to predict the relationship of sentences because the positive section used for training, are rare positive [35].

Since each compound comprises of numerous parts, at that point pick up benefits by the enhanced heuristics for massive scale information. This makes existing methodologies perform ineffectively on extensive scale word datasets. Cholwich Nattee et al. make the framework that handle a huge search space effectively by determining particular data into search heuristics. Presently, heuristic capacities utilized as a part of MIL frameworks are construct just with respect to quantitative data. After this attention on a sort of information comprising of a few sections. Finally, introduced an enhanced heuristic function for an information comprising of numerous parts [36]. Numerous biomedical relation extraction approaches depend on supervised machine learning, requiring an explicated corpus however biomedical relation extraction technique in light of MIL depends on semi-supervised Method. Distinguishing Biomedical Relations, which does not require actually explained corpus is difficult task. Andre Lamurias et al. purposed a technique which depended on the sparse MIL algorithm, used to prepare on a consequently produced corpus of 4,000 documents identified with miRNAs. [37].

***Bootstrapping Method (BM):*** BM used to extracting relationships, have gained considerable attention in recent years. These approaches are constructed with a fundamental assumption that if you know a couple of words that relate in a certain way, phrases containing these words, these type of words likely make relationships to express. Therefore, the sentences containing the word pair are used as training data for the relation extraction [38].

Eugene Agichtein and Luis Gravano consider a technique for extracting tables from unstructured plain text that requires just an uncertain group of training cases from users. At every iteration of the extraction procedure, Snowball assesses the nature of these examples and tuples without human intervention and keeps just the most dependable ones for the following cycle. Likewise build up an adaptable assessment system and measurements for the task and present a careful test assessment of Snowball and practically identical strategies over a collection of more than 300,000 daily paper archives [39]. Xuezhong Zhou et al. used BM for relation extraction in bio medical text. To keep the BM robust while

extracting the new seed tuples, utilize a dynamic bubble-up evaluation to guarantee that the excellent examples and seed tuples will add to the iterative procedure. This bootstrapping strategy is called bubble-bootstrapping [40].

Marco Pennacchiotti and Patrick Pantel purposed a method which is named as ESPRESSO, a semi supervised iterative BM joined with the web-based learning extension method, for extricating parallel semantic relations in a given chemical corpus. Espresso is proposed to extract different semantic relations showed by a given slight training of seed instances. For instance, hydrogen gas responds with oxygen gas and zinc responds with hydrochloric acid. At that point assess this relation on the CHEM corpus [41]. Jinghang Gu Et Al. purposed a method for relation extraction in social media networks. Family relations are the ones between people in a similar family, including spouse-wife, father-child, mother-child, fellowship and sisterhood. BM is used for decreasing the measure of training set. When separating examples and instances, there are five sorts of family relations accessible for bootstrapping [42].

***Graph Based Method (GBM):*** BM are well-known for the extraction of relationships, mainly since they require just a little amount of human assessment. Graphs can represent complex relationships between classes and instances. An ambiguous instance; For example, Usman Khawaja might be among the class of pilots and players.

Today social media like Facebook beats the web. Researchers are working on the analysis of social media for the sake of relation extraction from comments, tweets as well as posts. Meesun Song et al. used Directed graph for relation extraction. The proposed method is a dynamic radial graph to cope with the limitations of previous visualization techniques. Dynamic radial graph with the extracted social relationship form a Facebook dataset. The basic layout is similar to the radial layout [43]. The ability to precisely grasp both semantic and syntactic structures in biomedical data turns out to be progressively basic and allow intense understanding of scientific papers and clinical data. Yuan Luo et al. used graph method for relation extraction, in which nodes are entities and edges demonstrate relations or different entities associated by numerous edges can be viewed as one relation [44]. Yifan Peng et al. propose Extended Dependency Graph (EDG) by combining a couple of basic verbal concepts and incorporate data beyond composition or syntax as well as hope that the practice of EDG will empower machine learning frameworks to sum up more effortlessly. Results affirm that EDG gives up better results over KLM [45]. Gokhan Bakal and Ramakanth Kavuluru purposed a framework to report introductory results on calculating treatment relations between biomedical entities, absolutely in light of semantic patterns over biomedical knowledge graphs. The purposed instinct is genuinely clear – entities that take an interest in a treatment connection might be associated utilizing comparable way designs in biomedical knowledge graphs separated from scientific data [46].

G. J. Postma et al. purposed a framework which depends on the hypothesis of Government and Binding for the syntactic part and Conceptual Graphs for the analysis. The data is consisting of chemical abstracts which are online. Their semantic and important investigation depends on Conceptual Graphs [47]. Kundan Kumar and Siddhant Manocha implement a semi supervised method for relation extraction when the data is unstructured. Developing a knowledge graph include extracting relations from unstructured content taken after by effective storage in graphical databases. The framework which is purposed for extracting relations utilizing semantic normality in the distributed word vector embedding space. Therefore, additionally construct a question answering framework that parses the normal dialect inquiries utilizing using regular expressions and extracts answers from the knowledge graph. [38].

### 3) Unsupervised Learning Methods

In unsupervised learning there is no such supervised or structured data and we just have input information. The objective is to discover predictabilities at the occurrence. There is a structure for the information space so that specific examples happen more as often as possible than others. Unsupervised learning can be thought of as discovering patterns in the information well beyond what might be viewed as clean unstructured noise.

***Clustering Method (CM):*** Due to the large number of relationships between entities, it may be costly to cover a sufficiently large amount of training data to effectively mark each type of relationship in each new domain of interest. There are few challenges that why we selected CM. Firstly, the same semantic relationship expressed between two entities can be used different lexical or syntactic patterns. CM can be considered important learning without overseeing the problem; just like any such problem, try to find a structure in a collection of unlabeled data as well as it is very important unsupervised technique. Entities are grouped according to their categorical information. With the goal of using social media networks for the Semantic Web, a few reviews have analyzed automatic relation extraction of social media. Junichiro Mori et al. purposed a method for extraction of relation in social media incorporates the following steps. Gather co-occurrence data and local context of pair of entities. The purposed technique requires a list of instances (e.g., individual name, area name) to frame a social media network as the information [49]. The abundance of support data given in biomedical articles aroused the execution of information extraction extra ways to deal with biomedical relation extraction. Changqin Quan et al. displays an unsupervised technique which consists of clustering and sentence parsing to manage biomedical relation extraction. The proposed unsupervised method combining pattern clustering, dependency parsing and phrase structure parsing rules. Two methodologies on two unique tasks are accessed: Protein–protein relation extraction, and Gene–suicide relation [6]. Wikipedia is a most important database. It is a multilingual, brilliant information base, and an essential source of organized learning, which can be extracted and examined by machine. Song Liu and Fuji Ren using CM for the relation extraction of Wikipedia which is type of unstructured data. The entities are clustered as well as labeled, with the providing clusters giving data repetition that advantages the relation extraction. This method requires negligible manual

help and concentrates numerous relations. The outcomes are practically identical with different works [50].

***Open Information Extraction (Open IE):*** IE systems try to extract the semantics of text in natural language relationships, but most systems use supervised specific examples of the relationship to learn and therefore partial by the ease of use of training data. On the other hand, open IE systems such as TextRunner, show unlimited number of relationships found on the Web to handle [51].

Ndapandula T. Nakashole purposed a framework for extracting arrangement of relations from social media corpus. PATTY (the proposed method) learns textual patterns that mean combined relations. From Wikipedia, PATTY learned 350,569 example synsets [52]. To concentrate on a general chart of relations, there is no labeled corpus reasonable for taking in the extraction model. Open IE framework, depends entirely on the information content and its linguistic qualities. All the more particularly make examples to catch these qualities of content and afterward separate relations. Nhung T. H. Nguyen et al. propose an Open IE framework, which goes to concentrate relations of any sort from biomedical content. The framework from existing Open IE frameworks is that it utilizes predicate argument structure examples to identify the competitors of conceivable biomedical facts as well as physically assessed the performance of framework and found that it is sensibly exact. This framework connected with the entire MEDLINE and open that the relations between "Amino Acid, Peptide or Protein" substances are the most normally visible sort of relations [53]. Fei Wu and Daniel S. Weld worked on Open IE using Wikipedia. The underlying WOE has three main components: preprocessor, matcher, and learner. The preprocessor converts the plaintext of Wikipedia into a sequence of phrases, attributes annotations NLP and builds sentences of synonyms for the key elements. The preprocessor makes any Wikipedia article into HTML then splitting the article into sentences. A final step builds the pre-processor by equal phrases to find phrases according to infobox relationships. After this, builds training data for the component learner couple doing games heuristics values and Wikipedia article with infobox corresponding sentences in the article [54].

***Wrapper Induction Method (WIM):*** A wrapper is an extraction strategy that includes an arrangement of extraction rules and program code required to execute these rules. WIM learn the wrapper automatically. Given a sequence of training data, the induction algorithm figures out how to extract a wrapper target data [2].

Dingli Alexiei et al. purposed a method of Automatic semantic annotation using unsupervised information extraction and integration. The system keeps on iteration till there's no additional info to get or the user decides to interrupt the cycle. From the results, personal websites area unit recognized with easy heuristics. Essentially the system keeps on iteration till there's no additional info to get or the user decides to interrupt the cycle. Wrapper induction systems are extended to deal with less rigidly structured, free texts and even a mix of them [55]. Anna Lisa Gentile et al. propose a direct knowledge based strategy which is extremely adaptable as for various areas and does not require any training material,

but relatively activities Linked Data as essential learning source to construct basic learning resources. To demonstrate for areas that are secured and Linked Data fill in as an effective learning asset for Information Extraction. Analyze an openly accessible dataset exhibit under specific conditions; this basic un-directed approach can accomplish aggressive outcomes against some mind boggling best in class that dependably relies on upon training information [48].

The following Table I is listed with respect to method, domain as well as weakness of each technique in a specific domain. From table we extract that there were lot of work to be done in future as well as researchers may overcome the deficiencies of existing techniques especially in chemical and bio medical domain due to high production of complex data.

Table I: Weakness of each technique in a specific domain

| METHOD | DOMAIN | REF | WEAKNESS |
|---|---|---|---|
| PBM | Unstructured Data | [3] | Proposed method must require dictionary. |
| | Unstructured Data | [4] | Improvement required. |
| | Bio Medical data | [5] | There are some cases which involves false-negative results. |
| RBM | Social Media Data | [11] | Dataset is less and syntax features did not used. |
| | Biomedical Data | [9] | More than one species in single document cannot be Recognized. |
| | Chemical Data | [12] | Data is labeled manually which is not feasible when data is huge. |
| | Unstructured Data | [10] | The defined rules are less and does labeled all entities like time, Percentage and date. |
| MLEM's | Social Media Data | [16] | This method is only useful for English not any other languages |
| | Unstructured Data | [14] | The effectiveness of POS tagger is low. |
| | Biomedical Data | [15] | Suicide-related genes present in databases does not equate to finding the appropriate Relations within a document. |
| HMM | Social Media Data | [19] | Entity types such as products, movies and songs have not recognized. |
| | Biomedical Data | [18] | It is not evaluated that what happens if it could also be plugged into other existing Systems. |
| | Chemical Data | [20] | Survey is taken just of chemicals datasets. |
| | Unstructured Data | [17] | The domain of study is just NERC research which is too short. |
| SVM | Social Media Data | [23] | Retweets are not considered for NER. |

| Technique | Data Type | Ref. | Limitation |
|---|---|---|---|
| FBCM | Unstructured Data | [21] | Only three categories of NER are recognized which are too short. |
| | Social Media Data | [26] | The combinations of relations, contributions to feature selection and Hidden information in social media are not studied. |
| | Biomedical Data | [25] | Gene–disease relations are not extracted. |
| | Chemical Data | [27] | Did not attempt to explain the relation discussions in the document texts, which might have yielded stronger features. |
| | Unstructured Data | [28] | |
| KLM | Biomedical Data | [31] | The framework did not capture semantic similarity between words as well as named entities have been manually tagged by domain experts before relation extraction. |
| | Unstructured Data | [30] | The processing speed is low when kernels are applied as well as complexity becomes high. |
| CNN | Biomedical Data | [32] | Dictionary as well as training samples are not enormous. |
| | Unstructured Data | [33] | When the model is trained with the any order strategy, the choice of the evaluation strategy does not impact the performance. |
| MIL | Biomedical Data | [37] | The whole Database Created manually which is very time-consuming task. |
| | Chemical Data | [36] | The domain of the Research is not wide, only chemical data is considered. |
| BM | Social Media Data | [42] | The types of family relations are very short as well as the recall is low. |
| | Biomedical Data | [40] | The protein-protein interactions by using bootstrapping does not applied. |
| | Chemical Data | [41] | The value of recall is very low. |
| | Unstructured Data | [39] | Only two attributes are recognized which is too short. |
| GBM | Social Media Data | [43] | This algorithm does not take into account the complex relationship within message threads. |
| | Biomedical Data | [45] | Richer Features are not used as well as EDG does not use for the combinations of features. This method is only tested when datasets are large. |
| | Chemical Data | [47] | This algorithm is too much time consuming. |
| | Unstructured Data | [38] | Regular grammar may not be used to convert simple natural language questions. |
| Open IE | Social Media Data | [52] | Unbounded Relations as well as some entities relations like height, GDP etc. are not considered. |
| | Biomedical Data | [53] | |
| | Unstructured Data | [54] | Freebase does not used for training set. Domain explicit features are not implemented. |
| CM | Social Media Data | [49] | Just small number of social networks is analyzed. |
| | Biomedical Data | [6] | Suicide-related genes present in databases does not equate to finding the appropriate relations within a document, which is one of the limitations of the evaluation approach. |
| | Unstructured Data | [50] | Some other constraints such as dependency and web search results could be introduced to improve the extraction precision. |
| WIM | Social Media Data | [47] | Additionally, analyses are required for checking the conduct of the framework in various setups. |
| | Unstructured Data | [48] | Lack of robustness in the learnt wrappers. |

## III. CONCLUSION

In this paper, we examined the different techniques of IE in the literature and compare them. We have evaluated only for types of data from different domains like social media, biomedical, chemical and unstructured data. There are several points to make the survey. First, we see the trend of IE highly automated systems to be developed, which saves not only effort of programming, but also the marking effort. On the other hand, there is lot of improvement required to work on different techniques as well as on different data. Pattern based approach and Maximum likelihood estimation approach required much more attention because a lot of domains are existing where these approaches are not applied for the sake of NER. Kernel methods, CNN and wrapper Induction are also required considerably more consideration in different domains where these methodologies are not applied for relation extraction. Now we discuss the major outcome from this survey, there is lot work required to work on chemical

data because chemical data required a lot of formulas which is not easy to understand by a machine.

## REFERENCES

[1]. Nadeau, D., & Sekine, S. (2007). A survey of named entity recognition and classification. *Lingvisticae Investigationes*, *30*(1), 3-26.

[2]. Tang, J., Hong, M., Zhang, D., Liang, B., & Li, J. (2007). Information extraction: Methodologies and applications. *Emerging Technologies of Text Mining: Techniques and Applications*.

[3]. Talukdar, P. P., Brants, T., Liberman, M., & Pereira, F. (2006, June). A context pattern induction method for named entity extraction. In *Proceedings of the Tenth Conference on Computational Natural Language Learning* (pp. 141-148). Association for Computational Linguistics.

[4]. Solovyev, V., Gareev, R., Ivanov, V., Serebryakov, S., & Vassilieva, N. Dictionary and pattern-based recognition of organization names in Russian news texts. *in AWERProcedia– Information Technology and Computer Science*.

[5]. Amita, K., Shankar, S. V., Sanjay, M., & Sarvesh, B. M. (2016). Effectiveness of the Pattern-Based Approach in the Cytodiagnosis of Salivary Gland Lesions. *Acta cytologica*, *60*(2), 107-117.

[6]. Tang, Z., Jiang, L., Yang, L., Li, K., & Li, K. (2015). CRFs based parallel biomedical named entity recognition algorithm employing MapReduce framework. *Cluster Computing*, *18*(2), 493-505.

[7]. Dieb, T. M., & Yoshioka, M. (2015). Extraction of Chemical and Drug Named Entities by Ensemble Learning Using Chemical NER Tools Based on Different Extraction Guidelines. *Trans. MLDM*, *8*(2), 61-76.

[8]. Eltyeb, S., & Salim, N. (2014). Chemical named entities recognition: a review on approaches and applications. *Journal of cheminformatics*, *6*(1), 17.

[9]. Hanisch, D., Fundel, K., Mevissen, H. T., Zimmer, R., & Fluck, J. (2005). ProMiner: rule-based protein and gene entity recognition. *BMC bioinformatics*, *6*(1), S14.

[10]. Alfred, R., Leong, L. C., On, C. K., & Anthony, P. (2014). Malay named entity recognition based on rule-based approach. *International Journal of Machine Learning and Computing*, *4*(3), 300.

[11]. Li, B., Si, X., Lyu, M. R., King, I., & Chang, E. Y. (2011, October). Question identification on twitter. In *Proceedings of the 20th ACM international conference on Information and knowledge management* (pp. 2477-2480). ACM.

[12]. DIEB, M. (2013, October). Ensemble approach to extract chemical named entity by using results of multiple cner systems with different characteristic. In *BioCreative Challenge Evaluation Workshop* (Vol. 2, p. 162).

[13]. Mogotsi, I. C. (2010). Christopher d. manning, prabhakar raghavan, and hinrich schütze: Introduction to information retrieval.

[14]. Sun, J., Gao, J., Zhang, L., Zhou, M., & Huang, C. (2002, August). Chinese named entity identification using class-based language model. In *Proceedings of the 19th international conference on Computational linguistics-Volume 1* (pp. 1-7). Association for Computational Linguistics.

[15]. Quan, C., Wang, M., & Ren, F. (2014). An unsupervised text mining method for relation extraction from biomedical literature. *PloS one*, *9*(7), e102039.

[16]. Sarkar, K. (2015). A hidden markov model based system for entity extraction from social media english text at fire 2015. *arXiv preprint arXiv:1512.03950*.

[17]. Amarappa, S., & Sathyanarayana, S. V. (2013). Named entity recognition and classification in Kannada language. *International Journal of Electronics and Computer Science Engineering*, *2*(1), 281-289.

[18]. Zhao, S. (2004, August). Named entity recognition in biomedical texts using an HMM model. In *Proceedings of the International Joint Workshop on Natural Language Processing in Biomedicine and its Applications* (pp. 84-87). Association for Computational Linguistics.

[19]. Hsu, C. K. W. B. J. P., & Kıcıman, M. W. C. E. Simple and Knowledge-intensive Generative Model for Named Entity Recognition.

[20]. Umare, S. P., & Deshpande, D. N. A. A Survey on Machine Learning Techniques to Extract Chemical Names from Text Documents. *International Journal of Computer Science and Information Technology*, *4*, 1263-1266.

[21]. Mickelin, J. (2013). Named Entity Recognition with Support Vector Machines.

[22]. Ju, Z., Wang, J., & Zhu, F. (2011, May). Named entity recognition from biomedical text using SVM. In *Bioinformatics and Biomedical Engineering,(iCBBE) 2011 5th International Conference on* (pp. 1-4). IEEE.

[23]. Baksa, K., Dolovic, D., Glavaš, G., & Šnajder, J. (2014, January). Named Entity Recognition in Croatian Tweets. In *Ninth Language Technologies Conference, Information Society (IS-JT 2014)*.

[24]. Wang, J. (2015). *A Rule-based Methodology and Feature-based Methodology for Effect Relation Extraction in Chinese Unstructured Text* (Master's thesis, University of Sydney).

[25]. Bui, Q. C., Sloot, P. M., Van Mulligen, E. M., & Kors, J. A. (2014). A novel feature-based approach to extract drug–drug interactions from biomedical text. *Bioinformatics*, *30*(23), 3365-3371.

[26]. Tang, J., & Liu, H. (2014). An unsupervised feature selection framework for social media data. *IEEE Transactions on Knowledge and Data Engineering*, *26*(12), 2914-2927.

[27]. Pons, E., Becker, B. F. H., Akhondi, S. A., Afzal, Z., van Mulligen, E. M., & Kors, J. A. (2015). RELigator: chemical-disease relation extraction using prior knowledge and textual information. In *Proceedings of the Fifth BioCreative Challenge Evaluation Workshop* (pp. 247-253).

[28]. Tongtep, N., & Theeramunkong, T. (2009, April). A feature-based approach for relation extraction from Thai news documents. In *Pacific-Asia Workshop on Intelligence and Security Informatics* (pp. 149-154). Springer, Berlin, Heidelberg.

[29]. Moncecchi, G., Minel, J. L., & Wonsever, D. (2010, November). A survey of kernel methods for relation extraction. In *Workshop on NLP and Web-based Technologies*.

[30]. Zelenko, D., Aone, C., & Richardella, A. (2003). Kernel methods for relation extraction. *Journal of machine learning research*, *3*(Feb), 1083-1106.

[31]. Li, J., Zhang, Z., Li, X., & Chen, H. (2008). Kernel- based learning for biomedical relation extraction. *Journal of the Association for Information Science and Technology*, *59*(5), 756-769.

[32]. Sahu, S. K., Anand, A., Oruganty, K., & Gattu, M. (2016). Relation extraction from clinical texts using domain invariant convolutional neural network. *arXiv preprint arXiv:1606.09370*.

[33]. Lee, J. Y., Dernoncourt, F., & Szolovits, P. (2017). MIT at SemEval-2017 Task 10: Relation Extraction with Convolutional Neural Networks. *arXiv preprint arXiv:1704.01523*.

[34]. Wang, W., Ning, Y., Rangwala, H., & Ramakrishnan, N. (2016, October). A Multiple Instance Learning Framework for

Identifying Key Sentences and Detecting Events. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management* (pp. 509-518). ACM.

[35]. Chen, Y. J., & Hsu, J. Y. J. (2016, March). Chinese Relation Extraction by Multiple Instance Learning. In *AAAI Workshop: Knowledge Extraction from Text*.

[36]. Nattee, C., Sinthupinyo, S., Numaot, M., & Okadatt, T. Multiple-Instance Learning Based Heuristics for Mining Chemical Compound Structure.

[37]. Lamurias, A., Clarke, L. A., & Couto, F. M. (2017). Extracting microRNA-gene relations from biomedical literature using distant supervision. *PloS one*, *12*(3), e0171929.

[38]. Kumar, K., & Manocha, S. (2007). Constructing knowledge graph from unstructured text. *Self*, *3*(4).

[39]. Agichtein, E., & Gravano, L. (2000, June). Snowball: Extracting relations from large plain-text collections. In *Proceedings of the fifth ACM conference on Digital libraries*(pp. 85-94). ACM.

[40]. Zhou, X., Liu, B., Wu, Z., & Feng, Y. (2007). Integrative mining of traditional Chinese medicine literature and MEDLINE for functional gene networks. *Artificial intelligence in medicine*, *41*(2), 87-104.

[41]. Pennacchiotti, M., & Pantel, P. (2006). A bootstrapping algorithm for automatically harvesting semantic relations. In *Proceedings of the Fifth International Workshop on Inference in Computational Semantics (ICoS-5)*.

[42]. Gu, J., Hu, Y. N., Qian, L., & Zhu, Q. (2013). Research on building family networks based on bootstrapping and coreference resolution. In *Natural Language Processing and Chinese Computing* (pp. 200-211). Springer, Berlin, Heidelberg.

[43]. Kim, M. S. W. L. J. (2010). Extraction and Visualization of Implicit Social Relations on Social Networking Services.

[44]. Luo, Y., Uzuner, Ö., & Szolovits, P. (2016). Bridging semantics and syntax with graph algorithms—state-of-the-art of extracting biomedical relations. *Briefings in bioinformatics*, *18*(1), 160-178.

[45]. Peng, Y., Gupta, S., Wu, C., & Shanker, V. (2015). An extended dependency graph for relation extraction in biomedical texts. *Proceedings of BioNLP 15*, 21-30.

[46]. Bakal, G., & Kavuluru, R. (2015, December). Predicting treatment relations with semantic patterns over biomedical knowledge graphs. In *International Conference on Mining Intelligence and Knowledge Exploration* (pp. 586-596). Springer, Cham.

[47]. Postma, G. J., van Bakel, B., & Kateman, G. (1996). Automatic extraction of analytical chemical information. System description, inventory of tasks and problems, and preliminary results. *Journal of chemical information and computer sciences*, *36*(4), 770-785.

[48]. Gentile, A. L., Zhang, Z., Augenstein, I., & Ciravegna, F. (2013, June). Unsupervised wrapper induction using linked data. In *Proceedings of the seventh international conference on Knowledge capture* (pp. 41-48). ACM.

[49]. Mori, J., Ishizuka, M., & Matsuo, Y. (2007, January). Extracting Keyphrases to Represent Relations in Social Networks from Web. In *IJCAI* (Vol. 7, pp. 2820-2827).

[50]. Liu, S., & Ren, F. (2012, October). Relation extraction from wikipedia articles by entities clustering. In *Cloud Computing and Intelligent Systems (CCIS), 2012 IEEE 2nd International Conference on* (Vol. 3, pp. 1491-1495). IEEE.

[51]. Etzioni, O., Fader, A., Christensen, J., Soderland, S., & Mausam, M. (2011, July). Open Information Extraction: The Second Generation. In *IJCAI* (Vol. 11, pp. 3-10).

[52]. Nakashole, N. T. (2012). Automatic extraction of facts, relations, and entities for web-scale knowledge base population.

[53]. Nguyen, N., Miwa, M., Tsuruoka, Y., & Tojo, S. (2013, December). Open information extraction from biomedical literature using predicate-argument structure patterns. In *Proceedings of the 5th International Symposium on Languages in Biology and Medicine* (Vol. 51, p. 55).

[54]. Wu, F., & Weld, D. S. (2010, July). Open information extraction using Wikipedia. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*

[55]. Dingli, A., Ciravegna, F., & Wilks, Y. (2003, October). Automatic semantic annotation using unsupervised.