

## PATTERN-BASED SYSTEM TO DETECT THE ADVERSE DRUG EFFECT SENTENCES IN MEDICAL CASE REPORTS

SAFAA ELTYEB<sup>1</sup>, NAOMIE SALIM<sup>2</sup>

<sup>1,2</sup>Faculty of Computing, Universiti Teknologi Malaysia, Johor, Malaysia

<sup>1</sup>College of Computer Science and Information Technology, Sudan University of Science and Technology, Khartoum, Sudan

E-mail: <sup>1</sup>[safaa-82@hotmail.com](mailto:safaa-82@hotmail.com), <sup>2</sup>[naomie@utm.my](mailto:naomie@utm.my)

### ABSTRACT

#### Background

The detection of adverse drug effect sentences in medical text reduces the efforts required for the manual task of drug safety monitoring by decreasing the number of reports which need to be investigated by drug safety experts. Moreover it helps in compiling a highly accurate and machine-understandable drug-related adverse effects knowledge bases which can support pharmacovigilance research and aids computational approaches for drugs repurposing. In this study, we proposed a pattern-based method to detect the sentences containing drug-adverse effect causal relation from medical case reports.

#### Materials and methods

A set of 500 full abstracts from Medline medical case reports containing 988 adverse drug effect sentences from ADE corpus were used to evaluate the sentences detection task. Our method combined an outcome of a concept recognition system with a method for automatic generation of numerous patterns.

#### Results

Our method achieved recall of 72.8, precision of 93.6 and F-Score of 81.7 % in the adverse drug effect sentences detection task .

#### Conclusion

The results of this study can help database curators in compiling medical databases and researchers to digest the huge amount of textual information which is growing rapidly. Moreover, the mining algorithms developed in this study can be employed to detect sentences contain new associations between other medical entities in medical text.

**Keywords:** *Drug, Adverse effect, pattern-based, Relation extraction, Medical case reports .*

### 1. BACKGROUND

The rapid increase in the flow of published digital information in all disciplines has resulted in a pressing need for techniques that can simplify the use of this information. Manual information extraction (IE) from the literature by humans has become a business managed by information providers. However, manual IE is costly and obtaining the extracted information after publication is often time consuming and fallible [1].

Normally newly discovered biomedical knowledge is presented firstly as results in scientific publications before it is structured into data or knowledge bases[2]. Hence, IE helps database curators in compiling biomedical databases and assists researchers to digest the huge amount of textual information which is growing rapidly.

For instance, case reports published in medical literature have a basic role in the progress of medical science owing to their ample existence, rapid rate of generation and the essential information they include. They contribute to new discoveries of drugs and unpredicted effects. Case reports play an important role in detecting novelty and therefore contribute to medical progress [3].

In medical text, IE has achieved good results in named entity recognition (NER) tasks for entities (i.e. drugs, diseases, adverse effects, etc.). A later step is the extraction of relations between those recognized entities. In this paper, we describe our method for the identification of sentences containing drug-adverse effect causal relation .

An adverse drug reaction (ADR) (sometimes called adverse drug event or adverse drug effect) is defined as any appreciably noxious, unintentional, undesired or unpleasant reaction which results from the use of a dose of a medicinal product for the



purpose of prophylaxis, diagnosis or therapy. It predicts danger from future administration and asserts alteration of the dosage or withdrawal of the product (World Health Organization<sup>1</sup>-WHO; [4]. ADRs lead to further health complications or sometimes even death. The economic effect of ADRs is very substantial. For instance, about \$136 billion is spent annually on treating ADRs in the US, and other countries face similar problems [5-6]. In 1994, ADRs are ranked between the fourth and sixth leading causes of death in the US [7].

Usually the adverse effects profile of a drug is not known at the time of approval owing to the small size of samples, short period and limited applicability of pre-approval clinical trials [8]. Consequently several additional adverse effects appear later after a drug is used by a larger number of people for longer periods. As a result, drug manufacturers should monitor and report the adverse effects to the responsible authorities to decide on a suitable procedure (i.e. modifying drug usage, withdrawal from market, etc.) [9].

The main sources of ADRs information are clinical trials and post-marketing surveillance measures available to some authorities responsible for protecting and promoting public health such as the Food and Drug Administration (FDA<sup>2</sup>) and Centers for Disease Control and Prevention (CDC) in the US and similar governmental agencies in other countries [10]. Examples of other textual sources of ADRs information are patient health records, hospital discharge summaries, medical case reports, full text research articles, blogs and news reports [9].

Automatic IE can help regulatory authorities in rapid information screening and extraction, instead of manual inspection or traditional searching. As a result, this contributes to the acceleration of medical decision support, safety alert generation and risk factor estimation [9]. The extraction of ADRs information also helps in computational strategies for drug repurposing, such as drug side effect similarity strategy [11]. Moreover, relation extraction between drugs and adverse effects helps in indexing, accurate searching, visualization and rapid information tracing, and enhances the sensitivity of signal detection in pharmacovigilance [12].

This paper presents a pattern-based relation extraction system to detect the assertive sentences

contain drug-adverse effects CAUSE relation in Medline case reports abstracts.

The rest of this paper is organized as follows. Section 2 is a related work section that gives an idea about previous work in the extraction of ADRs from text. Details of the relation detection and evaluation methodology are presented in Section 3. Section 4 presents the results obtained, while Section 5 presents an outlook on the proposed method and results obtained. The conclusion and future work are set out in Section 6.

## 2. RELATED WORK

In recent years, many systems have been developed for the automatic extraction of relations between drugs and other specific entities from medical text, such as drug-drug interaction relation [13-15] and drug-disease treatment relation [11] [16]. Comparatively, few studies have addressed the extraction of drug-adverse effects CAUSE relation owing to the lack, until recently, of an open access gold standard corpora for evaluation purposes.

[17] adapted the Cancer Text Information Extraction System (CaTIES) for identifying terms suggestive of adverse drug events in documents. [18] used the natural language processing MedLEE system to identify medication events and entities which could be potential adverse drug events. A co-occurrence approach was used to detect associations between the two types of entities from discharge summaries of seven drugs. [10] proposed a lexical-based framework for mining relationships between drugs and adverse effects from user comments on health-related websites. [19] applied statistics and heuristic methods to build up a hierarchical ontology of side effects from patient-submitted drug reviews on health-related websites, focusing on the statin class of cholesterol-lowering drugs. [20] applied a Hidden Markov Model-based text mining system that can be used to extract the adverse side-effects of drugs from online medical forums. [21] developed a knowledge-based relation extraction system from Medline medical case reports. The knowledge base is a graph representation of concepts and relations between them, populated from the Unified Medical Language System (UMLS). A concept identification module was employed to identify drugs and adverse effects, and the knowledge-based module identifies the existence or non-existence of adverse effect relations between the found entities.

Another work in Medline case reports is a machine learning-based system by [12]. The Java

<sup>1</sup> <http://www.who.int>

<sup>2</sup> <http://www.fda.gov/>



Simple Relation Extraction (JSRE) tool is based on a Support Vector Machine used for identification and extraction of drug-related adverse effects. Closer to our work is [22], which used heuristic rule-based patterns to identify the relations between drugs and adverse effects from clinical records. Our work differs by using a large number of automatic generated patterns that do not depend on specific keywords, while the patterns in [22] were based on a small set of keywords designed manually.

In this study, firstly we automatically generate CAUSE relation-specific textual patterns from a training set of Medline medical case reports sentences that contain at least one drug-adverse effect CAUSE relation using known drug and adverse effects pairs. Then, we detect the potential sentences from another test set of full abstracts from Medline medical case reports.

The main characteristic of our work is that it does not require human effort to build patterns manually for the relation identification process, like all previous pattern-based systems, and consequently a large drug-adverse effects relationship knowledge base can be built from the large generated patterns. Furthermore, our system can discriminate between the type of relations (i.e. CAUSE, TREAT relations) between the drug and medical condition entities mentioned in the same sentence, while this discrimination issue is not explored by most of the previous studies.

### 3. MATERIALS AND METHODS

#### 3.1 Corpus

The data set used for training and testing is the ADE corpus [12]. The ADE corpus contains 2,972 Medline case reports which are manually annotated. The corpus contains annotations of 5,063 drugs, 5,776 conditions (e.g. diseases, signs, symptoms) and 6,821 relations between drugs and conditions comprising drug-adverse effect relations in 4,272 sentences taken from 1,644 abstracts. Drugs and conditions that do not comprise a potential adverse event relation are not annotated. Each relation is represented by a Medline identifier; the sentence contains a relation, the drug and its position, and the adverse effect and its position with regard to the Medline abstract

Owing to the sensitivity of our method in pattern generation, all names of drugs and conditions should be removed before generating patterns. Subsequently, all sentences containing drugs and conditions entities which obtained a suitable mapping from UML such as mapping to the broad UMLS semantic type finding (i.e.

*myotonia*, *hair loss*, *neutropenia*, etc.), or some entities not covered by UMLS (i.e. *pruritic bullous eruption*, *decrease in the D-dimers*, etc.) or are not mapped by UMLS because of differences in spelling between the UMLS metathesaurus and the corpus sentence (i.e. *interferon alpha* versus *interferon alfa*) are discarded. After the corpus is cleaned, there are 3,180 drug-adverse effect relations in 2,362 sentences.

For training and testing purposes, a 500 abstracts were selected randomly, using 10 cross validation. The next table shows the numbers of “positive sentences” and “negative sentences” in the selected abstracts

Table 1: Number Of Positive Sentences And Negative Sentences In The Selected Abstracts Of ADE Corpus

Items	Number
Sentences with at least one drug-adverse effect relation “positive sentences”.	988
Sentences with no relation “negative sentences”.	2196

#### 3.2 Adverse-drug effect Relationship Extraction.

Relation extraction is an important task in IE which aims to discover and describe the semantic relations between entities in text. Relations between entities can be inferred, after the entities are identified. Usually, the relations are binary. If it includes more than two entities, the relation said to contain complex associations. Relation extraction approaches range from applying the simple co-occurrences search to a syntactic analysis and parsing.

Case-Based Reasoning (CBR) is a methodology that depends on finding a past case similar to the new one for solving problems [23]. In IE systems, the cases from the training data (here, patterns for drug-adverse effects relation in sentences) are learnt in the training step, and then saved in a case base. During the testing step, the system searches the case base for cases most similar to the case problem.

The main idea for achieving the drugs- adverse effects sentence detection is that a set of automatic generated patterns are utilized to identify the existence or non-existence of a drug- adverse effect relation in a sentence. Constructing keywords or manual patterns to retrieve the most relevant case of drug- medical condition relation to specific sentence has main two disadvantages. Firstly, choosing keywords and constructing patterns requires domain experts and is often performed



manually. Secondly, the built patterns and keywords do not guarantee complete differentiation among cases. Furthermore, if the domain knowledge is built up, modifying these patterns can become very difficult.

In order to avoid these disadvantages, the Minimal Differentiator Expressions (MDE) algorithm proposed in [24] is adapted in this work. This algorithm automatically generates a set of linguistic patterns (expressions) used to retrieve the case most suitable to the input sentence. The main characteristic of those patterns is that they are composed of the simplest sets of words which permit differentiation among cases (consequently among the sentences contain another type of drug-condition relation - e.g TREAT relation-or sentences don't contain a relation between drugs and medical conditions entities). In this work, the automatically generated expressions are used to identify the assertive sentences containing a drug-adverse effect CAUSE relation and distinguish them from other sentences which do not contain a drug- adverse effect relation or contain another type of drug-condition relation. Examples of the generated expressions are: {*\*developed\* after\*starting\**}, {*\*caused\*by\**}, {*\*presented\*with\* after\**}, {*\*treated\*with\* developed\**}, {*\*presented\*after\* of\**}, {*\*occurring\*after\**}, {*\*developed\*after\*receiving\**}, {*\*induced\*after\**}, {*\*developed\*during\*treatment\**}, etc., where the character '\*' denotes to zero or more word, [see Additional file1.doc] for more examples of generated MDEs and the sentences detected through them.

If an expression unambiguously distinguishes a sentence from other sentences in the case base not necessarily on the sentence in its own case, it is called a Differentiator Expression (DE) of the sentence. MDE is a differentiator expression which does not contain any other differentiator expressions [24]. For instance, if we have the two DEs {*\*may\* cause\**} and {*\* cause\**}, the expression {*\*may\* cause\**} is not minimal because it contains the MDE {*\* cause\**} which has fewer terms and it also allows differentiation. The advantage of DEs is the prevention of linguistic interferences and overlapping between templates or expressions which is difficult to detect manually. All sentences contains drug -condition pairs induce a relation which corresponds to one case. The sentences in a case are represented by MDEs which are composed from the token(s) of each sentence that is mentioned between the first drug or condition and the last drug or condition in the sentence after removing any drug or condition

token (s) existing in the middle. The sentences used in training to generate the MDEs are preprocessed by removing some punctuation marks and numbers. Both sentences used for training and testing are converted to lower case.

MetaMap<sup>3</sup> [25], a Java API from the National Library of Medicine, is used to map the biomedical text to concepts in the UMLS metathesaurus. Currently, UMLS has 135 semantic types which are grouped into 15 semantic groups. To identify drugs and medical conditions in sentences, the 'Chemicals and Drugs' and 'Disorders' semantic type groups are used.

In brief, the modules of the system are divided into two parts. The first part aims to fill the case base with specific cases for the training phase and represented on: annotating sentences with MetaMap API; preprocessing the annotated sentences by removing the annotated entities and short sentences by taking specific token(s) to compose the MDEs as mentioned above; get MDEs for each sentence and calculate MDEs' CRelevance according to Eq. 1. The second part of the modules implemented the testing phase. First, each sentence is assigned to the most relevant case. The assignment tested according many criteria such as sentence match to specific number of generated MDEs or MDE relevance to specific case according to Eq. 1. CRelevance is calculated as the proportion of sentences S that satisfy expression e, according to the following equation:

$$CRelevance(e, C) = \frac{|\{S \in C \mid e \text{ exp } S \text{ wrt } B\}|}{|\{S \in C\}|} \quad (1) \quad [24]$$

Where B is the case base of cases C. Consequently, less important MDEs will not lead to a classification process[24].

#### 4. RESULTS

As mentioned above, a positive sentence is one that contains at least one clear definition of a drug causing a medical condition. Negative sentence does not contain. The sentences detection task performance was measured using the IE evaluation metrics precision (P), recall (R), and F-measure (F) over positive labeled sentences since they represented the focused class of sentences being studied. We used 500 abstracts selected randomly from the ADE corpus. While only the positive

<sup>3</sup> <http://metamap.nlm.nih.gov/>

sentences were used to generate the patterns in the training phase, the test set contains all positive and negative sentences of abstracts. We ran 10-fold cross validation to evaluate the performance of the method. At each run, 90% of the abstracts were used for the training and 10% for the testing.

Firstly, in the baseline experiment, we set the criteria of match between sentences and expressions (patterns) as follows: if a sentence from the testing set matches at least five patterns it is classified as positive sentence. To reduce false positive (FP) classified sentences and increase precision, we calculated the CRelevance of all generated MDEs according to Eq. 1. The MDE with the highest CRelevance for each sentence in the training set was saved. Then, if a sentence matched at least one MDE, it was classified as a positive sentence. This attempt resulted in improvement in P but degradation in R due to the absence of some important MDEs. In another attempt to improve R, we decreased the MDE number that should be at least matched by a sentence to one. We observed an increase in R and a sharp fall in P due to the increase in FP sentences. Further improvement to reduce the FP classified sentences was done by manual curation of all generated MDEs, which improved the F-score by 34.9 percentage points compared with the baseline experiment as, shown in table 2.

*Table 2: Performance Evaluation (In %) Of The Sentence Detection Task With Different Adjustments Of Sentence Match With Patterns Evaluated By 10-Fold Cross-Validation.*

Adjustments of sentence match	P	R	F
Sentence matches at least five MDEs from all generated MDEs	43.9	51.1	46.8
Sentence matches at least one MDE after selecting MDEs with the highest CRelevance	58.8	42.2	48.8
Sentence matches at least one MDE from all generated MDEs	37.2	74.9	49.5
Sentence matches at least one MDE after manual curation for all generated MDEs.	93.6	72.8	<b>81.7</b>

## 5. DISCUSSION

We have investigated the use of automatic generated patterns and case-based reasoning methodology to detect assertive sentences containing a drugs-adverse effects causal relation.

The main advantages of the presented method are achieving higher precision than other classification approaches such as machine learning and the automation of the patterns generation process which is characterized by less effort required to build patterns. At the same time, the results obtained are highly dependent on the richness of the training data and to what extent the test data include sentences that match the generated patterns. Therefore, there is a need for a large corpus of training data to significantly increase performance, and this is considered the main disadvantage of the presented method.

An error analysis was carried out on a sample of 100 randomly selected errors that were made by the sentence detection module. The largest source of errors (35 of FPs and 65 of FNs) was the FNs (65%) which occurred in our system when the target sentence has no matching pattern(s) among the available patterns. This problem might be relieved by increasing the size of training set and consequently increasing the generated patterns. FP classified sentences were generated because some sentences were matched with unsuitable MDE(s), and consequently the system classified the negative sentences as positive ones. However, the training corpus exhibits a higher class imbalance towards negative and positive sentences.

When we compared our sentence identification system with the [9] machine learning system, our system performed better by 4.7 percentage in term of F score. However, comparison of our system with the many other relation extraction and sentence classification systems reported in the literature is difficult, because of the large numbers of relation extraction tasks and evaluation datasets.

## 6. CONCLUSION AND FUTURE WORK

We developed a pattern-based medical relationship detection method. The system's value in the real world is represented by the ability of the sentence detection module to detect new cases of ADE or modify the existing statistics of ADEs. The automated detection of ADEs sentences reduces the efforts required in the manual task of drug safety monitoring by decreasing the number of final reports that need to be investigated by drug safety



experts [9] and assist drug safety professionals in collecting information from free text.

The main limitation of the present study is the dependence of the current system on Metamap mapping to the UMLS metathesaurus concepts for drugs and the identification of medical conditions. UMLS does not provide coverage for some phrases (*i.e pruritic bullous eruption, decrease in the D-dimers*) and descriptive phrases (*i.e behaving oddly, started sweating profusely, shivering, etc.*). Also, there were many entities that belonged to the broad UMLS semantic type group “*finding*” which were discarded and consequently the used corpus size was reduced. However, designing clean dictionaries for PubMed abstracts and incorporating other resources may remedy this drawback.

To further improve the study outcome, future studies should increase the training data by exploiting additional medical lexicons and resources for full coverage of entities recognition to incorporate more patterns, and machine-learning techniques to use these patterns in real case studies. Since the patterns generated by our current work depend on the underlying UMLS-based lexicons, we only used the patterns “*DRUG \*pattern\* ADVERSE EFFECT*” and “*ADVERSE EFFECT \*pattern\* DRUG*”. In our future studies we plan to use syntactic information with patterns such as “*NP1 \*pattern\* NP2*” and “*NP2 \*pattern\* NP1*” where NP1 and NP2 are noun phrases representing complete or parts of drugs and adverse effect term(s) to increase the generated patterns and consequently increase recall.

#### AUTHORS PROFILES:

SE: B. Comp. Sc. (SUST, Sudan), M. Sc. Comp. Sc. (SUST, Sudan), currently

Ph. D student at (UTM, Malaysia)

NS: Professor Dr. B. Comp. Sc. (UTM, Malaysia), M. Sc. Comp. Sc. (W. Michigan, US)

Ph. D Info. Sc. (Univ. of Sheffield, UK).

#### ACKNOWLEDGEMENTS

This work is supported by the Ministry of Higher Education (MOHE) and

Research Management Centre (RMC) at the Universiti Teknologi Malaysia (UTM) under Research University Grant Category (VOT: Q.J130000.2528.07H89) .We also would like to thank Sudan University of Science and Technology (SUST) for sponsoring the first author.

#### REFERENCES

- [1] L. Hawizy, *et al.*, "ChemicalTagger: A tool for semantic text-mining in chemistry," *Journal of cheminformatics*, vol. 3, p. 17, 2011.
- [2] P. Thomas, *et al.*, "Relation extraction for drug-drug interactions using ensemble learning," *1st Challenge task on Drug-Drug Interaction Extraction (DDIExtraction 2011)*, pp. 11-18, 2011.
- [3] J. P. Vandembroucke, "In defense of case reports and case series," *Annals of Internal Medicine*, vol. 134, pp. 330-334, 2001.
- [4] I. R. Edwards and J. K. Aronson, "Adverse drug reactions: definitions, diagnosis, and management," *The Lancet*, vol. 356, pp. 1255-1259, 2000.
- [5] R. Leone, *et al.*, "Drug-Related Deaths," *Drug Safety*, vol. 31, pp. 703-713, 2008.
- [6] C. S. van Der Hooft, *et al.*, "Adverse drug reaction-related hospitalisations," *Drug Safety*, vol. 29, pp. 161-168, 2006.
- [7] J. Lazarou, *et al.*, "Incidence of adverse drug reactions in hospitalized patients: a meta-analysis of prospective studies," *Jama*, vol. 279, pp. 1200-1205, 1998.
- [8] S. R. Ahmad, "Adverse drug event monitoring at the Food and Drug Administration," *Journal of general internal medicine*, vol. 18, pp. 57-60, 2003.
- [9] H. Gurulingappa, *et al.*, "Identification of adverse drug event assertive sentences in medical case reports," in *First international workshop on knowledge discovery and health care management (KD-HCM), European conference on machine learning and principles and practice of knowledge discovery in databases (ECML PKDD)*, 2011, pp. 16-27.
- [10] R. Leaman, *et al.*, "Towards internet-age pharmacovigilance: extracting adverse drug reactions from user posts to health-related social networks," in *Proceedings of the 2010 workshop on biomedical natural language processing*, 2010, pp. 117-125.
- [11] R. Xu and Q. Wang, "Large-scale extraction of accurate drug-disease treatment pairs from biomedical literature for drug repurposing," *BMC bioinformatics*, vol. 14, p. 181, 2013.
- [12] Gurulingappa, *et al.*, "Extraction of Adverse Drug Effects from Medical Case



- Reports," *Journal of Biomedical Semantics*, vol. Volume 3, 2012.
- [13] I. Segura-Bedmar, *et al.*, "Extracting drug-drug interactions from biomedical texts," *BMC bioinformatics*, vol. 11, p. P9, 2010.
- [14] L. Tari, *et al.*, "Discovering drug-drug interactions: a text-mining and reasoning approach based on properties of drug metabolism," *Bioinformatics*, vol. 26, pp. i547-i553, 2010.
- [15] S. Karnik, *et al.*, "Extraction of drug-drug interactions using all paths graph kernel," *1st Challenge task on Drug Drug Interaction Extraction*, 2011.
- [16] A. Behir and W. Ben Abdesslem Karaa, "Extraction of drug-disease relations from MEDLINE abstracts," in *Computer and Information Technology (WCCIT), 2013 World Congress on*, 2013, pp. 1-3.
- [17] M. Gysbers, *et al.*, "Natural language processing to identify adverse drug events," in *AMIA... Annual Symposium proceedings/AMIA Symposium. AMIA Symposium*, 2007, pp. 961-961.
- [18] X. Wang, *et al.*, "Active computerized pharmacovigilance using natural language processing, statistics, and electronic health records: a feasibility study," *Journal of the American Medical Informatics Association*, vol. 16, pp. 328-337, 2009.
- [19] J. Liu, *et al.*, "Automatic drug side effect discovery from online patient-submitted reviews: Focus on statin drugs," in *IMMM 2011, The First International Conference on Advances in Information Mining and Management*, 2011, pp. 91-96.
- [20] H. Sampathkumar, *et al.*, "Mining Adverse Drug Side-Effects from Online Medical Forums," in *Healthcare Informatics, Imaging and Systems Biology (HISB), 2012 IEEE Second International Conference on*, 2012, pp. 150-150.
- [21] N. Kang, *et al.*, "Knowledge-based extraction of adverse drug events from biomedical text," *BMC bioinformatics*, vol. 15, p. 64, 2014.
- [22] E. Aramaki, *et al.*, "Extraction of adverse drug effects from clinical records," *Stud Health Technol Inform*, vol. 160, pp. 739-43, 2010.
- [23] A. Aamodt and E. Plaza, "Case-based reasoning: Foundational issues, methodological variations, and system approaches," *Ai Communications*, vol. 7, pp. 39-59, 1994.[24] A. Moreo, *et al.*, "A high-performance FAQ retrieval method using minimal differentiator expressions," *Knowledge-Based Systems*, vol. 36, pp. 9-20, 2012.
- [25] A. R. Aronson, "Effective mapping of biomedical text to the UMLS Metathesaurus: the MetaMap program," in *Proceedings of the AMIA Symposium*, 2001, p. 17.