

VisualRDR: A general framework for creating, maintaining and learning of ripple down rules for Information Extraction

Delip Rao Sachindra Joshi Ganesh Ramakrishnan Sreeram Balakrishnan
Ashwin Srinivasan

Department of Computer Science and Engineering
Indian Institute of Technology Madras
Chennai-600036
India
email: delip@cse.iitm.ernet.in

Abstract

The problem facing the industry and the common user today is that of an information glut. Large amounts of useful information, in various forms, are being generated mostly for human consumption. This deluge of information requires us to find ways of gathering information from such unstructured sources with as little manual intervention as possible. This has been made possible by information extraction. In this paper we propose a novel knowledge-based approach to information extraction. We demonstrate how this knowledge for information extraction can be accrued in an incremental fashion.

1 Introduction

Information comes in various forms from various sources including, but not limited to, product manuals, FAQs, man-pages, websites, weblogs, email and news articles. The information thus generated is dynamic. Further to complicate matters, the information in most cases is generated for manual viewing where style and language occlude automatic perception of information. At the time of this writing, Google [1] alone indexes more than eight billion pages. Much more information than this exists in offline form. Manually perusing through this sea of unstructured sources for extracting information is an infeasible task. This motivates the need for techniques that automatically extract relevant information and present them in a structured format to enable querying or further processing. Information extraction is the task of identifying and

extracting relevant information from an unstructured source. The information extraction (IE) task can be formally defined as follows [18]:

“Given a set of structured elements E (target schema) and an unstructured source S , extract all instances of E from S .”

In general, there are two approaches to building information extraction systems: the *knowledge engineering approach* and the *trainable system approach* [3]. The knowledge engineering approach refers to manual development of rules (or grammars) for IE by “Knowledge Engineers”. The advantage of such a system is that, with skill and experience, it is possible to achieve good performance. Also the resulting knowledge base is very compact and comprehensible. But the flip side is that the manual rule development process, although easy, is very laborious.

The strengths and weaknesses of the automatic training approach are contrary to those of the knowledge engineering approach. The focus here is on producing training data instead of producing rules. This training data is used to learn a model that could be used to process novel instances. With the help of domain-experts to annotate texts, systems can be customized to a specific domain without any intervention from the developers. The success of this method depends on the availability of good training data. Annotated training data may be sparse, or difficult and expensive to obtain. Despite this, contemporary work focuses on the automatic learning aiming to alleviate the aforementioned problems.

Recent study by Ireson *et al* [15] shows that the adoption of a machine learning algorithm, in itself, does not provide a guaranteed advantage in information extraction. Hence we aim to bring the best of both worlds together by combining the trainable

system approach and the knowledge engineering approach. We propose a framework for acquiring rules for information extraction which can later be refined by a knowledge engineer to “tune” the system for high performance. This is based on Compton *et al*’s seminal work [6] on a knowledge acquisition methodology called the Ripple Down Rules (RDR). Our inspiration to use ripple down rules comes from its rich structure due to exceptions. Exceptions allow addition of knowledge with out breaking the previous rules.

Section 2 describes the prior art in Information Extraction systems. Section 3 introduces the knowledge acquisition framework for information extraction. Although the focus of this paper is on using ripple down rules for information extraction, in Section 4, we also give method of automatically learning such rules from raw corpus. The knowledge acquisition framework was implemented as a visual tool, called VisualRDR. Section 5 describes VisualRDR. The experiments and the results of using RDR for information extraction are detailed in Section 6.

2 Information Extraction

In this section we provide a synoptic view of some of the existing work on information extraction. Information extraction approaches for the purpose of this study is roughly classified into pattern based approaches and statistical approaches.

2.1 Pattern based approaches

The basic idea behind pattern based approaches is to learn patterns that extract the relevant information. The information extraction process is treated as a process of slot-filling. The patterns (or rules) can be used either for single slot extraction or multi slot extraction. In multi slot extraction the extraction rules are able to link together related information, as opposed to single slot rules that can only extract isolated data (e.g. in a document that contains several names and addresses, single-slot rules can not specify what the address of a particular person is). Section 3 lists some examples of such systems.

2.2 Statistical approaches

Machine learning methods using Naive Bayes, Support Vector Machines (SVM), Hidden Markov Models (HMM), Conditional Random Fields (CRF) and Semi Markov CRFs have been applied to information extraction. In the Naive Bayes method [9] we treat the document as a bag of words and totally ignore the linguistic structure of the document and a priori probabilities are estimated using weights taken as the TF-IDF measure of the words. Experiments have been done to study the possibility of using support vector machines in information integration. Features from

the document/sentence are extracted and classified using an SVM classifier. Hidden Markov Models are a powerful probabilistic tool for time series data and have been successfully applied to problems in speech and text processing. Recently, McCallum *et al* [14] have used HMMs to extract information from natural language text. They alleviate the large training data requirement by HMMs using a statistical technique called shrinkage. Peng *et al* [16] have used conditional random fields to extract information from scientific papers and Sarawagi *et al* [19] demonstrate the possibility of using semi Markov conditional random fields.

3 Knowledge acquisition for Information Extraction

Knowledge acquisition is “the transfer and transformation of potential problem-solving expertise from some knowledge source to a program.” [4] Knowledge acquisition is performed by knowledge engineering techniques. Knowledge engineering is “the process of reducing a large body of knowledge to a precise set of facts and rules.” [8]

The information extraction task can be solved using a knowledge based approach. Explicit rules could be elucidated for extracting information from unstructured sources. Several efforts in the past have made use of rules for information extraction. Some examples include, AUTOSLOG [17], LIEP [12] and CRYSTAL [22], WHISK [21], RAPIER [5] and SRV [10]. All these systems suffer from the maintainability problems pointed by [7].

We use the knowledge acquisition methodology proposed by Compton *et al* [6] called Ripple Down Rules (RDR). The main advantages of using RDR is that

1. RDR enable incremental acquisition of rules.
2. Rules are acquired in context
3. Rule based systems modeled on RDR are amenable to easy maintainance.

3.1 Ripple Down Rules

Ripple down rules (RDR) is a knowledge acquisition methodology and a way of structuring knowledge bases which grew out of long term experience of maintaining an expert system , GARVIN-ESI by Compton *et al* [6]. In the RDR framework, the human expert’s knowledge is acquired based on the current context and is added incrementally. Ripple Down Rules consist of rules which form a tree structure. Many rules in RDR are exceptions to other rules. Ripple down rules methodology allows incremental changes to knowledge base without causing unwanted side effects to the existing knowledge base. Compton *et al* have shown [7] that this approach allows clear separation of the knowledge engineering and domain expertise. This enables domain experts to change the knowledge base without

the need of a knowledge engineer. The root of the tree provides a default conclusion if the rule linked to the child is not satisfied for a particular case. Whenever an RDR incorrectly classifies a case or fails to classify a case, a rule is added. When a rule is added to the RDR structure, the case that prompted the rule is also stored with the rule. These cases are called “cornerstone cases.” Apart from providing contextual information, later we shall see that the cornerstone cases are also useful in on the fly validation. Inferencing in RDR is very simple. Whenever a case is satisfied by the rule which does not have a dependent, i.e. no branches, the conclusion associated with that rule is asserted. On the other hand, if the rule has a TRUE branch then that branch condition is also tested in a depth-first manner. The conclusion of the most deepest satisfying node is returned as the result. If no rule satisfies the case, except the default rule, then an exception branch is added to the default rule and a new rule is created. As always, with every new rule we also store the cornerstone case that actuated the creation of that new rule.

3.2 RDR for information extraction

Simple regular expression based patterns based on features in the text could be written for extracting information. We can organize the rules in the ripple down representation with the more generic rules at the top having child rules that further specialize them. We represent the rules for information extraction in XML format. The schema for our representation is shown below.

```
<!DOCTYPE KNOWLEDGE [
  <!ELEMENT KNOWLEDGE (RULE+)>
  <!ELEMENT RULE (CONTEXT, CONCLUSION,
    CSCASE+, EXCEPTION)>
  <!ELEMENT CONTEXT (#PCDATA)>
  <!ELEMENT CONCLUSION (#PCDATA)>
  <!ELEMENT CSCASE (LOCATION)>
  <!ELEMENT EXCEPTION (RULE+)>
  <!ELEMENT LOCATION (#PCDATA)>

  <!ATTLIST RULE ID CDATA #REQUIRED>
  <!ATTLIST CSCASE ID CDATA #REQUIRED>
]>
```

3.3 On the fly validation

It was realized that the error rate could be eliminated by validating the rules as they were added [6]. Every time a new rule is added, the cornerstone cases could be used to verify if the new rule did not inadvertently break previously existing rules. Thus each time a new rule is added, we verify if the cornerstone cases still continue to be satisfied by the rules which were originally created to cater to them. Instances where this does not happen is flagged to the user as potential

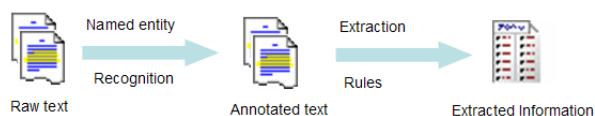


Figure 1: **Information Extraction Process**

problems. The user can then choose to modify the newly added rule or refine the existing rules. This ensures that the knowledge base is consistent at all times.

3.4 Our approach to information extraction

We identify and distinctly separate the tasks of entity identification and information extraction. The raw corpus is run through a rule based named-entity (NE) recognizer (Ref?). The NE annotator identifies entities like NAME, PERSON, DATE, AMOUNT, ORGANIZATION etc.

Ganesh feel free to add more here about the NE annotator.

We then write regular expression patterns (rules) over the annotated document. For example, the pattern `ORGANIZATION.*sold.*ORGANIZATION` would imply the organization on the left sold something to the organization on the right. This process is outlined in Figure 1.

4 Learning of RDR

It is clear that Ripple Down Rules is an efficient way to organize rule based systems. Creation of these rules however could be a time-consuming task, especially when we have large amounts of data and we need to write rules that cover them. This is true of information extraction systems which have to deal with large corpora. We propose to glean these rules from the corpus and construct an RDR tree from it. Several efforts in the past, including, INDUCT [11], Cut95 [20] have addressed the problem of learning RDR from training data.

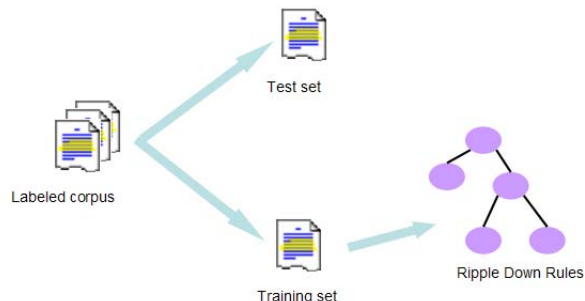


Figure 2: **Learning of Ripple Down Rules**

Figure 2 illustrates our methodology. In order to learn we start with an annotated corpus. We split the annotated corpus into a test and training data set. RDR is induced from the training data set and test set is used for validation of the learnt RDR. The RDR is said to cover a case if it accepts the case.

The steps in learning the RDR are as follows:

1. The corpus is tagged and entity annotated.
2. Extract features between interested entities.
3. These features are clustered on some criterion such that all similar patterns fall in the same cluster.
4. Find a regular expression for each cluster and build a first level RDR using the regular expressions corresponding to each cluster.
5. Now pass the training data again through this first level RDR.
6. At each node where misclassifications happen, take all the misclassified instances at that node and repeat the above algorithm (step 3 onwards) till there is no appreciable reduction in performance is observed.

5 VisualRDR

We implement information extraction in a visual framework, called VisualRDR, for easy creation, maintenance and learning Ripple Down Rules. The tool has to be general enough to support any application that uses RDR. To enable this a plugin based architecture is designed. Where the actual task to be performed by the rules is factored out into externally loadable plugins. Users of the tool can extend it to applications other than information extraction by writing appropriate plugins. Further, we provide facilities to learn RDRs if the application wishes to do so. Since ripple down rules are highly interactive, we envision knowledge engineers to manually add, modify and delete rules. Any modifications to the rule base should be reflected immediately in the GUI and the tool should be able to point out if any rules broke-down due to changes made. The following features were identified to be core for any RDR based application.

1. Create a new RDR
2. Edit an RDR
3. Use an RDR against a set of cases
4. On the fly validation

Optionally, we could also add capability to automatically learn the rules. Again, the exact method of rule learning is left to the user.

5.1 Creating new rules

New rules are added to the ripple down rule structure as exception. Every new rule added has to be an exception to a previously existing rule, which the newly added rule specializes. In case there is no rule relevant to the current case then the new rule is added as an exception to the default rule, which is the root of the ripple down tree. Figure 3 shows the VisualRDR in-

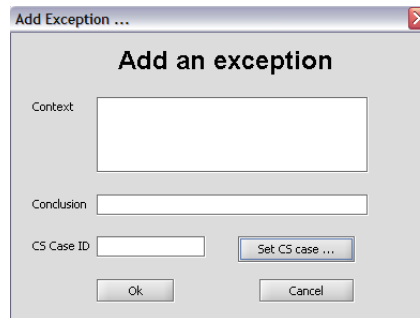


Figure 3: Adding exception

terface for adding rules.

5.2 Editing an RDR

VisualRDR supports manual editing of ripple down rules as any knowledge base system will inevitably involve a human in the loop. The interfaces are kept simple so that any domain-expert can easily modify, add or delete the rules without the need of a knowledge engineer. We ensure the consistency of the knowledge base is not perturbed during each edit operation by providing appropriate visual cues using on the fly validation.

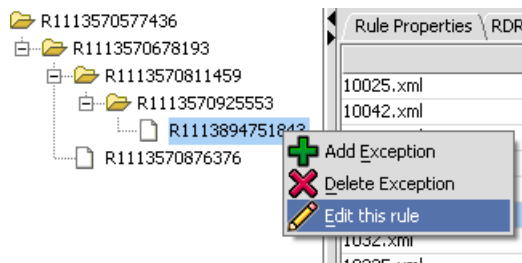


Figure 4: Editing of rules in VisualRDR

5.3 Using the RDR

VisualRDR operates in two modes, batch and single. In the batch mode, an entire collection of cases could be classified using the RDR. On the other hand, the single mode lets use validate a single case against the RDR. Figure 5 shows the VisualRDR interface for applying a rule to a case. The exact operation of the rule is left to the plugin writer. In our case, we extract information.

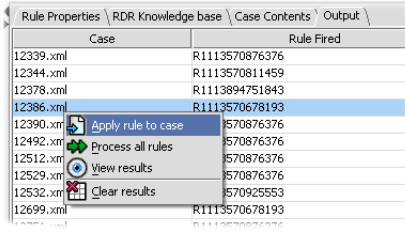


Figure 5: Using RDR

5.4 On the fly validation

Every time a change is made to the ripple down rule structure, the rule base is validated as mentioned in Section 3.3. *OTFValidate*, the procedure for on-the-fly validation is given below. Each time a rule is added to the RDR knowledge base, all cornerstone cases are reclassified using the new RDR. Cases where the new classification is different from the previous classification are flagged as possible errors.

procedure *OTFValidate*(*rdr*)

- 1: $cstones \leftarrow rdr.cscases$
- 2: **for all** $case \in cstones$ **do**
- 3: $newrule \leftarrow Classify(case, rdr)$
- 4: **if** $newrule \neq case.rule$ **then**
- 5: Flag $case$ as possibly missclassified.
- 6: **end if**
- 7: **end for**

Failed CS case	Expected rule	Triggered Rule
14582-R111357067819...	R1113570678193	R1113570577436
1032-R1113570811459...	R1113570811459	R1113570577436
10025-R111357087637...	R1113570876376	R1113570577436
12890-R111357092555...	R1113570925553	R1113570577436
12937-R111389475184...	R1113894751843	R1113570577436

Figure 6: On the fly validation

Figure 6 shows on the fly validation being done in VisualRDR.

5.5 Learning

VisualRDR leaves learning of RDR to be implemented by the plugin writer. In our case, we implement the learning algorithm mentioned in Section 4.

6 Experiments and results

6.1 The dataset

For the purpose of this experiment we use the acquisitions corpus [13] of Reuters news articles on mergers and acquisitions provided by RISE [2]. The corpus has 600 news articles about acquisitions providing information about buyers, sellers, the money involved and so on. The corpus is manually annotated so that it can be used for training and validation purposes.

6.2 The task

Our aim was to identify the names of the seller and purchaser involved from every news snippet in the acquisitions corpus. Although we only extract seller and purchaser instead of all information extracted by comparing techniques [10] [14], extracting other information is merely a process of adding more rules to the knowledge base. In this work we choose demonstrate the concept by extracting seller and purchaser information.

6.3 The knowledgebase

We performed our experiments with 6 seed rules that were refined over a period of time. These rules were based on simple regular expressions on data annotated by entity identification. All rules were hand-crafted. For details, please refer to the appendix A.

6.4 Results

The RDR mentioned listed in appendix A was used to extract information from the acquisitions corpus after entity identification. The performance for extracting seller and purchaser information was as shown below. The number of cases where the seller or buyer information was present is shown as “Total cases”. By precision we mean the proportion of the extracted information (seller or buyer) that are relevant and recall is the proportion of the relevant information that was extracted. The F1 measure combines the two by taking a harmonic mean of precision and recall.

	Seller	Purchaser
Total cases	471	447
Precision	64.54	58.38
Recall	77.70	80.98
F1 measure	70.51	67.84

Table 1: Result of RDR listed in appendix A

We compared our technique with two other methods, SRV [10] and another based on HMMs [14], for extraction of buyer and seller from the acquisitions corpus. Table 2 shows the comparison.

	Seller	Purchaser
SRV	34.30	42.90
HMM	30.90	42.90
RDR	70.51	67.84

Table 2: Comparison of RDR with other techniques

7 Conclusions and future work

In our current work we have shown the utility of a knowledge acquisition methodology, called Ripple Down Rules, for information extraction. We also described VisualRDR, a framework for creating, maintaining and learning of Ripple Down Rules. These rules could be fine-tuned by a domain-expert for achieving high levels of performance. We wish to add additional features to VisualRDR like highlighting redundant and contradictory rules.

8 Acknowledgements

Delip Rao would like to acknowledge the support by IBM India research labs during which this work was done.

References

- [1] <http://www.google.com>.
- [2] <http://www.isi.edu/info-agents/rise/>.
- [3] D. E. Appelt and D. J. Israel. Introduction to information extraction technology. *IJCAI*, 1999.
- [4] D. B. Buchanan, J. B. R. Betchel, and W. Clancey. Constructing an expert system. building expert systems, 1983.
- [5] M. E. Califf and R. J. Mooney. Relational learning of pattern-match rules for information extraction. In *AAAI/IAAI*, pages 328–334, 1999.
- [6] P. Compton and R. Jansen. Knowledge in context: a strategy for expert system maintenance. In *Proceedings of the 2nd Australian Joint Artificial Intelligence conference*, volume 406 of *Lecture Notes in Artificial Intelligence*, pages 292–306, Adelaide, 1988. Springer.
- [7] P. Compton and R. Jansen. A philosophical basis for knowledge acquisition, Aug. 16 2000.
- [8] E. A. Feigenbaum. Knowledge engineering. the applied side of artificial intelligence, 1984.
- [9] D. Freitag. Using grammatical inference to improve precision in information extraction. In *Workshop on Grammatical Inference, Automata Induction, and Language Acquisition (ICML'97)*, Nashville, TN, 1997.
- [10] D. Freitag. Multistrategy learning for information extraction, 1998.

- [11] B. R. Gaines and P. Compton. Induction of ripple-down rules applied to modeling large databases. *J. Intell. Inf. Syst*, 5(3):211–228, 1995.
- [12] S. B. Huffman. Learning information extraction patterns from examples. In *Learning for Natural Language Processing*, pages 246–260, 1995.
- [13] D. D. Lewis. Representation and learning in information retrieval. *Ph.D. Thesis*, 1992.
- [14] A. McCallum and D. Freitag. Information extraction with HMMs and shrinkage. In *AAAI'99 Workshop on Machine Learning for Information Extraction*, 1999.
- [15] F. C. Neil Ireson, D. F. Mary Elaine Califf, N. Kushmerick, and A. Lavelli. Evaluating machine learning for information extraction. In *International Conference on Machine Learning (ICML)*, 2005.
- [16] F. Peng and A. McCallum. Accurate information extraction from research papers using conditional random fields. In *HLT-NAACL*, pages 329–336, 2004.
- [17] E. Riloff. Automatically constructing a dictionary for information extraction tasks. In *AAAI*, pages 811–816, 1993.
- [18] S. Sarawagi. Automation in information extraction and integration. *VLDB*, 2002.
- [19] S. Sarawagi and W. W. Cohen. Semi-markov conditional random fields for, June 07 2004.
- [20] T. Scheffer. Algebraic foundation and improved methods of induction of ripple down rules, Dec. 05 1996.
- [21] S. Soderland. Learning information extraction rules for semi-structured and free text. *Machine Learning*, 34:233, 1999.
- [22] S. Soderland, D. Fisher, J. Aseltine, and Lenhart. CRYSTAL: Inducing a conceptual dictionary. *Proc. Int. Joint Conf. Artificial Intelligence (IJCAI-95)*, pages 1314–1319, 1995.

A The knowledge base

A listing of the RDR knowledge base used for extraction follows.

```

<?xml version="1.0" encoding="iso-8859-1"?>
<knowledge>
  <rule id="R1113570577436">
    <context>ANY</context>
    <conclusion>NOTHING</conclusion>
    <cscase id="CO">
      <location>NONE</location>
    </cscase>
  </rule>
  <exception>
    <rule id="R1113570678193">
      <context>
        acquiring<[^\.]*:neTag:NE:WHO:ORGANIZATION>>[^\.]*
        act<[^\.]*:neTag:NE:WHAT:STATUS>></context>
      <conclusion>extractResult()</conclusion>
      <cscase id="14582">
        <location>14582-R1113570678193.xml</location>
      </cscase>
    </rule>
    <exception>
      <rule id="R1113570811459">
        <context>
          acquired<[^\.]*:neTag:NE:WHO:ORGANIZATION>>[^\.]*
          act<[^\.]*sell[^\.]*:neTag:NE:WHAT:STATUS>></context>
        <conclusion>extractResult()</conclusion>
        <cscase id="1032">
          <location>1032-R1113570811459.xml</location>
        </cscase>
      </rule>
      <exception>
        <rule id="R1113570925553">
          <context>
            acquired<[^\.]*:neTag:NE:WHO:ORGANIZATION>>[^\.]*
            act<[^\.]*sell[^\.]*:neTag:NE:WHAT:STATUS>>[^\.]*its
            acquired_unit<[^\.]*:neTag:NE:WHO:ORGANIZATION>></context>
          <conclusion>extractResult()</conclusion>
          <cscase id="12890">
            <location>12890-R1113570925553.xml</location>
          </cscase>
        </rule>
        <exception>
          <rule id="R1113894751843">
            <context>
              acquired<[^\.]*:neTag:NE:WHO:ORGANIZATION>>[^\.]*
              act<[^\.]*sell[^\.]*:neTag:NE:WHAT:STATUS>>[^\.]*its
              acquired_unit<[^\.]*:neTag:NE:WHO:ORGANIZATION>>[^\.]*
              acquiring<[^\.]*:neTag:NE:WHO:ORGANIZATION>></context>
            <conclusion>extractResult()</conclusion>
            <cscase id="12937">
              <location>12937-R1113894751843.xml</location>
            </cscase>
          </rule>
        </exception>
      </rule>
    </exception>
  </rule>
  <rule id="R1113570876376">
    <context>
      acquiring<[^\.]*:neTag:NE:WHO:ORGANIZATION>>[^\.]*
      act<[^\.]*:neTag:NE:WHAT:STATUS>>[^\.]*
      acquired<[^\.]*:neTag:NE:WHO:ORGANIZATION>></context>
    <conclusion>extractResult()</conclusion>
    <cscase id="10025">
      <location>10025-R1113570876376.xml</location>
    </cscase>
  </rule>
</exception>

```