
An unsupervised method for learning probabilistic first order logic models from unstructured clinical text

Rahul Jha
Dragomir Radev

RAHULJHA@UMICH.EDU
RADEV@UMICH.EDU

Computer Science and Engineering, 2260 Hayward Street, Ann Arbor, MI 48109-2121

Abstract

We present a new unsupervised approach for learning probabilistic first order logic models from unstructured clinical text. We use Carroll, a system that generates a shallow semantic parse of natural language text, to extract predicates out of natural language text. These predicates are then used to learn a simple probabilistic first order logic model of the underlying data. We present our work on automatically learning these models and show some preliminary results obtained by modelling a public clinical database.

1. Introduction

Generating robust prediction models from clinical data is a widespread concern in the medical community. A lot of this data is buried in free text documents however, and is not available in a relational form that can be used for direct modelling. To be able to create sophisticated prediction models, we need to be able to extract relational data from this text. A lot of time and money would need to be invested to manually identify the relevant patterns in the textual data and convert the data to a structured form for analysis. Unsupervised approaches to find these correlations from structured data exist (Johnson, 1996), (Butte & Kohane, 1999), but no method for detecting such patterns from unstructured text in an unsupervised way was found. With supervised approaches (Visweswaran et al., 2003), we believe the need to generate enough labelled data can be a limiting factor in applying these methods to a wide range of clinical data.

In this paper, we present a general approach for learning prediction models from unstructured clinical text.

Very few assumptions are made about the nature of text, which will be elaborated further in section 5. This work can be applied successfully to a diverse range of clinical data.

We now present an outline of our system. After some normalization of the unstructured clinical text, we pass it through a semantic parser that produces an intermediate first order logic representation called logic forms (Moldovan & Rus, 2001). We use some heuristics to select a subset of the predicates discovered from the above process. This subset of predicates is then used to create a first order probabilistic model of the information in the text by using them as evidence variables in a generalized model created in BLOG, a language for defining probabilistic first order models (Milch & Russell, 2007). The complete workflow is as shown in Figure 1.

We use a first order probabilistic model for our system because it allows us to create a generalized model by expressing dependencies between classes of random variables. This model can be then made robust using evidence as and when it arrives. In a propositional probabilistic language such as Bayesian Networks, the number of random variables needed to define a scenario grows with the number of objects. First order probabilistic models, however, allow us to express these models concisely by abstracting over objects. For example, in BLOG, we can create a general model with the random variable classes: `Patient`, `Ailment`, and the random function: `Has`. We express a dependency between these by saying:

`Has(Patient, Ailment) = 0.01`

That is all patients have all ailments with a small probability. Now every time the model encounters an evidence like, `Has(John.Doe, Hypertension)`, where `John.Doe` is an object of type `Patient` and `Hypertension` is an object of type `Ailment`, it revises the probability values of the function `Has` to take them into account. Eventually, the probability values con-

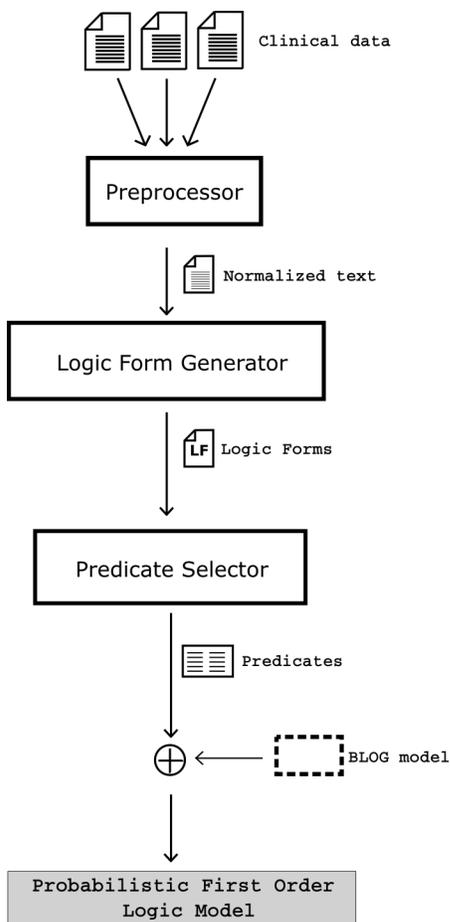


Figure 1. Outline of the workflow

verge to their actual values. This can automatically scale to as many patients and ailments as needed. Further details on BLOG are provided in Section 3.

The main contributions of this paper are:

1. A general first order probabilistic model for clinical data.
2. A method for learning evidence variables for the model directly from unstructured text.
3. A logic forms extractor, Carroll.

The rest of this paper is organized as follows. We provide an introduction to Logic Forms in section 2. Section 3 gives a short background of First order probabilistic languages and discusses BLOG in more detail. We present our model in section 4. Some preliminary results are presented in section 5, and section 6 ends with conclusion and pointers for future work.

2. Logic Forms

Logic forms are a shallow semantic representation of natural language that can be used as an intermediary step between syntactic parse and a deep semantic form. They follow a Davidsonian approach (Davidson, 1967) and represent semantic information in the form of predicates. As an example, consider the following sentences:

1. Some students like to study in the mornings.
2. When he handed in his homework, he forgot to give the teacher the last page.

They correspond to the following logic forms:

1. student (x_1) like (e_4, x_1, e_5) to (e_4, e_5) study (e_5, x_1, x_2) in (e_5, x_2) morning (x_2)
2. When (e_{11}, e_{10}) he (x_4) hand (e_{10}, x_4, x_1) in (e_{10}) his (x_1) homework (x_1) he (x_6) forget (e_{11}, x_6, e_{12}) to (e_{11}, e_{12}) give (e_{12}, x_6, x_3, x_2) teacher (x_2) last (x_3) page (x_3)

In this approach, all e-variables represent eventuality of the action, state or event stated by the verb to take place. All x-variables represent either the subject or the objects of these events. In the first example, the variable x_1 stands for the student and x_2 stands for morning. e_5 represents the fact that the event **study** relates **student**(x_1) and **morning**(x_2). e_4 represents the state of **like** between **student**(x_1) and **study**(x_5).

The benefit of using logic forms is that they are close enough to the syntactic representation to be generated by simple systems, yet contain important semantic relationships. They were first applied by (Moldovan & Rus, 2001) to create axioms from WordNet glosses that improve the performance of a Question Answering System. (Mulkar et al., 2007) used logic forms to create knowledge bases from Chemistry and Biology texts. They further extended this work to read syntactically complex biological texts (Mulkar et al., 2007) and discover causal and temporal ordering relations in biomedical articles (Rutu Mulkar-Mehta & Zhou, 2009).

For this work, we could not find a freely available logic form extractor that suited our needs. We therefore

Table 1. Evaluation results for Carroll on SenseEval 3 data. AMS, LCC, MITRE, SYD stand respectively for systems created by teams from University of Amsterdam, Language Computer Corporation, MITRE Corporation and University of Sydney as part of the SenseEval3 Logic Forms task.

TEAM	ARGUMENT LEVEL		PREDICATE LEVEL	
	PRECISION	RECALL	PRECISION	RECALL
AMS	0.729	0.691	0.819	0.783
LCC	0.776	0.777	0.876	0.908
MITRE	0.734	0.659	0.839	0.781
SYD	0.763	0.655	0.839	0.849
CARROLL	0.421	0.409	0.531	0.548

created a simple rule based system called Carroll to extract logic forms out of natural language. To extract the logical form out of a sentence, Carroll first extracts a dependency parse of the sentence using the stanford dependency parser (de Marneffe & Manning, 2008). It then use a set of simple rules to extract the logic forms from this parse. We tested Carroll’s performance on the test data used in the SenseEval 3 Logic Forms Task (Rus, 2004) and the results are summarized in Table 1 along with the results reported by other teams that participated in the SenseEval 3 Task.

Even though the accuracy for Carroll is low compared to other teams on this data set, it’s good enough for extraction of simple relational predicates that we require for this work. We are managing Carroll as an open source project and are currently working to improve its performance.

3. Probabilistic first order logic systems

3.1. Background

Probabilistic systems like Bayesian networks offer a practical tool for modelling probability distributions. However, the propositional nature of these systems makes them much less powerful and flexible than a first order logic representation. Inference over probabilistic first order logic representations quickly becomes intractable though, and the problem of efficiently reasoning over such models has received considerable attention in the research community. Many systems are now available that make it practical to express first order probabilistic models and perform inference over them. Specifically, the field of statistical relational learning (Getoor & Taskar, 2007) has generated a lot of work in this area.

There are several approaches for doing first order prob-

abilistic inference (Braz et al., 2008). Stochastic Logic Programs (Muggleton, S., 2000) achieve this by defining a distribution over proofs from a given logic program. Probabilistic Relational Models (PRMs) (Friedman et al., 1999) use the formalism of Frame Systems as a starting point. Markov Logic Networks (MLNs) (Richardson & Domingos, 2006) is another popular framework based on undirected models.

Another approach for doing inference over such models is by sampling. Two systems that take this approach are PRISM (Sato & Kameya, 1997) and BLOG (Milch & Russell, 2007). For this work, we have used BLOG, which we cover in further detail in the next section. Please refer to (Braz et al., 2008) for a complete survey of first order probabilistic systems.

3.2. Bayesian Logic (BLOG)

BLOG (Milch & Russell, 2007) or Bayesian Logic is a language that allows defining probability distributions over relational structures with varying sets of objects. A BLOG model defines a typed first-order language for a specific problem. It specifies certain non-random aspects of the scenario and specifies a probability distribution for the rest of the aspects. Inference in BLOG is done by a sampling based approximate inference algorithm. BLOG has well defined procedural and declarative semantics. Procedurally, a BLOG model defines a generative process that generates possible worlds. Declaratively, a BLOG model defines a contingent bayesian network over the basic random variables of the model.

We now introduce a few BLOG keywords that we will be using in our model. For a more complete description, the user is referred to the BLOG documentation.

3.2.1. TYPE

The keyword `type` introduces a new variable type. So, `type Color` introduces a new variable type called `Color`.

3.2.2. GUARANTEED

The variables introduced by the `guaranteed` keyword are guaranteed to exist in all possible worlds. The `guaranteed` keyword also asserts that the objects being declared are distinct objects, and declares constant symbols for those objects. So,

`guaranteed Color Blue, Green`

asserts that these two colors are guaranteed to exist in all possible worlds.

3.2.3. NONRANDOM

This keyword is used for functions, and declares the function as having the same interpretation in all possible worlds. Nonrandom function symbols can be used to represent known information about guaranteed objects. For example,

```
nonrandom Boolean Teaches(Professor,
Course);
```

3.2.4. RANDOM

This keyword is used to define functions whose values vary across possible worlds. Their values are set by steps in the generative process. Every random function should have a dependency statement that specifies its probability distribution. Here is an example,

```
random Boolean Likes(Professor, Professor);
```

```
Likes(p1, p2) {
```

```
if (p1 = p2) then = true
else TabularCPD[[0.2, 0.8]] };
```

TabularCPD just defines the initial conditional probability distribution for the function.

3.2.5. OBS

obs allows us to record an observation of a random variable. These statements can be used to provide evidence for the model.

```
obs GradeObtained(John, Stat10) = C;
```

3.2.6. QUERY

This allows a user to query the probability distribution for a function or a variable. For example

```
query GradeObtained(John, CS106);
```

This returns the probability distribution for John in CS106 over all the possible grades.

A complete BLOG model usually contains three files:

1. **Model file** This defines the basic model for the problem. This is where we define the variable types and functions. This is also where we write the dependency statements for the random variables.
2. **Evidence file** This is where we provide the guaranteed objects(although you could do this in Model file also), specify the non random functions and add observations for the random functions.

3. **Query file** This is where we write our query statements.

4. Modelling clinical data

4.1. Data

We used the MIMIC II clinical database from Physionet (Saeed et al., 2002). MIMIC II provides de-identified discharge summaries for about 25, 000 patients as flat files. For this experiment, we only used the “HISTORY OF PRESENT ILLNESS” sections of about 900 discharge summaries. This contains the data for patients in free text form. Here is a sample:

```
The patient is a 76-year-old male who had
been hospitalized at the [**Hospital 3723**]
from [**09-27**] through [**10-05**] of 2002
after undergoing a left femoral-AT bypass
graft and was subsequently discharged to a
rehabilitation facility.
```

We did some pre-processing on these files to make them suitable for our purposes. Certain special characters in the identifiers had to be normalized for the Stanford parser to function correctly on the data. Also, we used the Metamap (Aronson, 2001) system to identify UMLS concepts in the text and normalize them so that the dependency parser can recognize them as a single word. We followed a simple scheme of replacing each UMLS phrase of interest by a concatenated word representation of the UMLS gloss. For example, “Coronary Artery Bypass Graft” is replaced by `coronary_artery_bypass_graft`. This allows us to normalize over these concepts across the corpus and leads to better accuracy.

4.2. Predicate Extraction

After cleaning up and normalizing the data as above, we ran Carroll on the resulting files to obtain the logic forms for each of the sentences. The logic forms representation is still not suitable for creating models though. We used some simple heuristics to extract predicates from this form that we can use in our models.

A simplifying assumption in case of discharge summaries is that they always talk about the patient. Thus all pronominal references can be resolved to the patient variables a priori. We then find all the two argument predicates that the patient variables participate in. We extract these specific predicates and then sort them into buckets based on predicate names. So for example `report(e1, x1, x2)` and `report(e2, x3, e5)` would fall into the same bucket,

Table 2. Occurrence frequency for different predicates for clinical data

PREDICATE	COUNT
BE	253
HAVE	75
DENY	37
RECEIVE	14
UNDERGO	12
REPORT	5
USE	2
SEE	2
REQUIRE	2
NOTE	2
MOVE	2
DO	2
DEVELOP	2
COMPLAIN	2
WISH	1
TOLERATE	1
TAKE	1
SUFFERED	1

while `experience(e3, x1, x9)` would go into a different bucket. Upon organizing the predicates in this way, we got the distribution shown in Table 2.

As evident, most of the relations are located in the top 5 predicates (there are several more 1 count predicates that are not shown). We decided to ignore the copula predicate `be` since its generic nature would lead to an inconsistency in the final model. We then took the top four predicates from the distribution and built our model using data from these predicates only.

4.3. BLOG model

To model this data, we created the following BLOG model.

```

type Patient;
type Ailment;
type Event;

nonrandom Boolean have(Patient, Event);
nonrandom Boolean deny(Patient, Event);
nonrandom Boolean receive(Patient, Event);
nonrandom Boolean undergo(Patient, Event);

random Boolean LeadsTo(Event, Ailment);
random Boolean Has(Patient, Ailment);

LeadsTo(s,a) ~ TabularCPD[[0.3, 0.7]]();

Has(p, a) {

```

```

if (exists Event e ((have(p, e) |
                    deny(p, e) |
                    receive(p, e) |
                    undergo(p, e))
                    & LeadsTo(e, a)))
    then ~ TabularCPD[[0.8,0.2]]

else ~ TabularCPD[[0.2,0.8]]

}];

```

In this model, there are only three types of variables : Patients, Ailments and Events. For simplification, we combined the objects of `have`, `deny`, `receive` and `undergo` into the same variable called `Event`.

Each of the predicates described in the last section is represented by a nonrandom variable. We also have two random variables: `LeadsTo` is the probability of an event causing a particular ailment, `Has` is the probability of a patient having a particular ailment. `LeadsTo` is a hidden variable in the sense that we never directly observe this variable, but we believe it to influence the probability distribution for the other random variable, `Has`, as expressed in its dependency statement.

Once we have this model, the only thing left is to generate the evidence variables. Given the logic form predicates we have already extracted, this is a pretty straightforward operation. We present a small snapshot of the resulting evidence file below.

```

guaranteed Patient patient_00098, ... ;

guaranteed Ailment cardiothoracic_surgery, ... ;

guaranteed Event stenosis, ... ;

nonrandom Boolean receive(Patient, Event)
= ListInterp[2 , patient_04110, evaluation_renal,
              patient_01862, Solu_Medrol,
              ...
              patient_00554, nitroglycerin
              ];

nonrandom Boolean have(Patient, Event)
= ListInterp[2 , patient_00098, headache,
              patient_03445, history,
              ...
              patient_04271, recent_vomiting,
              ];

```

```

...      cardiothoracic_surgery)
...
0.9877429364259752      true
obs Has(patient_00098, neurosurgery) = true;      0.012257063574044977      false
obs Has(patient_02244, neonatology) = true;

```

... Thus within the limited amount of data that is available in the model, it performs as expected. Having a larger corpus of data would increase the accuracy of the predictions.

Using these model and evidence variables, we can now run queries on the system.

5. Results

We obtained results for a total of 137 patients. 38 different ailments and 110 different events were recorded. Out of the 137 patients, the ailments for 88 of them could be extracted from the discharge summaries.

We ran several queries on this model and verified the results. Here are a few sample query results:

```

LeadsTo(coronary_artery_bypass_graft,
        cardiothoracic_surgery)

```

```

0.6285624797326435      false
0.3714375202673559      true

```

```

LeadsTo(fevers,
        hepatobiliary_surgery)

```

```

0.5150551353697476      false
0.4849448646302515      true

```

```

LeadsTo(aortic_valve_replacement,
        neurosurgery)

```

```

0.9926608220001478      false
0.007339177999866096      true

```

We also tested the system by creating a dummy patient and adding the non random facts for this dummy patient reporting `coronary_artery_disease` and `haematocrit`. We chose these because there are two patients in our model who reported these and were observed to have `cardiothoracic_surgery`. We now ran a query for estimating the probability of the dummy patient having a `cardiothoracic_surgery` and obtained the following results:

```

Distribution of values for Has(patient_dummy,

```

We notice that several noise variables enter the system as a result of unsupervised semantic parsing. However, since they never occur with a high frequency, they do not affect the probabilities dramatically. The queries take about 15 minutes on an average on a machine with a 2 GHz Intel Core 2 Duo and 2GB main memory.

For further evaluation, we plan to use another probabilistic system to model the same data, and then compare the probability estimates for equivalent assertions. This system could either be another first order probabilistic system or a propositional system like Bayesian Networks. We are currently working on putting such an evaluation in place.

6. Conclusion and Future work

We present a simple approach for generating first order predicate logic models from unstructured clinical text. Even though this was a limited experiment, we got several key insights. From our experience, we believe this is a promising area for further research.

There are several improvements that can be made to our present system. We are not able to extract complex lexical chains from the logic forms at this time, which is something that would increase the accuracy of the models. We would also want to be able to generate more sophisticated BLOG models than is possible at the moment. A third concern is to increase the performance of the logic form generator itself so that more robust systems can be built on top of it. Our final concern is the creation of good evaluation metrics and gold standards for such systems that can be used to evaluate future work in this area.

Acknowledgments

This work was supported in part by the NIH Grant U54 DA021519. We would like to thank Dr. Brian Athey, Prof. H.V. Jagadish and Dr. K.P. Unnikrishnan for their guidance. Thanks to Professor Zee-shan Syed for some excellent suggestions. Thanks also to Amjad Abu Jbara, Shiwali Mohan and Samuel Smolkin for their help and support.

References

- Aronson, A. R. Effective mapping of biomedical text to the UMLS Metathesaurus: the MetaMap program. *Proceedings / AMIA ... Annual Symposium. AMIA Symposium*, pp. 17–21, 2001. ISSN 1531-605X.
- Braz, Rodrigo, Amir, Eyal, and Roth, Dan. A survey of first-order probabilistic models. In Holmes, Dawn and Jain, Lakhmi (eds.), *Innovations in Bayesian Networks*, volume 156 of *Studies in Computational Intelligence*, chapter 12, pp. 289–317–317. Springer Berlin / Heidelberg, Berlin, Heidelberg, 2008. ISBN 978-3-540-85065-6.
- Butte, A. J. and Kohane, I. S. Unsupervised knowledge discovery in medical databases using relevance networks. *Proceedings / AMIA ... Annual Symposium. AMIA Symposium*, pp. 711–715, 1999. ISSN 1531-605X.
- Davidson, Donald. The Logical Form of Action Sentences. In Rescher, Nicholas (ed.), *The Logic of Decision and Action*. University of Pittsburgh Press, 1967.
- de Marneffe, Marie-Catherine and Manning, Christopher D. The Stanford typed dependencies representation. In *Coling 2008: Proceedings of the workshop on Cross-Framework and Cross-Domain Parser Evaluation*, pp. 1–8, Manchester, UK, August 2008. Coling 2008 Organizing Committee.
- Friedman, Nir, Getoor, Lise, Koller, Daphne, and Pfeffer, Avi. Learning probabilistic relational models. In *In IJCAI*, pp. 1300–1309. Springer-Verlag, 1999.
- Getoor, Lise and Taskar, Ben. *Introduction to Statistical Relational Learning (Adaptive Computation and Machine Learning)*. The MIT Press, 2007. ISBN 0262072882.
- Johnson, S. B. Generic data modeling for clinical repositories. *Journal of the American Medical Informatics Association : JAMIA*, 3(5):328–339, 1996. ISSN 1067-5027.
- Milch, Brian and Russell, Stuart. First-Order probabilistic languages: Into the unknown. *Inductive Logic Programming*, pp. 10–24, 2007.
- Moldovan, Dan and Rus, Vasile. Logic form transformation of wordnet and its applicability to question answering. In *In Proceedings of ACL 2001*, pp. 394–401, 2001.
- Muggleton, S. Learning stochastic logic programs. *Proceedings of the AAAI2000 Workshop on Learning Statistical Models from Relational Data*, 5, 2000.
- Mulkar, R., Hobbs, J. R., Hovy, E., Chalupsky, H., and Lin, C. Y. Learning by reading: Two experiments. In *3rd International Workshop on Knowledge and Reasoning for Question Answering*, Hyderabad, India, 2007.
- Richardson, Matthew and Domingos, Pedro. Markov logic networks. *Machine Learning*, 62(1):107–136, February 2006. ISSN 0885-6125.
- Rus, Vasile. A first evaluation of logic form identification systems. In *SENSEVAL-3: Third International Workshop on the Evaluation of Systems for the Semantic Analysis of Text*, 2004.
- Rutu Mulkar-Mehta, Jerry R. Hobbs, Chun-Chi Liu and Zhou, Xianghong Jasmine. Discovering causal and temporal relations in biomedical texts. AAAI Press, 2009.
- Saeed, M., Lieu, C. and Raber, G., and Mark, R. G. MIMIC II: a massive temporal ICU patient database to support research in intelligent patient monitoring. *Computers in Cardiology*, 29:641–644, September 2002.
- Sato, Taisuke and Kameya, Yoshitaka. Prism: a language for symbolic-statistical modeling. In *Proceedings of the 15th International Joint Conference on Artificial Intelligence (IJCAI'97)*, pp. 1330–1335, 1997.
- Visweswaran, Shyam, Hanbury, Paul, Saul, Melissa, and Cooper, Gregory F. Detecting adverse drug events in discharge summaries using variations on the simple bayes model. *AMIA Annu Symp Proc*, pp. 689–93, 2003. ISSN 1942-597X.