

# Combining Truth Discovery and Open Information Extraction with Active Ensembling

Mouhamadou Lamine BA  
Qatar Computing Research Institute  
Tornado Tower, West Bay  
Doha, Qatar  
mlba@qf.org.qa

Laure Berti-Equille  
Qatar Computing Research Institute  
Tornado Tower, West Bay  
Doha ; Qatar  
lberty.qf.org.qa

## ABSTRACT

Web search engines or open information extraction systems usually reply to users' queries with a set of candidate answers that are often conflicting because they are claimed by multiple information sources. In this context, estimating information veracity is difficult for the users especially when they have no prior knowledge about the trustworthiness of the sources. In this demonstration paper, we showcase a system that supports event/entity and relation extraction based on keyword-search from the Web, processes the conflicting outputs, combines multiple truth finding algorithms with active learning to provide the most likely true answers and determine the most trustworthy sources.

## 1. INTRODUCTION

[Lamine: Page allocation]

- 1.25 pages -> abstract + introduction
- 1.5 pages -> Open information extraction + Active Ensembling for Truth Discovery
- 1 pages -> Demonstration System + Scenario
- 0.25 pages -> References

## 2. OPEN INFORMATION EXTRACTION

- décrire le type d'information auquel on s'intéresse par exemple "factoid claim"
- décrire le système sur lequel on se base
- décrire comment on transforme l'output de OpenIE
- donner qq exemples

In this study, we are interested on truth discovering on the large number of "factoid" statements (or claims) made by multiple Web sources in the same real-world domain. A "factoid" claim, e.g., *Barack Obama was born in Kenya*, is a piece of unverified or inaccurate information that is presented as factual, often as part of a publicity effort. Such type of claims is usually accepted as true because of its frequent redundancy over multiple sources. We focus on conflicting factoid claims provided by typical open Web information extraction systems as answers to users' queries.

Concretely, given user input, we retrieve the set of candidate claims, together with the associated sources, returned by TextRunner<sup>1</sup> from a Web corpus. We then format this output in such a way that fits our truth discovering process.

**TextRunner extractor.** The extraction engine relies on an unsupervised extraction procedure which queries, using a single and data-driven pass, an entire Web corpus of unstructured texts and extracts a list of candidate relational tuples which might satisfy a given user input query. A user query, in this setting, consists typically of a sentence formed by three components: two real-world *entities* and a certain *relation*. The goal of TextRunner being to find and gather the set of possible claims on the Web which support or not the specified relation between the two entities. A user query  $q$  can be defined formally with the help of a triplet  $(e_1, r, e_2)$  where  $e_1$  and  $e_2$  are real-world entities and  $r$  is a relation;  $r$  indicates a possible relationship between the two given entities. If some components, i.e., the second entity, of the query are not provided by the user, we say that the query is *incomplete*. Incomplete queries are common in Web search as human often has a partial knowledge of the real world: this motivates us the use of information extractors in order to complete his knowledge. Note that the information extraction systems do also well in the case of *incomplete* queries.

The output of TextRunner is a set of candidate claims which are ranked according to their number of sources. A claim in this context can correspond to a tuple, a relation, or an real-world entity. The system also offer the possibility, through pointers, to access to these sources and the full corpus in which the each has been extracted. We would harness these pointers in order to format the result obtained from the extraction system in such the way that it fits the input of our different truth finding algorithms. We show in the following how such a formatting is realized.

**Data formatting.** Given a user input tuple  $q = (e_1, r, e_2)$ , we refer to the set of  $n$  claims  $c_1 \dots c_n$  extracted and returned by the extraction engine as potential answers. We assume that the query admits only one true answers, thereby we are in the presence of  $n$  conflicting answers. For each returned claim, we go through the linked pointer and extract the corresponding source by using hand-written mapping rule. We denote the set of  $m$  sources queried by the extraction system in order to answer  $q$  by  $s_1, \dots, s_m$ .

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

1. TextRunner is an online at <http://openie.allenai.org/>

### 3. ACTIVE ENSEMBLING FOR TRUTH FINDING

#### 3.1 Ensembling

- donner idée générale pour introduire ce qu'est l'ensembling
- on a besoin de le faire dans le contexte de truth discovery car aucune methode ne bat toutes les autres dans tous les cas de figure - donc on combine les methodes : il y a plusieurs façon de combiner par ex. consensus de méthodes, etc.
- expliquer quelles méthodes on combine avec leurs avantages et inconvénients

An ensemble is a supervised learning algorithm in the sense that it can be trained and then used to make predictions. An ensembling, or commonly an ensemble-based active learning, is a learning process selecting one classifier type, or appropriate combinations of multiple classifier types, to construct ensembles for a given task.

#### 3.2 Active Learning Process

- notre approche que l'on défend ici dans la démo est semi supervisée en impliquant de l'utilisateur de façon active en lui demandant s'il peut confirmer des faits (facts) - si on a une ground truth partielle on la "rejoue" cas par cas

## 4. OUR DEMONSTRATION SYSTEM

#### 4.1 System Architecture and User Interface

The architecture of our demonstration system, given in Figure 1, comprises the following three main components.

*User I/O Interface.* It represents the main entry point of our application for user interaction. The user I/O interface is composed by a text search area where a given user can enter its search keywords, in terms of a relation. The final result of the overall process will be also shown to the users through this component. Finally, the user gives it feedbacks via the user I/O interface through the button options or the form.

*Information extraction module.* This is the information extraction module which considers the input of the user and browses several Web sources in order to return the relevant answers. In our system, we rely on TextRunner in order to extract information from Web corpus.

*Truth Finding Engine.* It corresponds to AllegatorTrack which contains twelve truth finding algorithms with different accuracy according to the types of claims and the characteristics of sources.

*Learning Module.* We have also a learning method that uses our knowledge bases of users' feedbacks. It enables to learn about the best truth finding algorithms, among the twelve, to use with respect to the type of entities or relations searched by the user.

*Knowledge Base.* The knowledge base contains the information used for the learning phase of the truth finding procedure. These informations include the true facts for some relations which have been learnt based on the feedbacks of the users. In addition, our knowledge base could be enriched with ground truth about some facts from reliable sources

such as Wikipedia. Based on the knowledge base, our system has the ability to improve the accuracy of the truth finding process by learning about the best method to use or the best parameters, e.g., sources' accuracy scores, to consider for a better bootstrapping of the process.

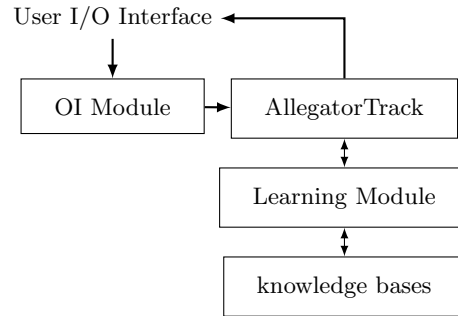


Figure 1: Architecture of our system

#### 4.2 Demonstration Scenario

A given user that wants to interact with our system must do it through the search form. Through the search form, she (or he) provides her searched relation, e.g., "Where is born Barack Obama?". The searched relation is then passed to the information extraction engine, TextRunner system in our case, which returns a set of answers considered to be relevant for the user's request. Each claim in the returned list is processed in order to extract the corresponding sources along with a detailed description of the claim which we format in a certain manner. The set of sources and the formatted versions of all claims are then passed to the truth finding module which integrates all the claims and computes the most probable answer together with the reliability scores of participated sources for the searched relation. Finally, the output of the truth finding process is returned to the user. The user can also want to review the output of our system by definitively validating it or not through its knowledge of the modeled world. For example when the system has totally wrong, it may be interesting to get such a kind of feedbacks from the user in order to change the used method, as there are many available with our system, and to enhance the process for the further search about the same world. The user gives feedbacks using the option buttons on the left-hand side of the outputted claims or the text form. The feedbacks given by the user are saved in knowledge bases within our system for further processes.

## 5. CONCLUSION

□

## 6. REFERENCES

- [1] Oren Etzioni, Michele Banko, Stephen Soderland, and Daniel S. Weld. Open Information Extraction from the Web. *Commun. ACM*, 51(12):68–74, December 2008.
- [2] Naeemul Hassan, Chengkai Li, and Mark Tremayne. Detecting check-worthy factual claims in presidential debates. In *Proc. CIKM*, pages 1835–1838. ACM, 2015.
- [3] Dalia Attia Waguih, Naman Goel, Hossam M. Hammady, and Laure Berti-Equille. AllegatorTrack: combining and reporting results of truth discovery from

multi-source data. In *Proc. ICDE*, pages 1440–1443, 2015.

- [4] Alexander Yates, Michael Cafarella, Michele Banko, Oren Etzioni, Matthew Broadhead, and Stephen Soderland. TextRunner: Open Information Extraction on the Web. In *Proc. NAACL*, pages 25–26. Association for Computational Linguistics, 2007.