

# Objectives for Information Extraction

Justin T. Chiu

June 27, 2019

## Abstract

Many recent information extraction systems predict the relationship between an entity and value given their locations in the text. This requires observed locations of mentions, which requires annotations at the word level. Any form of annotation at the word level does not scale as the size of the text and the number of labels increases, and even more so if there is ambiguity. In order to train a probabilistic information extraction without any supervision at the level of text, we specify a model that, for each word, either chooses a triple from a knowledge base to explain or chooses to explain nothing.

## 1 Problem Setup

In relation extraction (RE) we extract facts from a passage of text. The goal of RE is to convert facts expressed in natural language into a form amenable to computation. RE focuses on a relational representation: facts are extracted into a form that details how values are related to entities. This relational representation is useful for performing computation. Rather than dealing with noisy text, we may perform automatic deduction by composing facts in relational form.

In this work we build a relation extraction system with minimal supervision. We assume supervision only at the proposition level of a KB, with no annotation given at the level of text. In order to perform extraction, we must first identify where facts are. As sentences may contain many entities and values, and therefore possibly a large number of facts, we focus on identifying facts at the word level. In particular, we focus on identifying the location of value mentions and predicting who (the entity) and what (the relation type) the value mention is discussing.

Note on related work: Except for Zeng et al. (2018), prior work has either assumed that the locations of entities and values are given as input features or that the locations of entities and values are observed at training time.

We are primarily concerned with the scenario where we have an overcomplete KB schema with respect to a specific passage of text. This fits many scenarios in real world applications: we may have thousands of entities of interest if our KB was pulled from an external source such as Freebase, but the particular document we wish to analyze only discusses tens of entities, only a few of which are present in our KB.

Table 1: Notation for relation extraction. The most important terms are in the top section with the less important terms given in at the bottom.

$x$	$\triangleq$	A sequence of words $x_i, i \in \mathcal{I}$ corresponding to a sentence or document.
$r$	$\triangleq$	The knowledge base, an indexed set of records each consisting of entities, types, and values. Each $r_j = (e_j, t_j, v_j)$ .
$v$	$\triangleq$	The list of all values in a KB by index.
$f_*(x)$	$\triangleq$	A learned function whose output is of shape $ \mathcal{I}  \times D$ , where $ \mathcal{I} $ is the length of the word sequence and $D$ is the context-dependent.
$c_i$	$\triangleq$	A word level RV that indicates whether a word is content or not. We have $c_t \in \{0, 1\}$ .
$a_i$	$\triangleq$	A word level RV that indicates an the index of the aligned for word $x_t$ . We have $a_t \in J$ .
$z_i$	$\triangleq$	A word level RV that indicates the canonical value word $x_t$ is translated to, such that $z_t \in \mathcal{V}$
$e$	$\triangleq$	The list of all entities in a KB by index.
$t$	$\triangleq$	The list of all types in a KB by index.
$\mathcal{X}$	$\triangleq$	The vocabulary words come from.
$\mathcal{I}, \mathcal{J}$	$\triangleq$	The index set the text and KB respectively. Both are sets $\{0, 1, 2, \dots\} \subset \mathbb{N}$ .

The problem description is as follows: given a text  $x = [x_0, \dots, x_{|I|}]$ ,  $x_i \in \mathcal{X}$ , we must model the facts  $r = \{(e_j, t_j, v_j)\}_{j \in \mathcal{J}}$  expressed in that text with respect to a schema that details all entities  $e_j \in \mathcal{E}$ , relation types  $t_j \in \mathcal{T}$ , and all values  $v_j \in \mathcal{V}$ . We assume that the schema of the KB  $(\mathcal{E}, \mathcal{T}, \mathcal{V})$  is known at all times, and that the schema covers all facts of interest.

We refer to the set of facts  $r$  as a knowledge base (KB), and each fact  $r_j$  is referred to as a record. Each record  $r_j$  consists of an entity, type, and value triple. Let  $e, t, v$  be the projection of the entities, types, and values from the KB  $r$  in aggregate.

Given all entities  $e$  and types  $t$ , we reduce the construction of  $r$  to predicting, for every  $j \in \mathcal{J}$ , the value  $v_j$  corresponding to the entity  $e_j$  and type  $t_j$ . We propose to model  $p(v \mid x, e, t)$ , the distribution over values of corresponding entities and types given the text.

## 2 Model

We define a directed graphical model that jointly models word and KB level extraction. The model first predicts the word level information, then aggregates its predictions at the KB level.

We introduce the following random variables for each index  $i \in [I]$ :

1.  $c_i$  indicates whether a word is a value mention. We have  $c_i \in \{0, 1\}$ .
2.  $a_i$  indicates the index of the record aligned to word  $x_i$ . We have  $a_i \in J$ .
3.  $z_i$  indicates the canonical value word  $x_i$  is a mention of, such that  $z_i \in \mathcal{V}$ .

We propose the following factorization of the relation extraction system:

$$\begin{aligned}
 p(v \mid x, e, t) &= \sum_{z, a, c} p(v, z, a, c \mid x, e, t) \\
 &= \sum_{z, a, c} p(v \mid z, a, c, x, e, t) \prod_i p(z_i, a_i, c_i \mid x, e, t) \\
 &= \sum_{z, a, c} \prod_j p(v_j \mid z, a, c, x, e, t) \prod_i p(z_i \mid x) p(a_i \mid x, e, t) p(c_i \mid x)
 \end{aligned} \tag{1}$$

The model  $p(v \mid x)$  has the following generative story. Let  $f_*$  be arbitrary learned functions that maps each  $x_i$  in  $x$  to a point on a simplex. In all cases,  $f_*$  takes the form of a BLSTM over the sequence of words  $x$  with a transformation applied to the output at each time step  $i$ . The dimensionality of the simplex can be inferred from the domain of the distribution parameterized by  $f_*$ . For each index  $i \in \mathcal{I}$ , we have

1. Value mention  $p(c_i \mid x)$ : We choose whether word  $x_i$  is a value mention or not. Each  $c_i \sim \text{Bern}(f_c(x)_i)$ , where a value of 1 indicates that  $x_i$  is a value mention.
2. Alignment  $p(a_i \mid x, e, t)$ : We align the word  $y_i$  by choosing who (the entity) and what (the relation type) are being talked about. In particular,  $a_i \sim \text{Cat}(f_a(x, e, t)_i)$  denotes the alignment to the record  $r_{a_i}$  given by the index  $j = a_i$ .

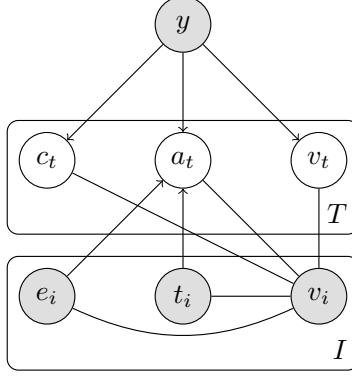


Figure 1: The inference network which predicts word-level values and alignments, then aggregates those choices at the KB level.

3. Translation  $p(z_i | x)$ : We translate the word  $x_i$  into a value  $z_i$ . The  $z_i \sim \text{Cat}(f_z(y)_i)$  is the canonical value from the KB schema associated with  $x_i$ .

Finally, we aggregate the word level information at the sequence level:

1. Aggregation  $p(v | z, a, c) = \prod_j p(v_j | z, a, c)$ : Given the word level  $z, a, c$ , we choose the KB level value  $v_j$  from a CRF. This choice is made independently from the other values in the KB given the word level choices. Let's take care with the conditioning on  $e, t$  for this definition: We parameterize

$$p(v_j | z, a, c, e, t) \propto \prod_i \exp(\psi(v_j, z_i, a_i, c_i, e, t)),$$

with

$$\psi(v_j, z_i, a_i, c_i, e, t) = \begin{cases} \phi(v_j, e, t), & c_i = 0 \\ \theta 1(v_j = z_i, a_i = j), & c_i = 1. \end{cases}$$

We model the marginal distribution  $p(z | x)$  rather than the full joint distribution  $p(v, z | x)$  with an undirected model because we are interested in extracting the facts expressed in the text without the influence of the values in a KB.

### 3 Training

Maximizing  $\log p(v | x)$  directly is very expensive. The aggregation model alone with  $z$  unobserved is a Restricted Boltzmann Machine, where exact inference is  $O(|\mathcal{V}|^{|\mathcal{I}|} + |\mathcal{V}|^{|\mathcal{J}|})$  since we must compute the partition function.

## 4 Three perspectives on training

We can either train  $q$  directly on the conditional task or train it to mimic the posterior of a suitable generative model.

Axes of objectives:

1. Proposal distribution: Learned or uniform (or some prior)
2. Probabilistic interpretation: Marginal likelihood or KL
3. Probabilistic interpretation 2: Approximate posterior of a generative model or learn directly.

## 4.1 Marginal loss

Our first loss, the marginal likelihood, takes the form of

$$\mathcal{L}_{\text{ML}} = \sum_t \log \sum_{a_t} \sum_{c_t} q(a_t, c_t, v_{a_t})$$

where the unobserved alignments and latent copy are marginalized over. This has the interpretation of maximizing the ‘softmax’ of the log probabilities, where softmax is from the physics usage as a smooth approximation to the max. However, the sharpness of is limited by the choice of  $q$ . With a uniform  $q(a, c)$ ,  $\mathcal{L} = \sum_t \sum_{a_t} q(v_{a_t}) + C$ , where  $C$  is the normalization term which we can assume to be constant here. The gradient of this term wrt the log probabilities is weighted by the posterior, so a high entropy  $q$  will result in a smaller gradient.

WRONG! The inference network is a mixture model consisting of  $q(v \mid y)$  and  $p(v)$ ! Derive from generative model.

## 4.2 KL

The second loss, a lower bound on the marginal likelihood, is the following

$$\mathcal{L}_{\text{KL}} = \sum_t \sum_{a_t} \sum_{c_t} q(a_t, c_t) \log q(v_{a_t})$$

The benefit of this loss is that we can approximate it via Monte Carlo sampling. Under a uniform  $q(a, c)$ , we must maximize the probabilities of all values. With a learned  $q$ , the lower bound can be made much tighter so there is less of a difference.

## 4.3 Approximating the posterior of a generative model

The information extraction model itself was inspired by the following generative process: For every word  $y_t$

1. Generate the values of the KB  $v \sim p(v)$  assuming the schema is fixed.
2. Choose  $c_t \sim \text{Bern}(f(y_{<t}))$ , which determines whether to generate word  $y_t$  from a language model or the KB.
3. If  $c_t = 0$ , generate  $y_t \sim \text{Cat}(g(y_{<t}))$ .
4. Otherwise pick an alignment to the KB  $a_t \sim \text{Cat}(h(y_{<t}))$  then generate  $y_t \sim \text{Cat}(f'(v_{a_t}))$ .



Figure 2: The generative model which produces words given a knowledge base.

We then train our model to approximate the posterior  $p(v \mid y)$ . Since this generative model does not use every record during generation, the posterior may only explain a subset of all records. This fits our hypothesis that the KB contains a superset of records expressed in text.

The loss is then given by

$$\mathcal{L}_{VI}$$

The primary benefit of this over a purely conditional approach is semi-supervised information extraction, where values are missing.

$$\arg \max_p \sum_{y'} \sum_{x'} p^*(y', x') \log \frac{p(y', x')}{p^*(y', x')} = \arg \max_p \sum_{y'} \sum_{x'} p^*(y', x') \log \sum_z p(y', z, x') \quad (2)$$

$$= \arg \max_p \sum_{y'} \sum_{x'} p^*(y', x') \log \sum_z p(y', z, x') \quad (3)$$

$$\geq \quad (4)$$

## References

Xiangrong Zeng, Daojian Zeng, Shizhu He, Kang Liu, and Jun Zhao. 2018. Extracting relational facts by an end-to-end neural model with copy mechanism. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 506–514, Melbourne, Australia. Association for Computational Linguistics.