
Alexandria: Unsupervised High-Precision Knowledge Base Construction using a Probabilistic Program

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Creating a knowledge base that is accurate, up-to-date and complete remains
2 a significant challenge despite substantial efforts in automated knowledge base
3 construction. In this paper, we present Alexandria – a system for unsupervised, high-
4 precision knowledge base construction. Alexandria uses a probabilistic program to
5 define a process of converting knowledge base facts into unstructured text. Using
6 probabilistic inference, we can invert this program and so retrieve facts, schemas
7 and entities from web text. The use of a probabilistic program allows uncertainty
8 in the text to be propagated through to the retrieved facts, which increases accuracy
9 and helps merge facts from multiple sources. Because Alexandria does not require
10 labelled training data, knowledge bases can be constructed with the minimum of
11 manual input. We demonstrate this by constructing a high precision (typically
12 97%+) knowledge base for people from a single seed fact.

13 1 Introduction & related work

14 Search engines and conversational assistants require huge stores of knowledge in order to answer
15 questions and understand basic facts about the world. As a result, there has been significant interest
16 in creating such knowledge bases (KBs) and corresponding efforts to automate their construction and
17 maintenance (see [15] for a review). For example, KnowledgeVault [5], NELL [1, 12], YAGO2 [8]
18 and many other systems aim either to construct a KB automatically or make an existing KB more
19 complete. Despite these efforts, there remain significant ongoing challenges with keeping KBs
20 up-to-date, accurate and complete. Existing automated approaches still require manual effort in
21 the form of at least one of the following: **a provided set of entities** used for supervised training of
22 components such as entity linkers/recognizers; **a provided schema** used to define properties/relations
23 of entities; or **a provided set of annotated texts** used to train fact extractors/part of speech taggers.
24 The holy grail of KB construction and maintenance would be a system which could learn and update
25 its own schema, which could automatically discover new entities as they come into existence, and
26 which could extract facts from natural text with such high precision that no human checking is needed.

27 With this goal in mind, we present Alexandria – a system for unsupervised, high-precision knowledge
28 base construction. At the core of Alexandria is a probabilistic program that defines a process of
29 generating text from a knowledge base consisting of a large set of typed entities. By applying
30 probabilistic inference to this program, we can reason in the inverse direction: going from text back to
31 facts. This inverse reasoning allows us to retrieve facts, schemas and entities from web text. The use
32 of a probabilistic program also provides an elegant way to handle the uncertainty inherent in natural
33 text. An important advantage of using a generative model is that Alexandria does not require labelled
34 data, which means it can be applied to new domains with little or no manual effort. The model is
35 also inherently task-neutral – by varying which variables in the model are observed and which are
36 inferred, the same model can be used for: learning a schema (relation discovery), entity discovery,

entity linking, fact retrieval and other tasks, such as finding sources that support a particular fact. In this paper we demonstrate schema learning, fact retrieval, entity discovery and entity linking. We will evaluate the former two tasks, while the latter two are performed as part of these main tasks.

An attractive aspect of our approach is that the entire system is defined by one coherent probabilistic model. This removes the need to create and train many separate components such as tokenizers, named entity recognizers, part-of-speech taggers, fact extractors, linkers and so on. A disadvantage of having such multiple components is that they are likely to encode different underlying assumptions, reducing the accuracy of the combined system. Furthermore, the use of a single probabilistic program allows uncertainty to be propagated consistently throughout the system – from the raw web text right through to the extracted facts (and back).

Related work – There has been a significant amount of work on automated knowledge base construction [15]. Because the Alexandria system can be used to perform many different tasks, it is related to a range of previous, task-specific systems. Here we describe the most relevant.

Unsupervised learning – Open IE (Information Extraction) systems, such as Reverb [6] and OLLIE [9] aim to discover information across multiple domains without having labelled data for new domains. Such systems do not have an underlying schema, but instead retain information in a lexical form. This can lead to duplication of the same fact stored using different words, or can even allow conflicting facts to be stored. Representing facts in lexical form makes them hard for applications to consume, since there is no schema to query against. In contrast, Alexandria aims to infer the underlying schema of new domains and to extract facts in a consistent form separate from their lexical representation.

Schema learning – the closest existing work is Biperpedia [7] which aims to discover properties for many classes at once. Biperpedia uses search engine query logs as well as text to discover attributes, in a process that involves a number of trained classifiers and corresponding labelled training data. Alexandria’s key differences are its unsupervised approach and the fact that schema learning is integrated into a single probabilistic model also used to perform other tasks.

Web scale fact extraction – several existing systems can extract facts against a known schema across the entire web. These include KnowledgeVault [5], NELL [12] and DeepDive [17]. Of these, KnowledgeVault is the largest scale and has performed KB completion (filling in missing values for entities where most values are known) of over 250M facts. However, KnowledgeVault makes extensive use of known values about an entity and so cannot discover new entities. NELL does not suffer from this limitation, but retrieves values with comparatively low precision ($\sim 80 - 85\%$). DeepDive is perhaps the most similar system to Alexandria in that it is based around a probabilistic model – a Markov Logic Network (MLN) [13]. DeepDive uses hand-constructed feature extractors to extract candidate facts from text for incorporation into the MLN. Because Alexandria uses a generative model of text, it can be applied directly to web text, without the need for feature extractors.

2 The Probabilistic Program

In recent years, there has been much interest in probabilistic programming as a powerful tool for defining rich probabilistic models [4, 3]. The Alexandria probabilistic program makes full use of this power, through language features such as polymorphism, objects and collections. To our knowledge, this probabilistic program is also the largest and most complex deployed in any application today. Using a rich, complex probabilistic program allows the process of generating text from facts to be modelled very precisely, leading to high accuracy. The complete program generates: a schema for an entity type, which consists of a set of named, typed properties; a knowledge base, with a large set of typed entities; and a number of text extracts, each describing an entity from the knowledge base in natural language. We use an extended form of the Infer.NET inference engine [10] to invert this probabilistic program, and so infer a schema and knowledge base given a very large number of web text extracts. A major contribution of this work is to show that such large scale inference can be performed in complex probabilistic programs, as will be discussed in Section 4.

We show probabilistic programs in C#, which allows for compact representation of the model. In these programs, the keyword `random` takes a distribution and returns an uncertain value with that distribution. The `Uniform` function takes a collection of objects and returns uniform distribution over these objects. Given these definitions, we now describe each section of the probabilistic program in detail.

```

// Loop over properties in the schema
for(int i=0;i<props.Length;i++) {
    // Pick number of names from geometric dist
    int numNames=random Geometric(probNames);
    // Allocate array for storing names
    var names=new string[numNames];
    // Pick names from property name prior
    for(int j=0;j<numNames;j++) {
        names[j]=random Property.NamePrior;
    }
    // Set generated strings as property names
    props[i].Names=names;

    // Pick one built-in type prior at random
    var typePrior=random Uniform(typePriors);
    // Draw a type instance from the type prior.
    props[i].Type=random typePrior;
}

// Create set of entities (of the same type)
Entity[] entities = new Entity[entityCount];
// Loop over entities of this type
for(int j=0;j<entities.Length;j++) {
    // Loop over properties in the schema
    for(int i=0;i<props.Length;i++) {
        // Pick number of alts from geometric dist
        int numAlts = random Geometric(probAlt);
        object[] alts = new object[numAlts];
        // Loop over alternatives
        for(int k=0;k<alts.Length;k++) {
            // Choose a property value from the prior
            alts[k]=random props[i].Type.Prior;
        }
        // Set alternatives as the property value
        entities[j][i]=alts;
    }
}

```

(a)

(b)

Figure 1: Programs for generating (a) a schema for an entity type (b) a knowledge base of entities.

90 **Generating a schema for an entity type** – The program in Figure 1a generates a schema for an
 91 entity type consisting of a set of properties, each with multiple names and a type. Names are drawn
 92 from a hand-constructed prior over property names. Types are drawn from a mixture over all built-in
 93 type priors. Types can have parameters, for example, a *Set* type has a parameter for the element
 94 type of the set – the type prior for each type defines a joint prior over these parameters, as described
 95 in Section 3. The core idea is that entity types are compound types learned from data, constructed
 96 from built-in primitive types whose parameters are also learned from data. These built-in types are
 97 the main way that prior knowledge is made available to the model.

98 **Generating a probabilistic knowledge base** – The program in Figure 1b generates an Alexandria
 99 knowledge base, consisting of a number of typed entities. Each typed entity has values for each
 100 property of the type – for example, an entity of type ‘person’ will have a value for the ‘DateOfBirth’
 101 property. To allow for disagreement over the value of a property, an entity can have many alternative
 102 property values – for example, multiple dates of birth for a person entity where there is disagreement
 103 about exactly when they were born. In the probabilistic program, first `entityCount` entities are
 104 created. For each entity and property, the number of alternative values is drawn from a (1-based)
 105 geometric distribution with parameter `probAlt`. Each alternative values is then drawn from the
 106 prior for that property type.

107 Importantly, the alternative values for a property represent different possible *conflicting* values, only
 108 one of which could actually be true. Such alternative values are very different to sets of values, where
 109 many values can be true for the same person (such a person who is a singer *and* a songwriter). For
 110 this reason, we model sets very differently to alternatives, using a specific set type (see Section 3).

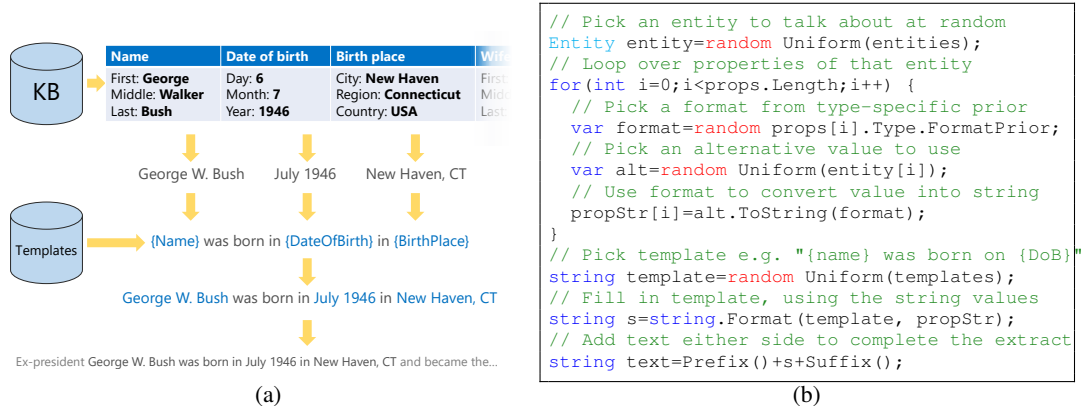


Figure 2: (a) Summary of the process of converting a typed entity from the KB into natural text containing facts about the entity. (b) Probabilistic program defining this process.

Generating text from KB values – Figure 2a summarises the process of converting typed values in the knowledge base into unstructured text describing those values. First, an entity to describe is selected at random from the knowledge base. Then each of the entity’s property values is converted into a string value using a format drawn from the type-specific format prior. For example, the date 6-July-1946 might get converted into the string “July 1946” using the format “MMMM yyyy”. The next task is to embed these string property values into a natural sentence, or part of one. This embedding is achieved using a template, such as “{name} was born on {date_of_birth} in {place_of_birth}” (a selection of templates learned in our experiments is given in the supplementary material). The template is selected at random from a provided set of templates. These templates are drawn from a suitable prior over strings that enforces that braces occur at word boundaries. The template is filled in by replacing each property placeholder with the corresponding value string. Finally, suitable prefix and suffix strings are concatenated, allowing the completed template to appear inside a larger section of text. It is this final variable that will have observed values coming from the web and will allow the system to learn from web text. The probabilistic program defining this process is given in Figure 2b.

The use of a probabilistic program to represent our knowledge base provides an elegant way to handle the uncertainty inherent in natural text. When we come to invert the probabilistic program and infer the values for a property, these values will be represented as probability distributions. As an example, if the text on the web page is “1972”, then the extracted distribution will be uniform across all dates in 1972. Another example is if the text is “scientist” then the value will be a distribution over all professions which are scientists. Preserving this uncertainty is essential for maintaining very high precision and for correctly merging information from multiple sources.

So far we have described the three sections of the core Alexandria probabilistic program. We made a number of extensions to this core program to handle lists of properties, model multiple texts occurring within a single page and to allow for noisy mentions of a property value – for reasons of space, these extensions are described in the supplementary material.

3 The Built-in Property Types

The Alexandria probabilistic program defines an entity type as a set of properties, where each property has one of the built-in property types. The built-in types used in this paper are `Date`, `PersonName`, `Place`, `Hierarchy` and `Quantity`. There is also a set type `Set<T>` whose element type `T` can have any of the above types – for example, a set of people’s names is written as `Set<PersonName>`. Finally, we have a ‘catch-all’ type that can generate any string value – this is used to discover properties which are not compatible with any of the built-in types (see Section 5).

These characteristics of each type are summarised in Table 1 – full details of each type are given in the supplementary material. Each built-in type has a set of type parameters (such as the element type of a `Set`). A manually-specified `TypePrior` defines a distribution over instances of the type, including all type parameters. Every type defines a prior over values of the type (`Prior`), which may depend on the type parameters and so can be learned from data. In addition, each type has a `ToString(value, format)` method, implemented as a probabilistic program, which converts (uncertain) values of the type into strings – these programs are listed in the supplementary material. Finally, each type has a prior over format strings (`FormatPrior`) which defines a distribution over the format that can be passed in to the `ToString(value, format)` method.

<i>Type</i>	<i>Type parameters</i>	<i>Prior</i>	<i>Example format string</i>
<code>Date</code>	-	Uniform over all dates between 1000 and 2100	“{dd} {MMMM} {yyyy}”
<code>PersonName</code>	-	Uniform over valid names	“{First} {M}. {Last}”
<code>Hierarchy</code>	A hierarchy	Learned node probs	“Lower3”
<code>Place</code>	-	Learned place prob	“Lower3”
<code>Quantity</code>	Quantity kind (e.g. Length), Prior mean and variance	Learned Gaussian or log Gaussian	“{feet:F0}’{sub_inch:F0}”
<code>Set<T></code>	Element type (T), Poisson, Beta parameters	Poisson over set size, Beta over mention prob.	“{0}, {1} and {2}”
‘catch-all’	-	Uniform over all strings	-

Table 1: Summary of the built-in property types.

152 The `Hierarchy` type takes a hierarchy as its type parameter. In this paper, we assume that a set
153 of hierarchies have been manually provided and the type prior is a uniform distribution over these
154 known hierarchies. For the experiments in Section 5, hierarchies were provided for occupations,
155 nationalities, star signs, genders, religions, causes of death and hair & eye colors. In future we intend
156 to learn these hierarchies by using a type prior over the structure and nodes of a hierarchy.

157 4 Distributed Approximate Inference

158 Infer.NET allows probabilistic programs to be compiled into efficient C# code for inferring posterior
159 distributions of specific variables in the program. This generated code applies a standard inference
160 algorithm, such as expectation propagation (EP) [11], using a library of probability distributions and
161 operations known as the Infer.NET runtime. For Alexandria, small examples of the core program can
162 be compiled in this way and executed locally on a single machine. Inference over string values is
163 handled using weighted automata, as described in [16].

164 We need to run three kinds of large scale inference queries on the probabilistic program:

165 **Template learning** – where the set of `templates` is inferred, given web texts (`text`), the schema
166 properties (`props`) and a set of `entities` where some of the property values are known;

167 **Schema learning** – where the set of schema properties (`props`) are inferred, given the web texts
168 (`text`), a set of minimal `entities` (each with just a name and one other property) and a corre-
169 sponding set of `templates`;

170 **Fact retrieval** – where the set of `entities` is inferred, given the entity names, the web texts
171 (`text`), the schema properties (`props`) and a set of `templates`.

172 In essence, each of these queries infers one of the three variables (`props`, `templates`, `entities`)
173 given the other two (and the web text). Fixing two out of these three is not essential, but helps to
174 keep the inference process efficient and the schedule straightforward to parallelize.

175 We apply these queries to billions of documents and millions of entities. To achieve this scale, it was
176 necessary to develop a distributed and optimised version of the inference algorithm. This distributed
177 algorithm also uses the Infer.NET runtime but differs from the automatically generated algorithm in
178 several important ways. Most importantly, it is written partially in SCOPE [2] allowing for large scale
179 execution on Microsoft’s Cosmos distributed computing platform. In addition, we use a manually
180 defined, message-passing schedule optimised for rapid convergence.

181 Each built-in property type used a corresponding Infer.NET distribution type to represent uncertain
182 values of the type and as messages in the message passing process. Values of these types can be
183 highly structured (such as strings, objects, hierarchies and sets) and so distributions were chosen to
184 be similarly structured. For example, distributions over entire objects and over entire sets were used.
185 Distributions over hierarchical values were designed to exploit the structure of the hierarchy.

186 To gain significant speed ups, we applied additional approximations to the EP messages. These
187 involved collapsing or removing uncertainty, where this could be done without adversely affecting
188 precision. We found that it was essential to preserve uncertainty in which templates was matched
189 (`template`), in the part of the matched text corresponding to a property (`propStr`) and in the
190 extracted value (`alt`). Conversely, we found that we could collapse uncertainty in the `entity` being
191 referred to in each piece of text, provided this was done conservatively and repeated at each iteration
192 of inference. Each text was assigned to the most probable existing entity, provided the probability
193 of this was above some high *mergeThreshold*, or otherwise to a new entity. Further speed ups were
194 achieved by caching and re-using the results of certain slow message passing operations.

195 To cope with billions of documents, we first use forward inference in the program to compute an
196 uncertain posterior over the `text` variable which defines all possible texts that could be generated
197 by the program for a particular query. This posterior distribution is converted into a search which
198 is executed at scale across the document corpus. The results of these searches are then returned
199 as observations of the `text` variable. We apply these observations for reverse inference in the
200 probabilistic program using the distributed inference algorithm. An overview of the inference pipeline
201 is given in the supplementary material.

5 Experiments and Results

We aim to demonstrate that Alexandria can infer both the schema for a type and high precision property values for entities of that type, from web scale text data. This inference is entirely unsupervised except for one single labelled example, used to bootstrap the system. The labelled example consists of the name “Barack Obama” and his date of birth “4 August 1961” which acts as the single textual label. Alexandria was also given a set of 3000 names used for schema learning and a separate test set of 4000 people’s names to retrieve property values for. These test set names were selected at random from the top 1 million most common names searched for on Bing. Alexandria also has access to around 8 billion English HTML pages cached from the public web. To create observed values of the `text` variable, these pages are processed to remove HTML tags and non-visible text (giving about 50TB of text). Some tags, such as bullets or line breaks, are converted into suitable text characters, to allow the model to make use of the layout of the text on the page. In addition, sections of text which are exact duplicates are assumed to be copies and treated as a single observation, rather than multiple independent observations.

To bootstrap the system from a single example, we ran a small-scale query that infers both the templates and the schema properties (`props`), with `entities` observed to a set containing only our single known entity (Barack Obama) with property values for his name and date of birth. This bootstrapping process outputs a small set of 2-property templates and a corresponding schema. The process of inferring a full schema and set of property values then used the large scale inference queries described in Section 4, as follows:

1. Run **fact retrieval** to retrieve date of birth for 10 names from the schema learning dataset;
2. Run **template learning** given these 10 entities and the 2-property schema;
3. Repeat **fact retrieval/template learning** for 100 and then all 3000 names;
4. Run **schema learning** given the 3000 entities and templates learned from these;
5. Run **template learning** for the top 20 properties in the new schema using these 3000 entities;
6. Finally, run **fact retrieval** to retrieve values for the 4000 names in the separate test dataset.

These queries were run with `probAlt` set to 0.005, the `mergeThreshold` set to 0.99 and `entityCount` set to 10 million. For the 4000 person fact retrieval, retrieving relevant texts from 8bn documents took around 25K compute hours. Running distributed inference on these retrieved texts took a further 900 compute hours (about 2-3 hours of clock time). The largest fact retrieval that we have tried was for 2 million people. In this case, the text retrieval time increased sub-linearly to 250K compute hours, whilst the inference time increased almost linearly to around 200K hours (the reduction in time-per-entity is due to the reduced data for rarer entities).

Results of schema learning – Table 2 shows the top properties discovered during schema learning. The first column in Table 2 shows the most common inferred name for each property (other inferred names are shown in the final column). The second column shows the most probable inferred type, out of the built-in types from Section 3. The third column gives the number of entities where property values were found during schema learning. The rows of the table have been ordered by this count, which was also used to select the 20 properties shown. This ordering focuses on properties that most people have, rather than properties present only for certain kinds of people, such as musicians or sportspeople. The fourth column gives the number of web domains that referred to the property at least once. This is an indication of how generally relevant the property is. For this table, we have excluded niche properties which are referred to by fewer than 20 domains to focus on generally relevant properties.

Using the catch-all type, the schema learning process discovers several properties whose values are not compatible with any of the built-in types. Such properties include descriptive text properties (‘best known for’, ‘biography’, ‘quotations’), times and time ranges (‘birth time’, ‘years active’), and a few rarer types. Mostly these would be straightforward to add as new built-in types.

Evaluating fact retrieval – Each retrieved value distribution needs to be evaluated against a ground truth value in the form of a string. Such evaluation is itself a challenging problem – for example, evaluating “scientist” against “physicist”, “184cm” vs “6ft1½”, “4/5/1980” against “May 4th 1980” and so on. To perform evaluation we re-use the probabilistic program and treat the ground truth text as an observed string value of the property. Inferring the property value gives a probability distribution of the same type as the retrieved value distributions. We then say that the value is correct if the expected probability of the ground truth value is higher under the retrieved value distribution

<i>Inferred name</i>	<i>Inferred type</i>	<i>Entities</i>	<i>Domains</i>	<i>Other inferred names</i>
name	PersonName	2,964	4,545	birth name, real name, birthname
born	Date	2,756	3,471	date of birth, birthday, ...
birthplace	Place	2,583	1,594	place of birth, birth place, ...
occupation	Set<Hierarchy(Occupations)>	2,569	801	profession, occupations, ...
nationality	Set<Hierarchy(Nationalities)>	2,485	505	citizenship
zodiac sign	Hierarchy(StarSigns)	2,336	328	sign, star sign, zodiac sign, ...
gender	Hierarchy(Genders)	2,110	247	sex
spouse	Set<PersonName>	2,058	665	spouse(s), wife, husband, ...
hair color	Set<Hierarchy(HairColors)>	2,050	340	hair, natural hair colour
height	Quantity(Length)	2,019	1,060	
age	Quantity(Time)	1,807	62	
eye color	Set<Hierarchy(EyeColors)>	1,705	405	eyes, eye colour
parents	Set<PersonName>	1,678	513	father, mother, father name, ...
died	Date	1,671	700	date of death, death, death date
religion	Set<Hierarchy(Religions)>	1,276	231	
siblings	Set<PersonName>	1,235	274	brother, sister
children	Set<PersonName>	1,121	368	
weight	Quantity(Weight)	594	325	
cause of death	Hierarchy(CausesOfDeath)	544	65	
place of death	Place	450	30	location of death, death place

Table 2: Results of schema learning for people. The top 20 discovered properties are shown, ordered by the number of entities with retrieved values for the property. See text for details of each column.

than under the prior for that property type. So for a retrieved distribution r , a ground truth distribution g and a prior p , the prediction is correct if $\sum_i g_i r_i > \sum_i g_i p_i$. where the sum is across all possible values of the property. For example, if the ground truth is “1975” and the prediction is “5 May 1975” then the left hand side of the equation equals $1/365$. Assuming a prior which is uniform over all dates between 1000 and 2500, the right hand side approximately equals $1/(365.25 \times 1500)$. Thus, the LHS is greater than the RHS and the prediction is considered correct. If the prediction had been “5 May 1976”, the LHS would equal zero and so the prediction would be considered incorrect. We aim to extract values at very high precision, around 97-99%. Evaluating at such high precision raises some specific challenges which we discuss in the supplementary material.

For each name in the test set, Alexandria outputs a discovered set of entities with that name. For common names, this set can be quite large! We only wish to evaluate the one that corresponds to the entity in the test set, or none at all if that entity was not discovered. To do this, we again re-use the probabilistic program and treat the ground truth as a new observation, given the set of discovered entities. If the inferred most probable `entity` for the observation is one of the discovered entities, then that entity is used for evaluation. If the inferred `entity` is a new entity, then we record that no prediction was made for that ground truth entity, that is, it was not discovered by Alexandria.

Results of fact retrieval – We evaluate our system against ground truth data from Microsoft’s knowledge graph (called Satori). Table 3 gives the evaluation metrics for fact retrieval for the

<i>Property</i>	<i>Prec@1</i>	<i>Prec@2</i>	<i>Recall</i>	<i>Alts</i>	<i>Property</i>	<i>Prec@1</i>	<i>Prec@2</i>	<i>Recall</i>	<i>Alts</i>
born	98.2%	99.4%	95.3%	1.38	eye				
birthplace	96.6%	97.4%	76.4%	1.09	color	94.7%	94.7%	88.4%	1.00
occupation	97.1%	97.3%	79.4%	1.19	parents	98.1%	98.1%	28.2%	1.03
nationality	98.2%	98.2%	83.3%	1.01	religion	97.6%	97.6%	57.5%	1.00
star sign	96.6%	97.3%	28.3%	1.08	siblings	100.0%	100.0%	16.4%	1.01
gender	99.6%	99.6%	38.9%	1.00	children	94.3%	94.3%	17.2%	1.05
height	98.6%	99.2%	79.7%	1.10	weight	97.1%	98.0%	68.2%	1.04
hair color	94.4%	94.4%	87.8%	1.00	cause				
spouse	95.5%	95.5%	44.5%	1.03	of death	98.4%	98.4%	63.4%	1.04
age	98.0%	98.0%	37.7%	1.03	place				
died	98.6%	99.2%	95.2%	1.32	of death	97.9%	97.9%	47.7%	1.00

Table 3: Fact retrieval metrics for the discovered properties.

discovered properties from Table 2. The fact retrieval process can result in more than one alternative value distribution for a particular name and property. For example, a date of birth property may have two retrieved alternatives, such as 5 May 1976 and 5 May 1977. Alternatives are ordered by the number of web pages that support them (which can be in the hundreds or thousands). The metric ‘precision@1’ refers to the percentage of retrieved values where the first alternative was correct, that is, the alternative with the most supporting web pages. The metric ‘precision@2’ gives the percentage of retrieved values where the first or second alternatives were evaluated as correct. Recall is defined as the percentage of entities with ground truth values where a prediction was made. We also report the average number of alternatives (‘Alts’), again where a prediction was made.

In considering Table 3, remember that these results were achieved without using *any* ground truth values, apart from a single date of birth value. Overall, the precisions of the first alternative (Prec@1) are high, with 13 of the 19 are in the 97%+ range and 9 of these above 98%. The lowest precisions are for children and hair and eye color properties, although these are still above 94%. Looking at the average number of alternatives, we see that most properties have just one alternative for the majority of predictions, so the precision@2 is the same as the precision@1. Exceptions to this include date of birth and date of death, which have 30-40% predictions with two alternatives (e.g. due to off-by-one errors). Considering the second alternative increases the precision of these properties by about 1%, bringing them to over 99% precision. We manually checked the causes of a sample of errors – these are ranked and discussed in the supplementary material. The recall of the properties varies widely, from 16.4% for siblings through to 95.3% for date of birth, with an average of 59.7%. Some of this variation is due to the maximum possible recall varying, that is, the variation in the fraction of values actually available on the web. Another factor is how well a set of templates can capture how such values are expressed in text. For example, dates of birth and death are often expressed in a standard form, whereas there is more variation in how siblings and children are described.

Comparison to previous results – most recent systems do not report results on any standard tasks, but instead use individual custom evaluation tasks and metrics, making rigorous comparison difficult. For example, YAGO2 [8] reports precisions in the range 92-98%, NELL [12] reports precisions in the range 80-85%, KnowledgeVault [5] reports a manually-assessed AUC of around 87%. The TAC-KBP Slot Filling track [14] is a competition which aims to provide a comparison for fact retrieval systems – the highest precision for TAC2014 was 59% from the DeepDive system [17]. Informal comparison to these reported results suggests that Alexandria’s accuracy is at least as good as existing supervised systems and superior to that of existing unsupervised systems.

6 Conclusions & future work

In this paper, we have shown how Alexandria can perform schema learning and high-precision fact extraction unsupervised, except for a single seed example. Whilst the results in this paper are for people entities, the system has been designed to be generally applicable to many types of entity. It’s worth noting that, in the process of learning about people, we have learned seed examples for other classes such as places, which we could use to do schema learning and fact extraction for these classes. By repeating this process recursively, we create the exciting possibility of using Alexandria like an Open IE system, to learn schemas, discovery entities and extract facts automatically across a large number of domains – this is our focus for future work. Our hope is that our high accuracy and strong typing will prevent ‘drift’ from occurring, which has reduced accuracy in previous Open IE systems.

The Alexandria model does not use any joint prior across property values – such as the graph prior and tensor factorization priors used in [5]. Incorporating such priors into the model has the potential to increase precision yet further.

Alexandria’s template-based language model is relatively simple compared to some NLP systems used in related work. In contrast, Alexandria’s model of types and values is in general more sophisticated, particularly in its handling and propagation of uncertainty. We believe that this has allowed the system to achieve very high precision. We expected to need a more sophisticated language model to achieve high recall – in fact, the ability to process the entire web means that we can still achieve good recall – a fact expressed in a complex way on one page is often expressed simply elsewhere.

We believe that Alexandria makes a step towards the holy grail of completely automatic KB construction and maintenance – we look forward to trying out the system in many new domains to see if the successful unsupervised learning of people can be replicated for other entity types.

References

- [1] A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. Hruschka, and T. M. Mitchell. Toward an Architecture for Never-Ending Language Learning. In *AAAI*, 2010.
- [2] R. Chaiken, B. Jenkins, P.-Å. Larson, B. Ramsey, D. Shakib, S. Weaver, and J. Zhou. SCOPE: easy and efficient parallel processing of massive data sets. *Proceedings of the VLDB Endowment*, 1(2):1265–1276, 2008.
- [3] Daniel Roy et al. PROBABILISTIC-PROGRAMMING.org, 2017. <http://probabilistic-programming.org>.
- [4] Defense Advanced Research Projects Agency. Probabilistic programming for advancing machine learning, 2013. DARPA-BAA-13-31.
- [5] X. L. Dong, E. Gabrilovich, G. Heitz, W. Horn, N. Lao, K. Murphy, T. Strohmman, S. Sun, and W. Zhang. Knowledge Vault: A Web-Scale Approach to Probabilistic Knowledge Fusion. In *KDD*, 2014.
- [6] A. Fader, S. Soderland, and O. Etzioni. Identifying relations for open information extraction. In *EMNLP*, 2011.
- [7] R. Gupta, A. Halevy, X. Wang, S. Whang, and F. Wu. Biperpedia: An Ontology for Search Applications. In *VLDB*, 2014.
- [8] J. Hoffart, F. M. Suchanek, K. Berberich, and G. Weikum. YAGO2: A spatially and temporally enhanced knowledge base from Wikipedia. *Artificial Intelligence*, 194:28 – 61, 2013.
- [9] Mausam, M. Schmitz, R. Bart, S. Soderland, and O. Etzioni. Open Language Learning for Information Extraction. In *EMNLP*, pages 523–534, 2012.
- [10] T. Minka, J. Winn, J. Guiver, S. Webster, Y. Zaykov, B. Yangel, A. Spengler, and J. Bronskill. Infer.NET 2.6, 2014. Microsoft Research Cambridge. <http://research.microsoft.com/infernet>.
- [11] T. P. Minka. Expectation propagation for approximate Bayesian inference. In *Proceedings of the 17th conference on Uncertainty in Artificial Intelligence*, pages 362–369, 2001.
- [12] T. Mitchell, W. Cohen, E. Hruschka, P. Talukdar, J. Betteridge, A. Carlson, B. Dalvi, M. Gardner, B. Kisiel, J. Krishnamurthy, N. Lao, K. Mazaitis, T. Mohamed, N. Nakashole, E. Platanios, A. Ritter, M. Samadi, B. Settles, R. Wang, D. Wijaya, A. Gupta, X. Chen, A. Saparov, M. Greaves, and J. Welling. Never-Ending Learning. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence (AAAI-15)*, 2015.
- [13] M. Richardson and P. Domingos. Markov Logic Networks. *Machine Learning*, 62(1):107–136, 2006.
- [14] M. Surdeanu and H. Ji. Overview of the English slot filling track at the TAC2014 knowledge base population evaluation. In *Proceedings of the 7th TAC*, 2014.
- [15] G. Weikum and M. Theobald. From Information to Knowledge: Harvesting Entities and Relationships from Web Sources. In *Proceedings of the 29th Symposium on Principles of Database Systems*, pages 65 – 76, 2010.
- [16] B. Yangel, T. Minka, and J. Winn. Belief Propagation with Strings. Technical Report MSR-TR-2017-11, Microsoft Research, February 2017.
- [17] C. Zhang. *DeepDive: A Data Management System for Automatic Knowledge Base Construction*. PhD thesis, University of Wisconsin-Madison, 2015.