

The *Ex* Project: Web Information Extraction using Extraction Ontologies

Martin Labský, Vojtěch Svátek, Marek Nekvasil, and Dušan Rak

Department of Information and Knowledge Engineering,
University of Economics, Prague, W. Churchill Sq. 4, 130 67 Praha 3, Czech Republic
e-mail: {labsky, svatek, nekvasim, rakdusan}@vse.cz

Abstract. Extraction ontologies represent a novel paradigm in web information extraction (as one of ‘deductive’ species of web mining) allowing to swiftly proceed from initial domain modelling to running a functional prototype, without the necessity of collecting and labelling large amounts of training examples. Bottlenecks in this approach are however the tedium of developing an extraction ontology adequately covering the semantic scope of web data to be processed and the difficulty of combining the ontology-based approach with inductive or wrapper-based approaches. We report on an ongoing project aiming at developing a web information extraction tool based on richly-structured extraction ontologies and with additional possibility of (1) semi-automatically constructing these from third-party domain ontologies, (2) absorbing the results of inductive learning for subtasks where pre-labelled data abound, and (3) actively exploiting formatting regularities in the wrapper style.

1 Introduction

Web information extraction (WIE) represents a specific category of web mining. It consists in the identification of typically small pieces of relevant text within web pages and their aggregation into larger structures such as data records or instances of ontology classes. As its core task is application of pre-existent patterns or models (in contrast to inductively discovering new patterns), it falls under the notion of ‘deductive’ web mining [16], similarly as e.g. web document classification. As such, some kind of prior knowledge is indispensable in WIE. However, the ‘deductive’ aspects of WIE are often complemented with inductive ones, especially in terms of learning the patterns/models (at least partly) from training data.

In the last decade, WIE was actually dominated by two paradigms. One—*wrapper-based*—consists in systematically exploiting the surface structure of HTML code, assuming the presence of regular structures that can be used as anchors for the extraction. This approach is now widely adopted in industry, however, its dependence on formatting regularity limits its use for diverse categories of web pages. The other—*inductive*—paradigm assumes the presence of training data: either web pages containing pre-annotated tokens or stand-alone examples of data instances. It is linked to exploration of various computational learning paradigms, e.g. Hidden-Markov Models, Maximum Entropy Models, Conditional Random Fields [10] or symbolic approaches

such as rule learning [1]. Again, however, the presence of (sufficient amounts of) annotated training data is a pre-condition that is rarely fulfilled in real-world settings, and manual labelling of training data is often unfeasible; statistical bootstrapping alleviates this problem to some degree but at the same time burdens the whole process with ‘heavy computational machinery’, whose requirements and side-effects are not transparent to a casual user of a WIE tool. In addition, both approaches usually deliver extracted information as rather weakly semantically structured; if WIE is to be used to fuel semantic web repositories, secondary mapping to ontologies is typically needed, which makes the process complicated and possibly error-prone.

There were recently proposals for pushing ontologies towards the actual extraction process as immediate prior knowledge.¹ *Extraction ontologies*² [3] define the concepts, the instances of which are to be extracted, in the sense of various attributes, their allowed values as well as higher level (e.g. cardinality or mutual dependency) constraints. Extraction ontologies are assumed to be hand-crafted based on observation of a sample of resources; however, due to their clean and rich conceptual structure (allowing partial intra-domain reuse and providing immediate semantics to extracted data), they are superior to ad-hoc hand-crafted patterns used in early times of WIE. At the same time, they allow for rapid start of the actual extraction process, as even a very simple extraction ontology (designed by a competent person) is likely to cover a sensible part of target data and generate meaningful feedback for its own redesign; several iterations are of course needed to obtain results in sufficient quality. It seems that for web domains that consist of a high number of relatively tiny and evolving resources (such as web product catalogs), information extraction ontologies are the first choice. However, to make maximal use of available data and knowledge and avoid overfitting to a few data resources examined by the designer, the whole process must not neglect available labelled data, formatting regularities and even pre-existing domain ontologies.

In this paper we report on an ongoing effort in building a WIE tool named *Ex*, which would synergistically exploit all the mentioned resources, with central role of extraction ontologies. The paper is structured as follows. Section 2 explains the structure of extraction ontologies used in *Ex*. Section 3 describes the steps of the information extraction process using extraction ontologies and other resources. Section 4 briefly reports on experiments in three different domains. Finally, section 5 surveys related research, and section 6 outlines future work.

2 Ex(traction) ontology content

Extraction ontologies in *Ex* are designed so as to extract occurrences of *attributes* (such as ‘age’ or ‘surname’), i.e. standalone named entities or values, and occurrences of whole *instances of classes* (such as ‘person’), as groups of attributes that ‘belong together’, from HTML pages (or texts in general) in a domain of interest.

¹ See [8] for general discussion of the types and roles of ontologies in the WIE process.

² In earlier work [7] we used the term *presentation ontology*.

2.1 Attribute-related information

Mandatory information to be specified for each attribute is: name, data type (string, long text, integer, float) and dimensionality (e.g. 2 for screen resolution like 800x600). In order to automatically extract an attribute, additional knowledge is typically needed. Extraction knowledge about the attribute *content* includes (1) textual value patterns; (2) for integer and float types: min/max values, a numeric value distribution and possibly units of measure; (3) value length in tokens: min/max length constraints or a length distribution; (4) axioms expressing more complex constraints on the value and (5) coreference resolution knowledge. Attribute *context* knowledge includes (1) textual context patterns and (2) formatting constraints.

Patterns in *Ex* (for both the value and the context of an attribute or class) are nested regular patterns defined at the level of tokens (words), characters, formatting tags (HTML) and labels provided by external tools. Patterns may be inlined in the extraction ontology or sourced from (possibly large) external files, and may include the following:

- specific tokens, e.g. 'employed by'
- token *wildcards*, which require one or more token properties to have certain values (e.g. any capital or uppercase token, any token whose lemma is 'employ')
- *character-level* regular expressions for individual tokens
- *formatting tags* or their classes, such as 'any HTML block element'
- *labels* that represent the output of external tools, such as sentence boundaries, part-of-speech tags, parsed chunks or output from other IE engines.
- *references* to other *matched patterns*; this allows for construction of complex grammars where rules can be structured and reused
- *references* to other *matched attribute candidates*: a *value* pattern containing a reference to another attribute means that it can be nested inside this attribute's value; for *context* patterns, attribute references help encode how attributes follow each other

For *numeric* types, default value patterns for integer/float numbers are provided. Tabular, uniform, normal and mixture distributions are available to model attribute values. Linking a numeric attribute to unit definitions (e.g. to various currency units) will automatically create value patterns containing the numeric value surrounded by the units. In case there are multiple convertible units the extraction knowledge is reused.

For both attribute and class definitions, *axioms* can be specified that impose constraints on attribute value(s). For a single attribute, the axiom checks the to-be-extracted value and is either satisfied or not (which may boost or suppress the attribute candidate's score). For a class, each axiom may refer to all attribute values present in the partially or fully parsed instance. For example, a price with tax must be greater than the price without tax. Axioms can be authored using the JavaScript³ scripting language. We chose JavaScript since it allows arbitrarily complex axioms to be constructed and also because the web community is used to it.

In addition, *formatting constraints* may be provided for each attribute. Currently, four types of formatting constraints are supported: (1) the whole attribute value is contained in a single parent, i.e. it does not include other tags or their boundaries; (2) the

³ <http://www.mozilla.org/rhino>

value fits into the parent; (3) the value does not cross any inline formatting elements; (4) it does not cross any block elements. We investigate how custom constraints could easily be added by users. By default, all four constraints are in effect and influence the likelihood of attribute candidates being extracted.

In many tasks, it is necessary to identify co-referring occurrences of attribute values. For each attribute, a coreference resolution script can be provided that determines whether two values of the same attribute (or of its specialized attributes) may co-refer to the same entity. By default, identical extracted values for the same attribute are treated as co-referring.

2.2 Class-related information

Each *class definition* enumerates the attributes which may belong to it, and for each attribute it defines a *cardinality* range and optionally a cardinality distribution. Extraction knowledge may address both content and context of the class. *Class content patterns* are analogous to the attribute value patterns, however, they may match *parts* of an instance and must contain at least one *reference* to a member attribute. Class content patterns may be used e.g. to describe common wordings used between attributes or just to specify attribute ordering. *Class context patterns* are analogous to attribute context patterns.

Axioms are used to constrain or boost instances based on whether their attributes satisfy the axiom. For each attribute, an *engagedness* parameter may be specified to estimate the apriori probability of the attribute joining a class instance (as opposed to standalone occurrence). Regarding class context, analogous *class context patterns* and similar *formatting constraints* as for attributes are in effect also for classes.

In addition, constraints can be specified that hold over the whole sequence of extracted objects. Currently supported are minimal and maximal instance counts to be extracted from a document for each class.

2.3 Extraction evidence parameters

All types of extraction knowledge mentioned above, i.e. value and context patterns, axioms, formatting constraints and ranges or distributions for numeric attribute values and for attribute content lengths, are essentially pieces of evidence indicating the presence (or absence) of a certain attribute or class instance. In *Ex*, every piece of evidence may be equipped with two probability estimates: precision and recall. The *precision* of evidence states how probable it is for the predicted attribute or class instance to occur given the evidence holds, disregarding the truth values of other evidence. For example, the precision of a left context pattern “person name: \$” (where \$ denotes the predicted attribute value) may be estimated as 0.8; i.e. in 80% of cases we expect a person name to follow in text after a match of the “person name:” string. The *recall* of evidence states how abundant the evidence is among the predicted objects, disregarding whether other evidence holds. For example, the “person name: \$” pattern could have a low recall since there are many other contexts in which a person name could occur.

Pattern precision and recall can be estimated in two ways. First, annotated documents can be used to estimate both parameters using simple ratios of counts observed in

text. In this case, it is necessary to smooth the parameters using an appropriate method. For a number of domains it is possible to find existing annotated data, e.g. web portals often make available online catalogs of manually populated product descriptions linking to the original sellers' web pages. When no training data is available or if the evidence seems easy to estimate, the user can specify both parameters manually. For the experimental results reported below we estimated parameters manually.

2.4 Complex examples

In order to illustrate most of the above features, we present and explain two relatively large real-world examples from two different domains.⁴

The first example is from the *product catalogue* domain. Fig. 1 shows the structure of an extraction ontology for the domain of computer monitor offers. Fig. 2 displays part of the corresponding code in the XML-based ontology definition language, dealing with the name of the monitor model, its price and derived prices with and without tax.

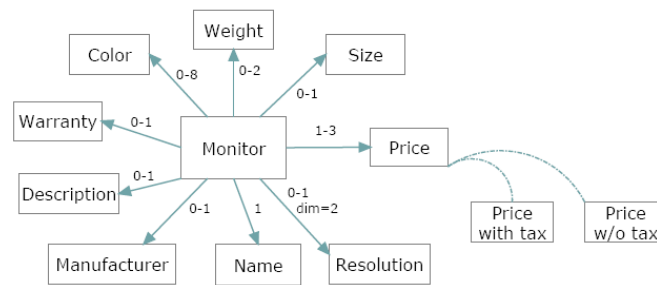


Fig. 1. General scheme of extraction ontology of monitors

In the global scope of the model, extraction knowledge affecting more than one attribute is defined. First, a pattern declares that in 75% of cases a *Monitor* instance starts with its *name* which is closely (within at most 20 tokens) followed by up to 4 *price* attributes. Second, an axiom says that in all cases the *price_with_tax* must be greater than *price_without_tax* (if both are specified).

The textual *name* attribute shows the usage of generic patterns (*model_id*) which can be reused within the following value and context sections. The first pattern in the value section describes a typical form of a computer monitor name and claims that 80% of its matches identify a positive example of monitor name. It also claims that 75% of monitor names exhibit this pattern. The value section further constrains the name's length and its position within HTML formatting elements; e.g. the last formatting pattern denies crossing of formatting block boundaries by saying that 100% of monitor names do not cross block tags. The context section's only pattern says that in 20% of

⁴ More details about experiments carried out in those domains are in Section 4.

```

<class id="Monitor">
<pattern id="name_first_and_price_follows" type="pattern" cover="0.75">
  ^ $name (<tok/>{0,20} $price){1,4} </pattern>
<axiom> $price_with_tax > $price_without_tax </axiom>
<attribute id="name" type="name" card="1" eng="0.70">
<pattern id="model_id">
<tok case="UC"/> | <tok type="ALPHA" case="CA"/> | <tok type="ALPHANUM|INT"/>
</pattern>
<value>
<pattern cover="0.5" p="0.8"> (LCD (monitor|panel)?)?
<pattern src="manuf.txt"/> <pattern ref="model_id"/>{1,2}
</pattern>
<length><distribution min="1" max="7"/></length>
<pattern cover="0.5" type="format"> has_one_parent </pattern>
<pattern cover="0.5" type="format"> fits_in_parent </pattern>
<pattern cover="0.8" type="format"> no_crossed_inline_tags </pattern>
<pattern cover="1.0" type="format"> no_crossed_block_tags </pattern>
</value>
<context>
<pattern p="0.3" cover="0.2"> model? name :? $ </pattern> </context>
</attribute>
<attribute id="price" type="float" card="1-4" units="euro, pound, dollar" eng="0.80">
<value>
<pattern cover="0.9" p="0.9">
<pattern ref="generic.price"/> $unit | $unit <pattern ref="generic.price"/> </pattern>
<distribution min="100" max="5000" />
<transform> $.replace(/s+/, " ") </transform>
</value>
<context>
<pattern cover="0.05" p="0.5"> price of ? :? $ </pattern> </context>
</attribute>
<attribute id="price_with_tax" card="0-1" type="float" extends="price">
<pattern id="intro">
price? ( (with|including|incl.) (tax|taxes) | \?(? (tax|taxes) included \)? )
</pattern>
<pattern id="outro">
( (with|including|incl.) (tax|taxes) | \?(? (tax|taxes) included \)? | \? (price with (tax|taxes) \) )
</pattern>
<context>
<pattern cover="0.5" p="0.8"> $ <pattern ref="outro"/> </pattern>
<pattern cover="0.1" p="0.1"> <pattern ref="intro"/> :? $ </pattern>
</context>
</attribute>

```

Fig. 2. Fragment of code of extraction ontology for computer monitors

cases, monitor names are preceded by labels like 'model name' and that observing such labels identifies monitor names with a 30% precision only.

The *price* attribute and its two extensions are numeric. The first pattern of the value section utilizes a pattern 'generic.price' and currency units, both of which are imported from a generic datatypes model (not shown). The numeric value constraint uses by default the first unit (Euro). The value transformation is a script applied to the value after extraction. The *price_with_tax* attribute, and similarly, *price_without_tax* (not shown) inherit all extraction knowledge from *price* and specify additional context patterns. Observing these patterns will cause the price being extracted as one of the two extensions.

The second example is taken from the *contact information* ontology developed within the EU project MedIEQ;⁵; the goal of the project is to ease expert-based accred-

⁵ <http://www.medieq.org>

```

<class id="Contact">
<script src="contact.js" />
<pattern id="title_name_together_and_first" type="pattern" cover="0.7">
^ $title{0-3} .? $name ,? $title{0-3} </pattern>
<axiom cover="0.8"> nameMatchesEmail($name, $email) > 0 </axiom>
<classifier id="cls1" method="weka" classtype="attribute"
name="weka.classifiers.rules.JRip" features="contact_all_jrip.fea"
model="./ex/data/med/train/contact_all_jrip.bin" elements="*" />
<attribute id="degree" type="name" card="0-4" eng="0.60">
<pattern id="degree_pre"> ( Miss | Lady | Sir ) .? </pattern>
<pattern id="degree_suf"> ( MSc | MA | MPH ) .? </pattern>
<value>
<pattern cover="0.7" p="0.8" ignore="case">
<pattern ref="degree_pre" /> | <pattern ref="degree_suf" /> </pattern>
<pattern cover="1" type="format"> has_one_parent </pattern>
</value>
</attribute>
<attribute id="name" type="name" card="1" eng="0.80">
<pattern id="init"> (A|B|C|D|E|F|G|H|I|J|K|L|M|N|O|P|R|S|T|U|V|W|X|Y|Z) .? </pattern>
<value>
<pattern cover="0.5" p="0.8">
<pattern src="first.txt" type="pattern" ignore="lemma" case="CA|UC" />
<pattern ref="init" />?
<pattern src="last.txt" type="pattern" ignore="lemma" case="CA|UC" />
</pattern>
<pattern cover="0.25" p="0.4">
<pattern src="first.txt" type="pattern" ignore="lemma" case="CA|UC" />
<pattern ref="init" />?
<tok type="alpha" case="CA|UC" />
</pattern>
<pattern p="0.7" cover="0.5"> <phr label="cls1.name" /> </pattern>
<pattern cover="0.3" type="format"> fits_in_parent </pattern>
<length> <distribution min="1" max="6" /> </length>
<axiom> checkPersonName($) </axiom>
<refers> nameRefersTo($, $other) </refers>
</value>
<context> <pattern cover="0.05" p="0.6"> person? name :? $ </pattern> </context>
</attribute>
<attribute id="email" type="name" card="0-1" eng="0.60">
<value> <pattern cover="0.98" p="0.9"> <pattern ref="generic.email" /> </pattern>
</value>
</attribute>

```

Fig. 3. Fragment of code of extraction ontology for contact information

itation of medical websites by automatically extracting information that is critical for the evaluation. This includes, among other, *contact information* of website responsible. A simplified version of the code related to persons' degrees and emails is shown.

In the global scope of the model, we can see the extraction knowledge referring to more than one attribute. The first pattern states that in 70% of cases, a *Contact* starts with its name or degree, and that these are only separated by punctuation. The axiom claims that in 80% of cases, the person's name and email exhibit some string similarity which can be intercepted by the referenced *script function* *nameMatchesEmail()* which returns non-zero if it thinks the given name corresponds to the given email. This function has been imported from a *script* "contact.js" above. Finally, a *classifier link* contracts an external trained classifier to classify all attributes of the *Contact* class. Classifications made by this classifier can be used in all patterns of this ontology.

The *degree* attribute demonstrates how generic patterns (*degree_pre* and *degree_suf*) can be reused by other patterns. The value section declares that in 80% of cases, observing the specified pattern really identifies a degree. It also claims that 70% of all degree occurrences will match this pattern (much larger enumeration would be needed to satisfy this number in reality). The *formatting pattern* requires all degrees to be enclosed in a single formatting element.

The *name* attribute's value section draws extraction knowledge from several sources. Its first pattern uses large first name and surname lists and inserts an optional initial in between. It claims to be relatively precise: 80% of its matches should correctly identify a person name. However, it is only expected to cover about 40% of all person names as the lists can never be exhaustive. The second pattern is less precise as it allows any alphabetic surname which is either capital or uppercase. When this pattern is matched, the extraction decision will strongly depend on the observed values of other extraction evidence. The third pattern takes into account the classifier's positive decision and it expects the classifier to have a 70% precision and a 50% recall. In addition, the name length is constrained and an axiom can perform a final check of the to-be-extracted value using a script function. As there are often multiple occurrences of a single person name on a web page, the *refers* section uses a script function to determine whether two names may co-refer (e.g. John Smith to Smith). The context section describes the typical surrounding phrases. As the context pattern's precision is not enough to create new person names, it will boost person name candidates for which it is observed (thanks to its precision) and suppress a little those for which it is missing (thanks to its coverage).

For the *email* attribute, just a single pattern is used within its value section which refers to a generic pattern defined by an imported datatype extraction ontology.

3 The extraction process

The inputs to the extraction process are the extraction ontology and a set of documents. Extraction consists of six stages depicted in Fig. 4.

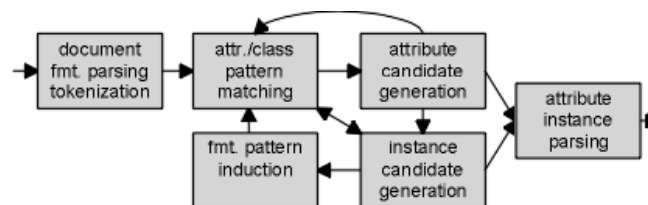


Fig. 4. Extraction process schema

3.1 Document preprocessing

First, the analysed document is loaded and its formatting structure is read into a simplified DOM tree of formatting objects. Currently *Ex* supports plain text, HTML and

XHTML documents. To robustly read web pages containing invalid HTML (most of today’s web pages) we employ the CyberNeko HTML parser,⁶ which fixes HTML syntax problems while parsing (adds missing end tags etc.), allowing for easier creation of a formatting element tree. Any text found in the formatting elements and their attributes is tokenized using a configurable tokenizer. A flat array of tokens is created for the document with each token linking to its parent formatting element. As part of tokenization, new words are registered in a common vocabulary, lemmatized and linked to their lemmas (if available), and classified by token type (e.g. alphanumeric) and case (e.g. capital). Sentence boundaries are detected⁷ and added as labels that can be exploited by extraction ontology patterns. Finally, any utilised external tools are executed to produce further labels, such as third-party named entity recognizers or a part-of-speech tagger. During this phase, *NLP tools* like part-of-speech taggers or chunk parsers could also be employed to create information useful for the definition of extraction patterns (future work).

3.2 Attribute candidate generation

After loading a document, all attribute value and attribute context patterns of the ontology are matched against the document’s tokens. Where a value pattern matches, the system attempts to create a new candidate for the associated attribute (attribute candidate – AC). If more value patterns match at the same place, or if there are context pattern matches for this attribute in neighbouring areas, then the corresponding evidence is turned on as well for the AC. Evidence corresponding to all other non-matched patterns is kept off for the AC. Also, during the creation of the AC, all other evidence types (axioms, formatting constraints, content length and numeric value ranges) are evaluated and set. The set of all evidence Φ_A known for the attribute A is used to compute a *conditional probability estimate* P_{AC} of how likely the AC is given all the observed evidence values:

$$P_{AC} = P(A|E \in \Phi_A) \quad (1)$$

The full formula is described and derived in [6]. We assume conditional independence of evidence given that the attribute holds or does not hold. The AC is created only if P_{AC} exceeds a pruning threshold defined by the extraction ontology.

In places where a context pattern matches and there are no value pattern matches in neighbourhood, the system tries to create ACs of various length (in tokens) in the area pointed to by the context pattern. Again, all the evidence values are combined to compute P_{AC} for each new AC. For patterns which include other attributes, we run the above process until no new ACs are generated.

The set of (possibly overlapping) ACs created during this phase is represented as an *AC lattice* going through the document, where each AC is scored by $score(AC) = \log(P_{AC})$. Apart from the ACs which may span multiple tokens, the lattice also includes one ‘background’ state for each token that takes part in some AC. A background state BG_w for token w is scored as follows:

⁶ <http://people.apache.org/~andyc/neko/doc/html/>

⁷ <http://www.comp.nus.edu.sg/~qiul/NLPTools/SentenceSplitter>

$$score(BG_w) = \min_{AC, w \in AC} \log\left(\frac{1 - P(AC)}{|AC|}\right) \quad (2)$$

where $|AC|$ is the length of the AC in tokens. The extraction process can terminate here if no instance generation or formatting pattern induction is done, in which case all ACs on the best path through the lattice are extracted.

3.3 Instance candidate generation

At the beginning of the instance candidate (IC) generation phase, each AC is used to create a simple IC consisting just of that single AC. Then, a bottom-up IC generation algorithm is employed to generate increasingly complex ICs from the working set of ICs. At each step, the highest scoring (seed) IC is chosen and its neighbourhood is searched for ACs that could be added to it without breaking ontological constraints for the IC class. Only a *subset* of the constraints is taken into account at this time as e.g. some minimum cardinality constraints or axioms requiring the presence of multiple attributes could never get satisfied initially. Each added AC is also examined to see whether it may co-refer with some AC that is already present in the IC; if yes, it is only added as a reference and it does not affect the resulting IC score.

After adding ACs to the chosen seed IC, that IC is removed from the working set and the newly created larger ICs are added to it. The seed IC is added to a *valid IC set* if it satisfies *all* ontological constraints. As more complex ICs are created by combining simpler ICs with surrounding ACs, a limited number of ACs or AC fragments is allowed to be skipped (AC_{skip}) between the combined components, leading to a penalization of the created IC. The IC scores are computed based on their AC content and based on the observed values of evidence E known for the IC class C :

$$sc_1(IC) = \exp\left(\frac{\sum_{AC \in IC} \log(P_{AC}) + \sum_{AC_{skip} \in IC} (1 - \log(P_{AC_{skip}}))}{|IC|}\right) \quad (3)$$

$$sc_2(IC) = P(C|E \in \Omega_C) \quad (4)$$

where $|IC|$ is the number of member ACs and Ω_C is the set of evidence known for class C ; the conditional probability is estimated as in Eq. 1. By experiment we chose the Prospector [2] pseudo-bayesian method to combine the above into the final IC score:

$$score(IC) = \frac{sc_1(IC)sc_2(IC)}{sc_1(IC)sc_2(IC) + (1 - sc_1(IC))(1 - sc_2(IC))} \quad (5)$$

The IC generation algorithm picks the best IC to expand using the highest $score(IC)$. The generation phase ends when the working set of ICs becomes empty or on some terminating condition such as after a certain number of iterations or after a time limit has elapsed. The output of this phase is the set of valid ICs.

As the generated IC space can get very large for some extraction ontologies, it can be constrained by several configurable pruning parameters: the IC probability can be thresholded so that only promising hypotheses are generated; the absolute beam size limits the number of most probable ICs to be kept for each span in the document; the relative beam size is used to prune all ICs within a span whose probability relative to the best IC in that span is less than the specified ratio.

3.4 Formatting pattern induction

During the IC generation process, it may happen that a significant part of the created valid ICs satisfies some (a priori unknown) *formatting pattern*. For example, a contact page may consist of 6 paragraphs where each paragraph starts with a bold person name together with scientific degrees. A more obvious example would be a table with the first two columns listing staff first names and surnames. Then, if e.g. 90 person names are identified in such table columns and the table has 100 rows, the induced patterns make the remaining 10 entries more likely to get extracted as well.

Based on the lattice of valid ICs, the following *pattern induction* procedure is performed. First, the best scoring sequence of non-overlapping ICs is found through the lattice. Only the ICs on the best path take part in pattern induction. For each IC, we find its nearest containing formatting block element. We then create a subtree of formatting (incl. inline) elements between the containing block element (inclusive) and the attributes comprising the IC. This subtree contains the names of the formatting elements (e.g. paragraph or bold text) and their order within parent (e.g. the first or second cell in table row). Relative frequencies of these subtrees are calculated over the examined IC set (separately for each class if there are more). If the relative and absolute frequencies of a certain subtree exceed respective configurable thresholds, a new formatting pattern is induced and the subtree is transformed into a new *context pattern* indicating the presence of the corresponding class. This induced formatting context pattern is an example of ‘local’ evidence only useful within the currently analysed document (or a set of documents coming from the same source). The precision and recall of the induced context patterns are based on the relative frequencies with which the patterns hold in the document (or document set) with respect to the observed ICs.

The newly created context patterns are then fed back to the pattern matching phase, where they are matched and applied. This extra iteration rescores existing ACs and ICs and may as well yield new ACs and ICs which would not have been created otherwise. With our current implementation we have so far only experimented with pattern induction for ICs composed of a single attribute. Using this feature typically increases recall but may have adverse impact on precision. One approach to avoid degradation of precision is to provide attribute evidence which will prevent unwanted attributes from being extracted.

3.5 Attribute and instance parsing

The purpose of this final phase is to output the most probable sequence of instances and standalone attributes through the analysed document. The valid ICs are merged into the AC lattice described in Section 3.2 so that each IC can be avoided by taking a competing path through standalone ACs or through background states. In the merged lattice, each IC is scored by $\log(\text{score}(IC))|IC|$.

The merged lattice is searched using *dynamic programming* for the most probable sequence of non-overlapping extractable objects. Path scores are computed by adding the state scores as they consist of log probability estimates. To support *n-best* results, up to *n* back pointers, ordered by their accumulated scores, are stored for each visited state. The *n* best paths are retrieved by backtracking using an *A**-based algorithm [4].

Extractable objects (attribute values and instances) can be read from each extracted path and sent to output.

In order to support constraints over the whole extracted sequence of objects (such as the minimal and maximal instance counts in a document), the search algorithm has been extended so that back pointers accumulated for each visited state include information about how the particular partial path leading into the visited state (inclusive) satisfies the constraint. In case of the instance count constraint, each back pointer of each visited state includes the instance counts observed on the corresponding incoming partial path. Paths that violate the constraints in a way that cannot be undone by further expansion of the path are not prolonged by the search. All paths that reach the final state of the document lattice and do not completely satisfy the constraints are discarded.⁸

3.6 Incorporating third party tools

In practical WIE tasks it often happens that some of the attributes of interest are relatively easy to extract using manually specified evidence, some require *machine learning* algorithms such as CRFs [10] in order to achieve good extraction results, and some may benefit from a combination of both. To support all three cases, *Ex* allows named entity candidates identified by other engines to be included in all types of textual patterns described above. For example, suppose our task is to extract instances of a Person class composed of a person name and a scientific degree. Let's also suppose we have training data for person names but no data for degrees. A viable approach would then be to train e.g. a CRF classifier to identify person names in text and to specify evidence for degrees manually. To incorporate the CRF classifier's suggestions into the extraction ontology, a simple *attribute value pattern* like "`#{crf:personname}`" can be added to the person name attribute. Here, `#{}` denotes a reference to an external named entity, *crf* is the source component name and *personname* is the identifier output by the CRF classifier. The precision for this value pattern can either be derived from the CRF classifier confidence score, or we can use the expected precision of the classifier for this attribute. To estimate the recall of the pattern, we can use the expected recall achieved by the classifier. Additionally to this pattern, the user may specify more patterns to correct (limit or extend) the classifier's suggestions.

4 Experimental Results

4.1 Contact Information on Medical Pages

In the EU (DG SANCO) MedIEQ project⁹ we experiment with several dozens of *medical website quality criteria*, most of which are to be evaluated with the assistance of IE tools. One of them is the presence and richness of contact information. Results are presented for 3 languages. A meta-search engine [15] was used to look up relevant medical

⁸ Since the extended search algorithm may increase memory requirements, it is only used when some constraints are specified by the extraction ontology.

⁹ <http://www.medieq.org>

web sites; these websites were then spidered, their contact pages were manually identified, contact information was annotated and the pages were added to contact page collections for each language. In total, 109 HTML pages were assembled for English with 7000 named entities, 200 for Spanish with 5000 and 108 for Czech with 11000 named entities. A contact extraction ontology was developed for each language, with language-independent parts reused. The extraction ontology developer was allowed to see 30 randomly chosen documents from each collection for reference and to use any available gazetteers such as city names, frequent first names and surnames. Results were evaluated using the remaining documents. On average, each ontology contained about 100 textual patterns for the context and content of attributes and of the single extracted 'contact' class, attribute length distributions, several axioms and co-reference resolution rules. The effort spent on developing and tuning the initial English ontology was about 2-3 person-weeks, and its customization to a new language amounted to 2 person-weeks.

Table 1 shows contact extraction results for 3 languages.¹⁰ Two evaluation modes are used: strict (the first number) and loose. In the strict mode of evaluation, only exact matches are considered to be successfully extracted. In the loose mode, partial credit is also given to incomplete or overflowed matches; e.g. extracting 'John Newman' where 'John Newman Jr.' was supposed to be extracted will count as a 66% match (based on overlapping word counts). Some of the performance numbers below may be impacted by relatively low inter-annotator agreement (English and Spanish datasets are being cleaned to remove inconsistencies). Fig. 5 shows a sample of automatically annotated data.

Table 1. Contact IE results

attribute	English			Spanish			Czech		
	prec	recall	F	prec	recall	F	prec	recall	F
title	71/78	82/86	76/82	-	-	-	87/89	88/91	88/90
name	66/74	51/56	58/64	71/77	81/86	76/81	74/76	82/83	78/80
street	62/85	52/67	56/75	71/93	46/58	56/71	78/83	66/69	71/75
city	47/48	73/76	57/59	48/50	77/80	59/61	67/75	69/79	68/77
zip	59/67	78/85	67/75	88/91	91/94	89/93	91/91	97/97	94/94
country	58/59	89/89	70/71	67/67	78/78	72/72	64/66	87/96	74/78
phone	97/99	84/87	90/93	84/89	91/96	87/92	92/93	85/85	88/89
email	100/100	99/99	100/100	94/95	99/99	96/97	99/99	98/98	98/98
company	57/81	37/51	44/63	-	-	-	-	-	-
dept.	51/85	31/45	38/59	-	-	-	-	-	-
overall	70/78	62/68	66/72	71/76	81/86	76/80	81/84	84/87	82/84

¹⁰ At the time of writing, degrees were not annotated as part of the Spanish collection and results for company and department names for Spanish and Czech were still work in progress.

<p>CENTRAL OFFICE 460 Capitol Avenue Hartford, Connecticut 06106 Fax: 860-418-6003</p> <p>Linda Goodman Birth to Three Director 860-418-6147 linda.goodman@po.state.ct.us</p> <p>Kathy Granata Administrative Assistant 860-418-6146 kathy.granata@po.state.ct.us</p> <p>Eileen McMurrer Child Find & Public Awareness Coordinator 860-418-6134 eileen.mcmurrer@po.state.ct.us</p> <p>Deb Resnick CSPD Coordinator 860-418-6151 deb.resnick@po.state.ct.us</p> <p>Alice Ridgway QA Manager 860-496-3073 Torrington 860-496-3087 Fax 860-418-6141 Hartford alice.ridgway@po.state.ct.us</p>	<p>Contact 366 (0.0009) 860 - 418 - 6003 Linda Goodman 860 - 418 - 6147 linda.goodman@po.state.ct.us</p> <p>Contact 443 (0.0011) Kathy Granata 860 - 418 - 6146 kathy.granata@po.state.ct.us</p> <p>Contact 440 (0.0011) Eileen McMurrer 860 - 418 - 6134 eileen.mcmurrer@po.state.ct.us</p> <p>Contact 437 (0.0011) Deb Resnick 860 - 418 - 6151 deb.resnick@po.state.ct.us</p> <p>Contact 434 (0.0011) Alice Ridgway 860 - 496 - 3073 Torrington 860 - 496 - 3087 860 - 418 - 6141</p>
--	---

Fig. 5. Sample of automatically annotated data; extracted instances on the right.

4.2 Product Catalogues

Another application of the *Ex* system is to extract information about products sold or described online. Our experiments have so far been limited to TV sets, computer monitors and bicycles. For each product type, an initial version of extraction ontology has been created to extract instances composed of typically a model name, price and multiple product-specific attributes. For example, the monitor extraction ontology contains 11 attributes, about 50 patterns and several axioms. The effort spent so far is about 2 person weeks. We experiment with a data set of 3,000 partially annotated web pages containing monitor ads. The average F-measure is now around 75% but complete evaluation has not been completed yet.

In order to obtain annotated data, we cooperate with one of the largest Czech web portals. The annotated data come from the original websites some of which send structured data feeds to the portal in order to get included in their online product database. The portal, on the other hand, can use such data to develop and train extraction models to cover the remaining majority of sites.

4.3 Weather Forecasts

Finally, we experimented with the domain of weather forecasts. Here our goal was to investigate the possibility to assist the ontology engineer in reusing existing *domain ontologies* in order to develop the extraction one/s. An advantage of this domain was the fact that several OWL ontologies were available for it. We analysed three of them by means of applying generic rules of two kinds:

1. Rules suggesting the *core class/es* for the extraction ontology. As the extraction ontology for extraction from HTML-formatted text¹¹ is typically more class-centric and hierarchical than a properly-designed domain ontology, only few classes from the domain ontology are likely to become classes in the extraction ontology, while others become attributes that are dependent on the core class/es. For example, ‘Day’ is typically an attribute of a ‘Forecast’ class in an extraction ontology, while in the domain ontology they could easily be two classes connected by a relationship. One of such core class selection rules is, in verbal form, e.g. “Classes that appear more often in the domain than in the range of object properties are candidates for core class/es.”
2. Rules performing the actual *transformation*. Examples of such rules are e.g. “A data type property D of class C may directly yield an attribute of C.” or “A set of mutually disjoint subclasses of class C may yield an attribute, whose values are these subclasses.”

Most such independently formulated selection and transformation rules appeared as performing well in the initial experiment in the weather domain; details are in [11]. Transformation rules seemed, by first judgement, to suggest a sensible and inspiring, though by far not complete, skeleton of an extraction ontology. Testing this ontology on real weather forecast records is however needed for proper assessment.

In general, although the first experiments look promising, extensive usage of domain ontologies as starting point for extraction ontologies seems to be hindered by unavailability of high-quality domain ontologies for most domains, e.g. in relation to different categories of products or services, judging by the results of Swoogle-based¹² retrieval. This obstacle is likely to disappear in the not-so-distant future, as the semantic web technology becomes more widespread.

4.4 Discussion

The main practical advantage of the approach based on extraction ontologies was expected to be the *rapid start* of the whole process while, at the same time, the growing body of extraction knowledge remains manageable in long term thanks to support for high-level *conceptual modelling*. This assumption seems to have been verified by our experience from the mentioned applications.

More complicated is, however, the comparison of the numerical results with systems based on other (especially, purely inductive) grounding. Their common assumption is that a corpus of *labelled training data* is available; then the overall task reduces to the invocation of a learning component on training data followed by application on the learnt model on unseen data. The success of such a setting can be best approximated by measuring the extraction *accuracy* on the labelled collection using n-fold *cross-validation*. Measuring purely by classification accuracy under the assumption of labelled data availability (which is a pre-requisite for inductive approaches while only

¹¹ This is not the case for extraction from free text, which is more relation-centric.

¹² <http://swoogle.umbc.edu>

minor help for extraction ontology building), our approach is in most cases slightly below the top-reported systems for common benchmark tasks.¹³ This observation follows, in particular, from our new experiments in the *seminar announcements* domain,¹⁴ for which a labelled corpus and benchmark results for several systems exist.¹⁵

However, extraction ontologies are assumed to be used, at first place, in the situations when there are very few or no labelled training data but enough domain knowledge. Manual building of the extraction ontology is then a (possibly quicker and less tedious) alternative to manual labelling of a sufficiently large sample of data. In many applications, including those mentioned above, the effect of saving the domain experts from labelling hundreds of pages outweighs the slightly lowered accuracy, especially if the output of extraction merely serves as support for a human examining the source documents (as in the medical website quality evaluation scenario).

5 Related Work

Most state-of-the-art WIE approaches focus on identifying structured collections of items (records), typically using inductively learnt models. Ontologies are often considered but rather as additional structures to which the extracted data are to be adapted after they have been acquired from the source documents, for the sake of a follow-up application [5]. There is no provision for directly using the rich structure of a domain-specific ontology in order to guide the extraction process.

Though our system has an optional wrapper-like feature, it also significantly differs from mainstream *wrapper* tools such as Kapow¹⁶ or Lixto¹⁷, which focus on building wrappers for each page/site separately, while our extraction ontologies can be reused within the whole domain. Among the wrapper-oriented approaches, the most similar to ours seems to be HiLeX [13], which allows to specify extraction ontologies as trees of linguistic and structural elements and evaluate them using a powerful logical language. Due to its assumption of nested rectangular semantic portions of web pages, it is however tuned to extraction from tabular data (rather than from more linearly-structured data), although it relaxes the dependence on HTML formatting proper.

The approach to WIE that is inherently similar to ours (and from which we actually got inspiration in the early phase of our research) is that developed by Embley and colleagues at BYU [3]. The main distinctive features of our approach are: (1) the possibility to provide the extraction patterns with probability estimates (plus other quantitative info such as value distributions), allowing to calculate the weight for every attribute candidate as well as instance candidate; (2) the effort to combine hand-crafted extraction ontologies with other sources of information—HTML formatting and/or known data instances (3) the pragmatic distinction between extraction ontologies and domain

¹³ Such comparison is however skewed by the impossibility to carry out cross-validation for manual creation of extraction knowledge, as a human would not be able to forget about the already seen data samples between the individual ‘runs’ as a software system is.

¹⁴ Initial results of the experiments are in a working paper [9].

¹⁵ <http://www.cs.umass.edu/~mccallum/code-data.html>

¹⁶ <http://www.kapowtech.com>

¹⁷ <http://www.lixto.com>

ontologies proper: extraction ontologies can be arbitrarily adapted to the way domain data are typically *presented* on the web while domain ontologies address the domain as it is (but can be used as starting point for designing extraction ontologies). For similarly pragmatic reasons (easy authoring), we also used a proprietary XML syntax for extraction ontologies. An objective comparison between both approaches would require detailed experiments on a shared reference collection.

An approach to automatically discover new extractable attributes from large amounts of documents using statistical and NLP methods is described in [12]. On the other hand, formatting information is heavily exploited for IE from tables in [17]. Our system has a slightly different target; it should allow for fast IE prototyping even in domains where there are few documents available and the content is semi-structured. While our system relies on the author to supply coreference resolution knowledge for attribute values, advanced automatic methods are described e.g. in [19]. The system described in [18] uses statistical methods to estimate the mutual affinity of attribute values.

Our ideas and experiments on domain ontology selection and transformation to extraction ontology are related to the generic research in ontology selection and content evaluation [14], especially with respect to the notion of intra-ontology concept centrality; this relationship deserves further study.

6 Conclusions

The *Ex* system attempts to unify the often separate phases of WIE and ontology population. Multiple sources of extraction knowledge can be combined: manually encoded knowledge, knowledge acquired from annotated data and knowledge induced from common formatting patterns by the means of wrapper induction. An alpha version of *Ex* (incl. extraction ontology samples) is publicly available¹⁸.

Future work will evaluate the system integrated with *trainable* machine learning algorithms and exploit some more coarse-grained web mining tools including *web page classifiers*. For the latter, a rule-based *post-processing engine* is under development; it will, following a higher-level domain-specific website model, filter and transform the extraction results based on the context—semantic class of the given page. Both *instance parsing* and *formatting pattern induction* algorithms themselves also need improvement in accuracy and speed. Furthermore, we plan to investigate how *text mining* over the extraction results could help us identify ‘gaps’ in the ontology, e.g. non-labelled tokens frequently appearing inside a ‘cloud’ of annotations are likely to be unrecognised important values. Finally, we intend to provide support for semi-automated transformation of *domain ontologies* to extraction ones.

7 Acknowledgments

The research leading to this paper was partially supported by the EC under contract FP6-027026, Knowledge Space of Semantic Inference for Automatic Annotation and Retrieval of Multimedia Content - K-Space. The medical website application is carried out in the context of the EC-funded (DG-SANCO) project MedIEQ.

¹⁸ <http://eso.vse.cz/~labsky/ex>

References

1. Ciravegna, F.: (LP)², an Adaptive Algorithm for Information Extraction from Web-related Texts. In: Proc. IJCAI-2001 Workshop on Adaptive Text Extraction and Mining, Seattle.
2. Duda, R.O., Gasching, J., and Hart, P.E. Model design in the Prospector consultant system for mineral exploration. In: Readings in Artificial Intelligence, pp. 334–348, 1981.
3. Embley, D. W., Tao, C., Liddle, D. W.: Automatically extracting ontologically specified data from HTML tables of unknown structure. In: Proc. ER '02, pp. 322–337, London, UK, 2002.
4. Huang, X., Acero, A., Hon, H. W.: Spoken Language Processing: A Guide to Theory, Algorithm and System Development. Prentice Hall, New Jersey, 2001.
5. Kiryakov, A., Popov, B., Terziev, I., Manov, D., Ognyanoff, D.: Semantic annotation, indexing, and retrieval. *J. Web Sem.*, volume 2, pp. 49–79, 2004.
6. Labský, M., Svátek, V.: Information extraction with presentation ontologies. Technical report, KEG UEP, <http://eso.vse.cz/~labsky/ex/ex.pdf>.
7. Labský, M., Svátek, V.: On the design and exploitation of presentation ontologies for information extraction. In: ESWC 2006 Workshop on Mastering the Gap: From Information Extraction to Semantic Representation. CEUR-WS, Vol.187, 2006.
8. Labský, M., Svátek, V., Šváb O.: Types and Roles of Ontologies in Web Information Extraction. In: ECML/PKDD04 Workshop on Knowledge Discovery and Ontologies, Pisa.
9. Labský, M., Svátek, V., Nekvasil, M.: Information Extraction Based on Extraction Ontologies: Design, Deployment and Evaluation. Working draft.
10. Lafferty, J., McCallum, A., Pereira, F.: Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In Proc. 18th International Conf. on Machine Learning, pp. 282–289. Morgan Kaufmann, San Francisco, CA, 2001.
11. Nekvasil, M., Svátek, V., Labský, M.: Transforming Existing Knowledge Models to Information Extraction Ontologies. In: 11th International Conference on Business Information Systems (BIS'08), Springer-Verlag, LNBI, Vol.7., pp. 106–117.
12. Popescu, A., Etzioni, O.: Extracting Product Features and Opinions from Reviews. In: Proceedings of the conference on Human Language Technology and Empirical Methods in Natural Language Processing (EMNLP 2005), Vancouver, BC, Canada, 2005, 339–346.
13. Ruffolo, M., Manna, M.: HiLeX: A System for Semantic Information Extraction from Web Documents. In: Proc. Enterprise Information Systems, Springer, LNBI, 2008, 194–209.
14. Sabou, M., Lopez, V., Motta, E.: Ontology selection for the real semantic web: How to cover the queen's birthday dinner? In: Proc. EKAW 2006. Springer LNCS, 2006, 96–111.
15. Stamatakis, K., Metsis, V., Karkaletsis, V., Růžička, M., Svátek, V., Amigó, E., Pöllä, M.: Content collection for the labeling of health-related web content. In: Proceedings of the 11th Conference on Artificial Intelligence in Medicine (AIME 07), LNAI 4594, Amsterdam, 7-11 July, (2007), 341–345.
16. Svátek, V., Labský, M., Vacura, M.: Knowledge Modelling for Deductive Web Mining. In: Proc. EKAW 2004, Springer Verlag, LNCS, 2004, 337–353.
17. Wei, X., Croft, B., McCallum, A.: Table Extraction for Answer Retrieval. In: Information Retrieval Journal, vol. 9, issue 5, pp. 589-611, 2006.
18. Wick, M., Culotta, A., McCallum, A.: Learning Field Compatibilities to Extract Database Records from Unstructured Text. Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing, Sydney, Australia, July 2006, 603–611.
19. Yates, A., Etzioni, O.: Unsupervised Resolution of Objects and Relations on the Web. In: Proc. Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics, 121–130.