

# Helping Novices Avoid the Hazards of Data: Leveraging Ontologies to Improve Model Generalization Automatically with Online Data Sources

Sasin Janpuangtong, Dylan A. Shell

■ The infrastructure and tools necessary for large-scale data analytics, formerly the exclusive purview of experts, are increasingly available. Whereas a knowledgeable data miner or domain expert can rightly be expected to exercise caution when required (for example, around misleading conclusions supposedly supported by the data), the nonexpert may benefit from some judicious assistance. This article describes an end-to-end learning framework that allows a novice to create models from data easily by helping structure the model-building process and capturing extended aspects of domain knowledge. By treating the whole modeling process interactively and exploiting high-level knowledge in the form of an ontology, the framework is able to aid the user in a number of ways, including in helping to avoid pitfalls such as data dredging. Prudence must be exercised to avoid these hazards as certain conclusions may only be supported if, for example, there is extra knowledge that gives reason to trust a narrower set of hypotheses. This article adopts the solution of using higher-level knowledge to allow this sort of domain knowledge to be used automatically, selecting relevant input attributes, and thence constraining the hypothesis space. We describe how the framework automatically exploits structured knowledge in an ontology to identify relevant concepts, and how a data extraction component can make use of online data sources to find measurements of those concepts so that their relevance can be evaluated. To validate our approach, models of four different problem domains were built using our implementation of the framework. Prediction error on unseen examples of these models show that our framework, making use of the ontology, helps to improve model generalization.

A variety of well-established supervised learning methods produce a model from a set of examples. Despite the maturity of these algorithms, decisions that result from models are unlikely to be correct if data have been used indiscriminately. This is part of the so-called data-dredging problem (Smith and Shah 2002). Figure 1 shows a bemusing example: U.S. spending per annum on science, space, and technology is highly correlated with suicides by hanging, strangulation, and suffocation. Logically, we know that this correlation does not imply causation (that is, higher spending on technology cannot cause more suicides by hanging or vice versa). Unfortunately, without the capacity to distinguish real and spurious correlations, learning methods are prone to picking up such correlations in producing models (Tukey 1977). The onus to be judicious ultimately falls on the person building the model. Otherwise, data dredging may lead to specious models, which may overfit, generalize poorly, and suggest conclusions that are fallacious.

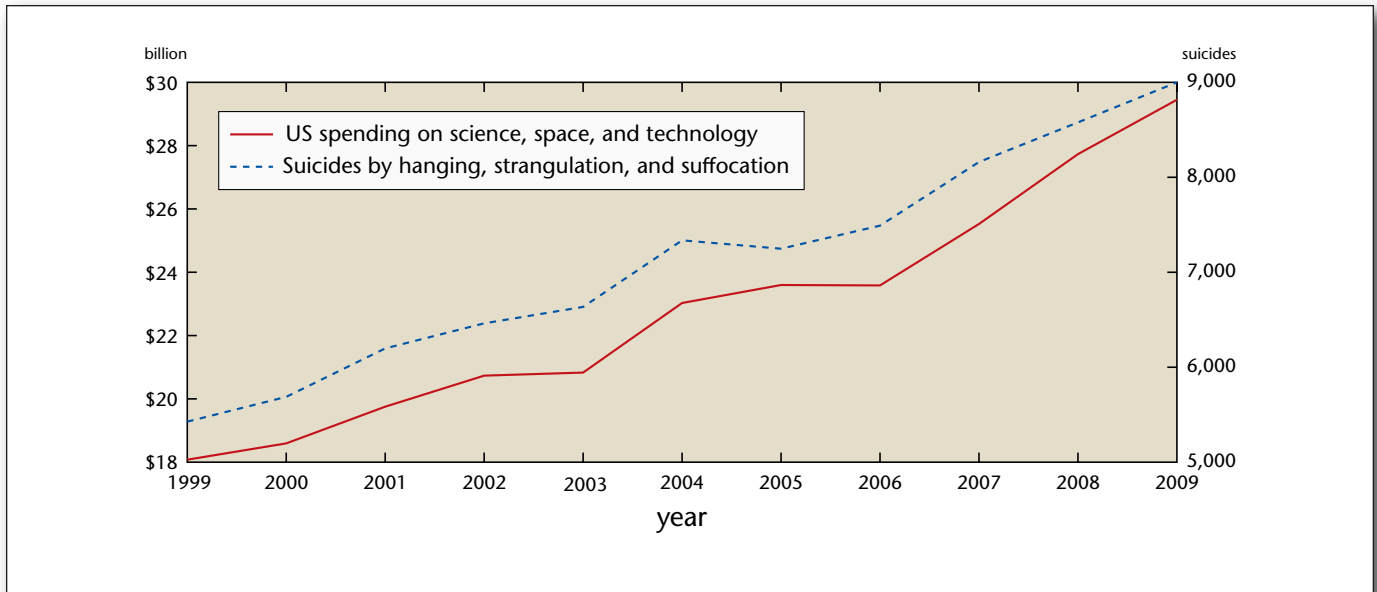


Figure 1. A Plot Showing Correlation Between Data of Two Attributes: U.S. Spending on Science, and Suicides by Hanging.

Data of these attributes are highly correlated with correlation = 0.992 (source: [www.tylervigen.com](http://www.tylervigen.com)).

The increased availability of data and the existence of easy-to-use procedures for regression and classification in commodity software allows inexperienced users to search for correlations among a large set of variables with scant regard for their meaning. Indeed, data dredging has been democratized and anyone may use seemingly sophisticated tools to arrive at unsound conclusions. This motivates the present work in developing a software framework to treat the whole modeling process rather than merely the model fitting stage. This framework should allow a non-expert user to (1) easily or even automatically create a model from existing data, (2) avoid the pitfalls of data dredging, and (3) build an accurate model, which can correctly predict unseen data.

Our work is inspired by the modeling process employed by Nelson and Sprecher (2008) to build a model of nuclear power use in a given country to help understand civil nuclear proliferation. Figure 2 shows the key ideas distilled from their modeling process. As experts, they used their knowledge of the nuclear domain to identify factors that are relevant to nuclear power. Critically, the choices made regarding the data inputs for the model were directed by domain knowledge at a conceptual level rather than correlations in the data themselves. Moreover, most of the data that was used to build their model came from existing online data sources, such as government and public organization websites. While this modeling process has broad promise, it involves two challenges for a nonexpert. First, it relies heavily on domain knowledge from the person building the model. Second, the data collection itself was performed manually. Better data collection methods

must be developed to help anyone collect data from different sources across a variety of formats.

We have set out to develop a semiautomated model-building framework that adopts key ideas from, but also addresses the challenges of, the modeling process previously described.<sup>1</sup> The framework operates as follows: (1) An existing ontology is used as a source of background knowledge rather than relying on knowledge from the person building a model. (2) Since ontologies are precise machine-manipulable representations of a priori structured relationships among concepts in a problem domain, they also enable a machine to explore the knowledge in an ordered way to determine the relevance of domain concepts. Such concepts are then used to construct an hypothesis space, and data are used to find the best model in this space using a learning method. (3) An ontological concept can be operationalized to measurements by finding a data source corresponding to that concept, then looking for measurements from information published on that source. (4) A data-extraction component is introduced into the framework to help a nonexpert extract desired data from a source easily.

This article is an attempt to bridge the gap between information extraction (IE) and learning from data, helping a user easily accomplish the learning task but also ensuring an accurate model is built. The IE research aims to automatically extract structured information from unstructured or semistructured data sources so that a machine can interpret and automatically make use of those data (Etzioni et al. 2008). Linking these two disciplines through semantic relationships underlying the data enables a

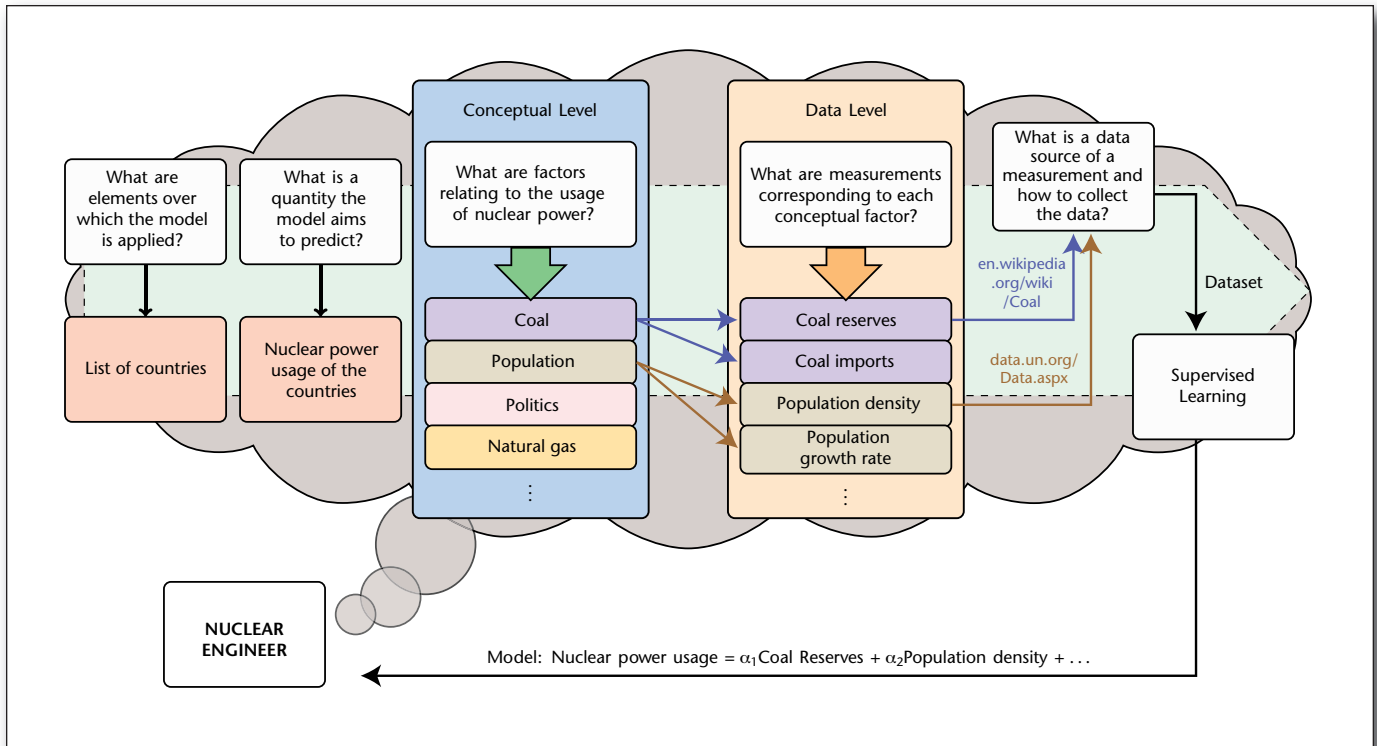


Figure 2. Modeling Process Extracted from Nelson and Sprecher (2008).

The expert carries out the process by progressing from left to right to build a model for predicting nuclear power use. The expert starts by specifying a list of countries and then collecting nuclear power data of those countries from an existing data source. Background knowledge is used to come up with several hypotheses about the factors that may influence nuclear power usage. For example, a country with large coal reserves might be expected to use coal to produce electricity rather than using nuclear power. Thus, coal may affect nuclear power usage. Next, the expert moves to find measurements that correspond with the factors. For instance, coal reserves and coal imports are possible measurements relating to the coal in a country. The expert selects one of these measurements for each factor and uses a data source to provide data for the selected measurement. The collected data from these sources are used to construct a data set for learning. Finally, the expert uses the resultant data set with a designated learning method to build a model.

machine automatically to build a model of the domain. Using relationships of this form to interpret and evaluate attribute relevance helps to impose structure on the model to reduce overfitting. Thus, the framework is a way to help inexperienced users to cope with the data deluge by making use of existing data.

## Learning Framework

Our end-to-end framework is illustrated in figure 3. It employs knowledge at two levels, high-level concepts and their relationships in an ontology, and low-level data from existing data sources. Model building proceeds from left to right and consists of three main components: finding relevant concepts, operationalizing concepts, and learning. Each component is associated with an ontology, online data sources, and a learning library to perform its task, respectively. For example, using this framework to build a model for predicting nuclear power use of countries, a nonexpert user starts by giving the first component a query string “Nuclear power” that will be the output of the

model, representing the quantity he or she aims to predict. This component uses structured relationships in an ontology (which we assume is given) to automatically retrieve a set of concepts, denoted  $\Theta$ , that are relevant to the input query.

So far, relationships are only captured at a high level between concepts. To evaluate its predictive value, concepts in  $\Theta$  must be operationalized to corresponding data. To perform this step the user describes the set of elements over which the model is applied (the model’s domain), and what the model aims to predict (the model’s output). Tabular input containing two columns is used for this purpose. In the example, the first column contains a list of countries and the second column contains nuclear power usage data. We assume that the user already has this tabular input (for example, it could be collected from a web page publishing these data in tabular format). The input query and tabular input are used to pose a question, in this case essentially asking “Which attributes of countries are relevant to nuclear power?” This tabular input is also stored in the framework as an initial data set.

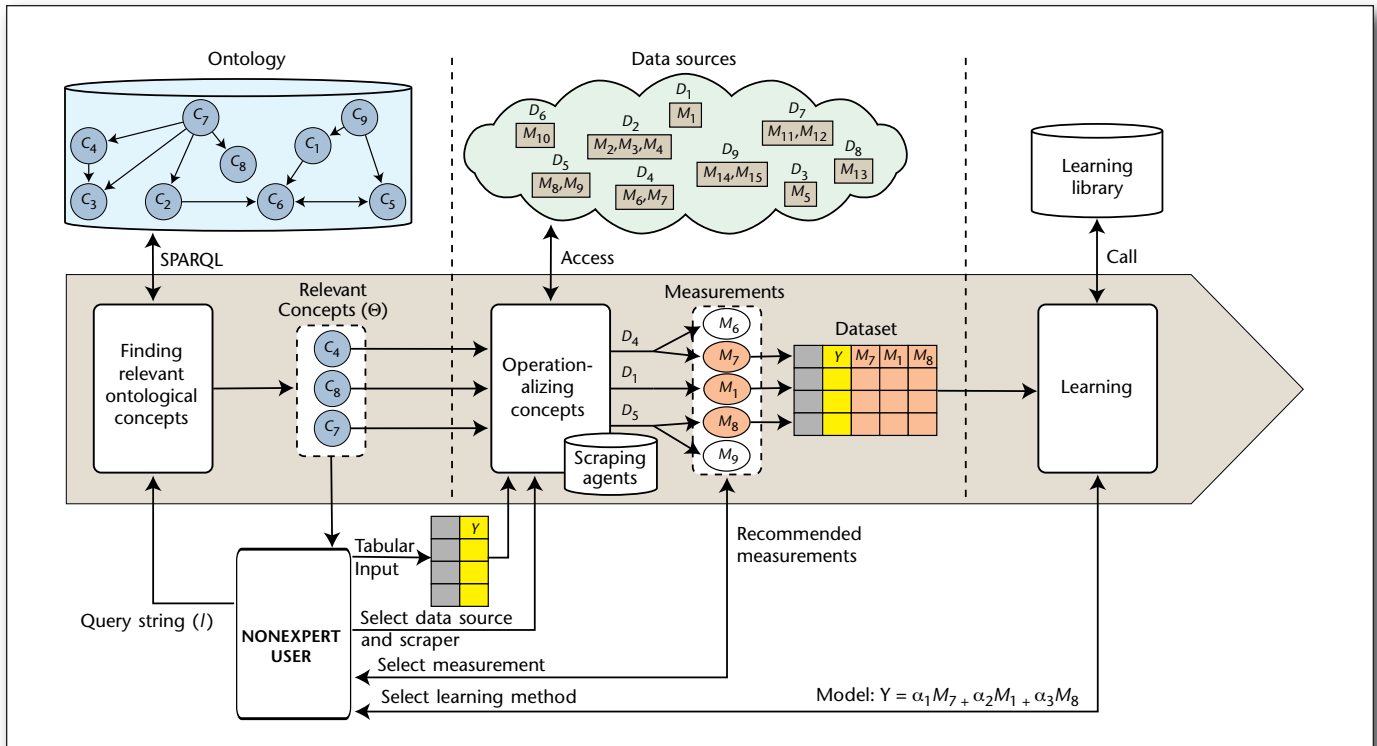


Figure 3. Proposed Learning Framework.

This framework incorporates both an existing ontology and data sources to build a model from data. Knowledge in the ontology is used to construct an initial model that is composed of relevant ontological concepts and data, the latter being associated with the concepts and is retrieved from existing data sources. The model is built and validated on the data using a selected learning method.

For each concept in  $\Theta$ , a data source (for example, a web page, Excel file, and others) that provides measurements or values is specified. The user selects a suitable scraping module (for example, table scraping, list scraping) to extract contents. If one of the concepts is “Coal,” representing the energy resource, then the user may provide the Wikipedia article,<sup>2</sup> which contains several tables with data related to this concept. He or she selects the table-scraping tool from the framework. The second component accesses the article and uses the selected scraper and data from the tabular input to extract all tables containing data about countries. These tables are represented as possible measurements related to the coal concept. The user chooses a table (the one providing coal reserves) that has columns (for example, subbituminous, lignite, total-coal, and others) containing data related to different aspects of coal reserves. He or she selects one of these column. Data from the selected column are extracted and then added into the data set. This step is repeated until all concepts in  $\Theta$  are operationalized to data and added to the data set. At this point, the data set is ready for use in the learning process (that is, the third component). The user selects a learning method (for example, linear regression, decision tree). The framework calls the selected

method from an existing library (for example, scikit-learn, R) to build a model from the data set. The model is returned to the user so it can be used for predictions, and to determine (for example, from coefficients) the relative importance of the various concepts. The user might conclude, for instance, that large coal reserves reduce the likelihood of a country building a nuclear power station.

The framework exploits relationships between three components: concepts in an ontology, existing data sources, and measurements of the concept and the given examples. Next we emphasize how the components involved help minimize the human effort required.

## Finding Relevant Concepts from an Ontology

We propose an algorithm for retrieving and ranking concepts that are relevant to a given query string  $\ell$  from a given ontology  $O$ . This algorithm makes use of  $O$ 's taxonomic hierarchy, which is a broadly applicable knowledge representation, useful across many ontologies. A taxonomic hierarchy is typically composed of two main elements: categories and con-

cepts. These elements are used to construct the hierarchy. Closely related concepts are organized in the same category. Closely related specific categories are also organized into the same broader category forming super- and subcategory relationships. We use SPARQL (the query language for RDF data) queries to acquire taxonomic knowledge from  $O$  without any preprocessing effort.

We adopt the ideas of the hyperlink induced topic search (HITS) algorithm (Kleinberg 1999), originally for searching and ranking relevant web documents on a given topic by considering two different notions of relevance: hubs and authorities. Representing ontological categories as hubs and concepts as authorities, HITS can be employed on the ontology to find and rank relevant concepts for a given input query. The intuition underlying our algorithm is that if a category is relevant to the input query, then all concepts in that category should also be relevant to the query. Likewise, if a concept is relevant to the query then all categories containing the concept should be considered as relevant as well. Since an ontological element is explicitly defined either as a category or concept, it has only one score associated with it (that is, the hub score for a category and the authority score for a concept), which also means that the scoring computation from HITS can avoid iterative updates.

### Scoring Schemes

Two separate scoring schemes are used to quantify the relevance and utility of finding other related information for categories and concepts. Both schemes are the product of two components: voting and frequency. Voting captures the relationships between categories and concepts, while frequency tracks how often categories or concepts reappear during execution of the algorithm.

#### Category Scoring

A category scores well if it links to many rare concepts that are relevant to an input query. Thus, category  $c$ 's score comes from voting from concepts in  $c$  and how many relevant concepts appear in  $c$  so that

$$Score(c) = \sum_{i \in I_R} IF(i, c) \times \sum_{i \in I_R} IVote(i, c), \quad (1)$$

where  $I_R$  denotes a set of relevant concepts,  $IF(i, c)$  is a function with value 1 if a concept  $i \in c$  or otherwise 0.  $IVote(i, c)$  gives the vote based on rarity of a concept  $i$ , such that

$$IVote(i, c) = \begin{cases} 1/|\text{categories containing } i| & \text{if } i \in c, \\ 0 & \text{otherwise.} \end{cases}$$

The  $IVote$  function states that if  $i$  appears in many categories, it is not a rare concept and its vote is shared among the many categories, so a category containing many common concepts is penalized.

#### Concept Scoring

A concept scores well if it is linked by many relevant categories. Thus, concept  $i$ 's score comes from the votes of categories containing  $i$ , and how many relevant categories  $i$  appears in, such that

$$Score(i) = \sum_{c \in C_R} CF(i, c) \times \sum_{c \in C_R} CVote(i, c), \quad (2)$$

where  $C_R$  denotes a set of relevant categories,  $CF(i, c)$  is 1 if  $i \in c$  or otherwise 0, and  $CVote(i, c)$  returns the score of  $c$  as:

$$CVote(i, c) = \begin{cases} Score(c) & \text{if } i \in c, \\ 0 & \text{otherwise.} \end{cases}$$

### Algorithm Details

The algorithm performs three steps to find and rank concepts that are relevant to the input query  $\ell$  in  $O$ . Details appear in figure 4. The algorithm starts by finding a concept  $\ell$ , whose label matches (textually)  $\ell$ , to represent the query. If more than one concept is found, the first is selected. Next, all concepts that link to (inlinks of) or receive a link from (outlinks of)  $\ell$  are retrieved from  $O$  to construct a set of initial relevant concepts,  $I_R$ .

Finding a set of relevant categories is performed in the second step (lines 5–24). For each concept in  $I_R$ , its categories are discovered and scores calculated using equation 1. The categories are sorted by score and the top  $n$  selected. Heuristics are used to further select categories from these top categories, finding those that (1) contain few concepts (that is, discarding categories that list names of films, animals, or scientists), and (2) contain neither very few nor many subcategories (that is, categories that are too specific or too general). The resulting set is denoted by  $C_R$ . Input parameters  $m$ ,  $min$ , and  $max$  are used to adjust this behavior. Suitable values depend on the ontology and problem domain.

The final step (lines 25–34) retrieves all concepts from each category in  $C_R$ , scoring each with equation 2. All concepts are sorted by score before being returned as the output.

## Implementation: DBpedia and Wikipedia

Our implementation leverages an existing ontology and online data sources; some implementation details are worth discussing. We used DBpedia (Bizer et al. 2009), the ontology counterpart of Wikipedia, as a source of background knowledge and used Wikipedia articles as data sources for corresponding DBpedia concepts. The vast amount of general knowledge in this ontology made testing on multiple case studies feasible. Moreover, each DBpedia concept has a corresponding Wikipedia article often containing detailed information associated with that concept in different formats, such as text, list, or

**Algorithm 1:**

**input:**  $\ell$  = Query string from user,  $O$  = SPARQL endpoint of ontology,  
 $n$  = Number of categories,  $m$  = Maximum concepts in category,  
 $min$  = Minimum sub-categories,  $max$  = Maximum sub-categories  
**Output:**  $\Theta$  = A set of ranked relevant concepts

```

1  $\epsilon \leftarrow \text{getConceptFromStr}(\ell, O)$ 
2  $Out \leftarrow \text{FindOutLink}(\epsilon, O)$ 
3  $In \leftarrow \text{FindInLink}(\epsilon, O)$ 
4  $I_R \leftarrow \epsilon \cup Out \cup In$ 
5 foreach  $i \in I_R$  do
6    $Cats \leftarrow \text{getCategories}(i, O)$ 
7    $vote \leftarrow 1/Length(Cats)$ 
8   foreach  $c \in Cats$  do
9      $C_{vote}[c] \leftarrow C_{vote}[c] + vote$ 
10     $C_f[c] \leftarrow C_f[c] + 1$ 
11 foreach  $c \in C_{vote}$  do  $C_{score}[c] \leftarrow C_{vote}[c] \times C_f[c]$ 
12  $C_{score\_sort} \leftarrow \text{SORT}(C_{score})$ 
13 while  $k < t$  do
14    $c \leftarrow C_{score\_sort}[k]$ 
15    $mb \leftarrow \text{CountConcepts}(c, O)$ 
16    $sb \leftarrow \text{CountSubCategories}(c, O)$ 
17   if  $mb < m$  AND  $(min < sb < max)$  then  $\text{Add}(C_R, c)$ 
18    $k \leftarrow k + 1$ 
19 foreach  $i \in I_R$  do  $I_f[i] \leftarrow 1$ 
20 foreach  $c \in C_R$  do
21    $Cons \leftarrow \text{getConcepts}(c, O)$ 
22   foreach  $i \in Cons$  do
23      $I_{vote}[i] \leftarrow I_{vote}[i] + C_{score}[c]$ 
24      $I_f[i] \leftarrow I_f[i] + 1$ 
25 foreach  $i \in I_{vote}$  do  $I_{score}[i] \leftarrow I_{vote}[i] \times I_f[i]$ 
26  $\Theta \leftarrow \text{SORT}(I_{score})$ 
27 Return  $\Theta$ 

```

Figure 4. Finding and Ranking Relevant Ontological Concepts.

table. Our implementation exploits this connection automatically to get an article of a concept and use it as a primary data source for finding possible measurements of the concept. Thus, the user needn't specify a data source for a concept, and the system can build a model with minimal human effort (as highlighted in the sequence in figure 5). Our system focuses on data that are published in tabular format since each table often encapsulates a complete,

nonredundant set of facts, and tables structure data for easy automatic interpretation and extraction (Bhagavatula, Noraset, and Downey 2013).

As shown in figure 5, suppose the user would like to build a model for predicting the GDP of countries. He or she can use this system by providing a query string "Gross Domestic Product." The system retrieves relevant concepts from DBpedia. In the second step the user selects the tabular input, with the

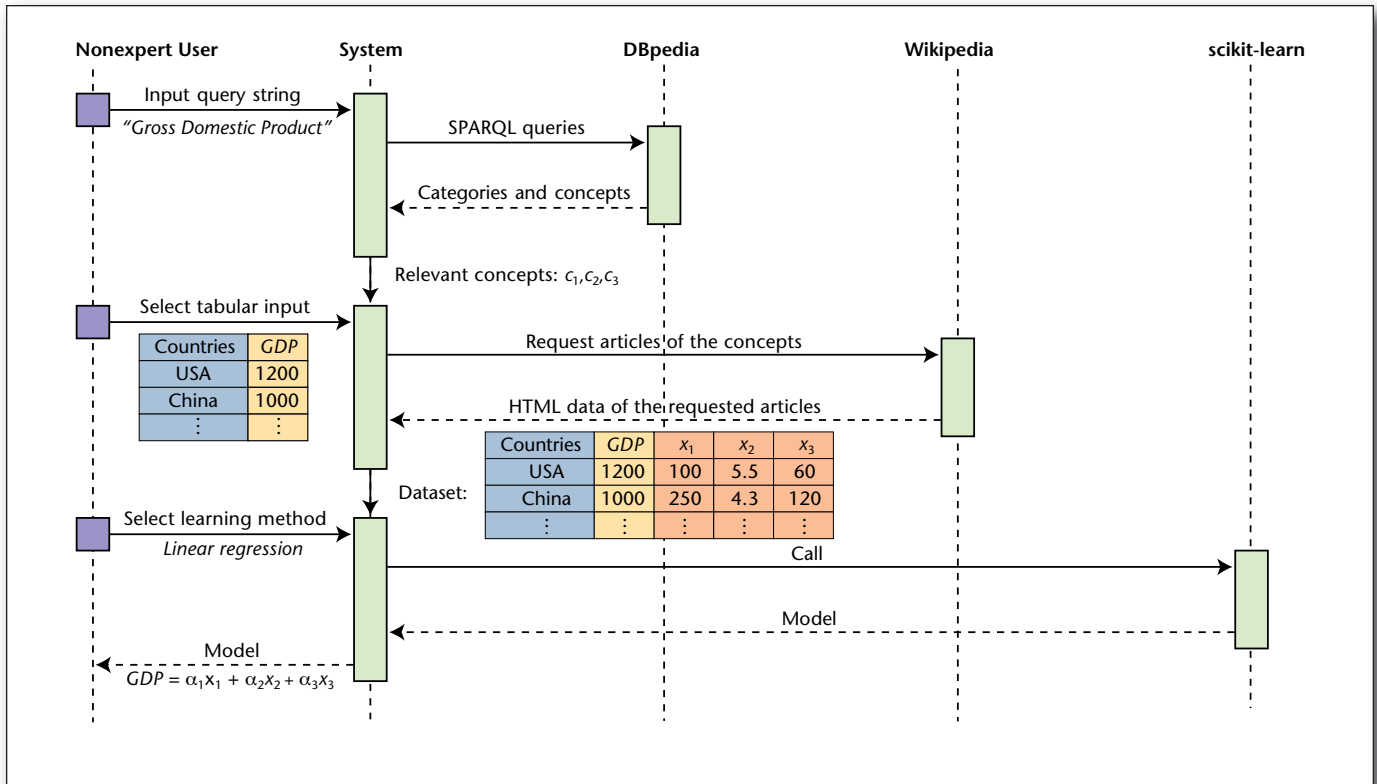


Figure 5. Sequence of Interactions with the Implementation.

A diagram showing an interaction use case with our DBpedia and Wikipedia implementation highlighting how a model can be built with little human effort. Light (green) rectangles indicate operations that are done automatically and dark (blue/purple) ones show where user intervention is required.

first column containing a list of countries and the second column containing GDP data for each country. Then the system requests HTML data of Wikipedia articles corresponding to concepts output from the first step. The system processes HTML data of these articles to construct a data set. In the third step, the user selects a learning method. The system calls that method from the Scikit-learn library to build a model from the data set. We can see that the user gives the system only three inputs and then lets the system carry out the remainder to build a model. This resonates with our earlier motivation of producing a system to help the nonexpert user to build a model from existing data easily. Implementation details of the first two steps are described in the following sections.

### Finding Relevant Concepts from DBpedia

There are three important points when implementing Algorithm 1 with DBpedia. Firstly, internal links among Wikipedia articles are used to find inlinks and outlinks of concept  $\epsilon$ . DBpedia already contains these internal links as RDF triples through the *wikiPageWikiLink* predicate. Secondly, suitable values parameters  $n$  and  $m$  for DBpedia were found to be in the ranges

200–250 and 120–200, respectively, depending on a problem domain, while the values of *min* and *max* were set to 6 and 30, respectively. Lastly, we added a DBpedia-specific condition at the end of the algorithm to further select only concepts whose name (after removing all prefixes) starts with the term *List of*. We found that corresponding articles for these concepts usually have tables providing data about a specific aspect of the concept. Focusing on these types of concepts allows our implementation automatically to find and collect data as needed. For instance, the Wikipedia article “list of countries by GDP” has a table that contains data for GDP by country (useful for the preceding example). Figure 6 shows how this implementation works to find relevant concepts from DBpedia for an input query. The sample results obtained through this implementation using input queries; Poverty and Gross Domestic Product are shown in table 1.

### Collecting Data from Wikipedia Tables

Algorithm 2 (see figure 7) automatically extracts data from a table on a Wikipedia page. For each concept in  $\Theta$ , this algorithm starts by eliminating the concept that is redundant with  $\epsilon$  by checking whether the

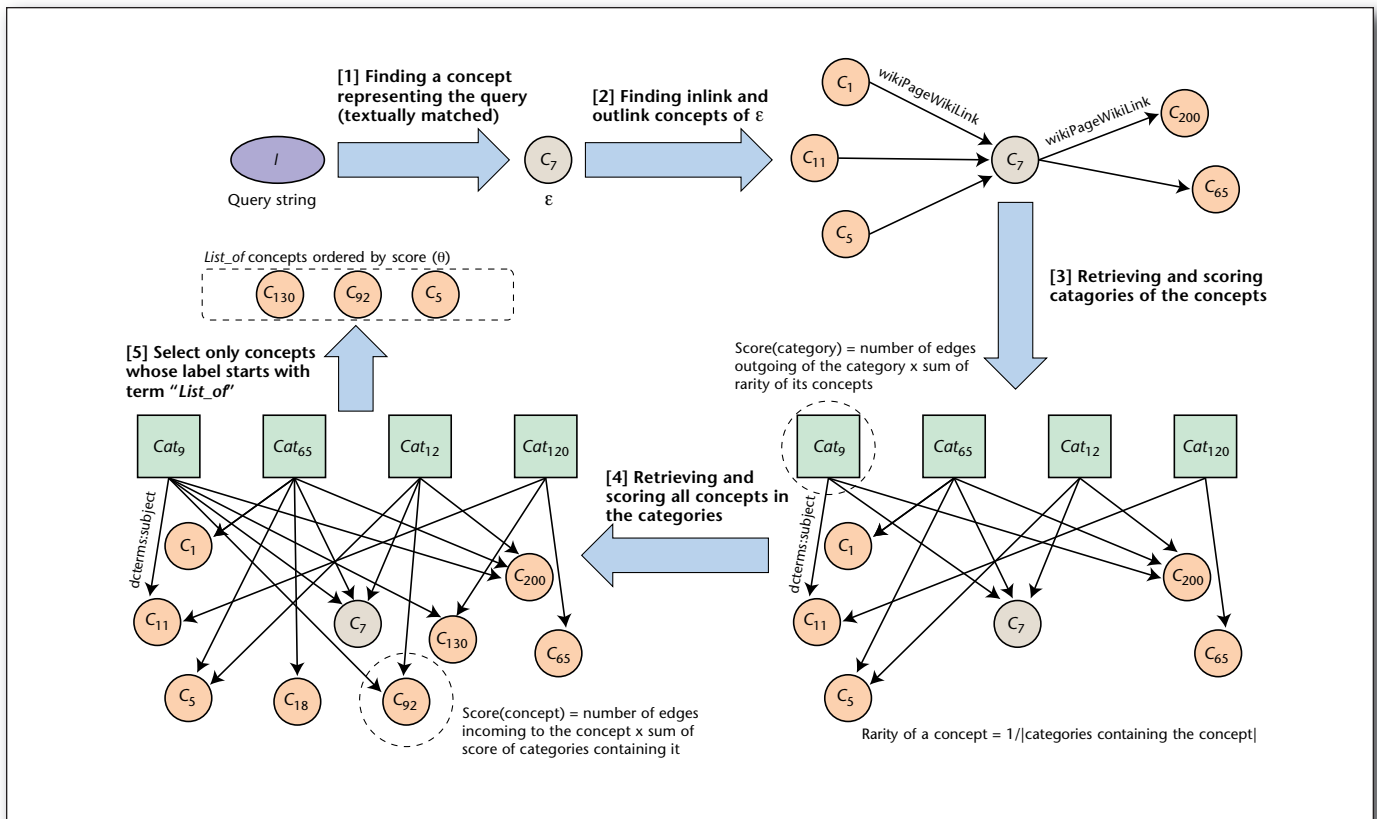


Figure 6. Steps of Running Algorithm 1 with DBpedia.

Running each step of algorithm 1 with taxonomic knowledge in DBpedia grows the graph as illustrated: (1) Finding a seed node,  $\epsilon$  representing the input query,  $l$ . (2) Outlinks and inlinks of  $\epsilon$  (red-filled circles) are added to the graph. (3) Categories of each added concept are discovered. Relevant categories (filled rectangles) are selected to add to the graph. (4) Extra concepts from the added categories are retrieved and then added. (5) Concepts whose label starts with the term *List of* are selected and then returned as output.

string  $l$  appears in the concept's label. It then acquires the URL of a Wikipedia article associated with the concept from DBpedia, requests the article in HTML, and extracts all tables from the result.

Heuristics are used to select a table that (1) has one column, which we call an example column, that partially matches to elements in the first column of the tabular input (countries in the nuclear example), and (2) has this example column appear first or second. If no table is selected, the concept is discarded. The algorithm then seeks columns in the selected table that contain numerical data. The numerical column closest to the example column is selected; if no numerical column is found, the example column is used to construct a new column that contains binary data (that is, value 1 is assigned to indicate that an element from the first column of the tabular input appears in the example column, otherwise value 0 is assigned).

At the end of this step, the system produces a tabular data set containing data from attributes of the model's domain (countries, in the nuclear example) and these attributes are also relevant to the input query.

Although algorithm 2 can automate table detection and data scraping of a Wikipedia article to operationalize a concept, this heuristic approach may select the wrong table or column to retrieve data yielding a poor final outcome. In a manual approach, by contrast, the system only retrieves tables from an article and shows them to a user. The user then selects a table and a column to add to a data set. Doing so, however, relies on the user having enough knowledge to select the correct table and column. A hybrid approach could be a better option, where the system initially chooses a table and a column for users and then lets them decide whether to accept that choice or change to another more appropriate table or column. The system can also provide a score for each table or column based on its contents (for example, number of matched countries in a table) to help the user make a wise choice.

### Implementing the Framework with Other Ontologies and Data Sources

It is worth discussing implementations with ontologies and data sources other than DBpedia and

**$\ell$  = Poverty**

countries\_by\_percentage\_of\_population\_living\_in\_poverty

countries\_by\_unemployment\_rate

countries\_by\_employment\_rate

countries\_by\_Sen\_social\_welfare\_function

permaculture\_project

sovereign\_states\_and\_dependent\_territories\_by\_fertility\_rate

global-manpower\_t\_for\_military\_service

wars\_and\_anthropogenic\_disasters\_by\_death\_toll

countries\_by\_sex\_ratio

countries\_by\_infant\_mortality\_rate

 **$\ell$  = Gross Domestic Product**

Australian\_states\_and\_territories\_by\_gross\_state\_product

sovereign\_states\_by\_external\_assets

countries\_by\_economic\_freedom

freedom\_indices

countries\_by\_Sen\_social\_welfare\_function

countries\_by\_percentage\_of\_population\_living\_in\_poverty

countries\_by\_percentage\_of\_population\_suffering\_from\_undernourishment

countries\_by\_energy\_intensity

countries\_by\_future\_gross\_government\_debt

countries\_by\_public\_debt

Table 1. Top 10 “List of” Concepts of Input Queries: “Poverty” and “Gross Domestic Product.”

Wikipedia. First, in algorithm 1 we assume that there is a concept,  $\epsilon$ , whose label matches an input query textually. An ontology used by the user could contain terminology that is more variable (for example, different acronyms and synonyms that do not directly match an input query). An additional step that can either automatically enrich the input query with synonyms or acronyms or find concepts that are semantically related to the query (Freitas et al. 2011, Pirro 2012) should precede the algorithm to resolve this issue. The user selects one of the concepts returning from this step and uses it as  $\epsilon$ . Second, other common predicates, such as owl:sameAs or rdfs:seeAlso could be used to find inlinks and outlinks of  $\epsilon$  rather than using the wikiPageWikiLink predicate, which is specific to DBpedia. Third, for algorithm 2, choosing a corresponding data source for a concept can be done directly by a user (for example, an Excel file containing data corresponding to the concept). A suitable user interface should be supplied to allow a user to specify the location of that data source. A search engine service, such as Google search service, could be employed in the sys-

tem to help a user easily find a suitable data source. By preprocessing the label for concept and then inputting it as a query to the search service, the user can get locations of online data sources corresponding to the concept and then select one of them to retrieve the data. Finally, scrapers for different data formats, such as lists, texts, files, should be added to the data-scraping toolbox, along with an interface that allows a user to select a suitable scraper. Inconsistent or poorly structured or formatted data pose difficulties for scraping. One way to handle this problem is to have the data-scraping module resort to an interactive mode where a user can see scraped results immediately and adjust the module to resolve any scrapping errors on the fly.

In this work, we describe only the back end of the system. To deploy this system, a suitable interface needs to be developed so that even inexperienced users without programming background can use the system. The system should also be able to record each decision made at each step by the user, so that the information can be used to improve the system. Moreover, the system should provide a mechanism

**Algorithm 2:**

**input:**  $\Theta$  = Output from Algorithm 1,  $\ell$  = Query from Algorithm 1,  $O$  = DBPedia's SPARQL endpoint,  
 $t_0$  = Tabular input  
**Output:** A dataset  $t$  in tabular format

```

1  $t \leftarrow t_0$ 
2 foreach  $i \in \Theta$  do
3   if Contain( $i, \ell$ ) then continue
4    $e\_idx \leftarrow 100$ 
5    $p \leftarrow \text{GetWikiArticle}(i, O)$ 
6    $T \leftarrow \text{ExtractTables}(p)$ 
7   foreach  $\tau \in T$  do
8      $c\_idx \leftarrow \text{GetExampleCol}(\tau, t_0[0])$ 
9     if  $c\_idx \neq -1$  AND  $c\_idx < e\_idx$  then
10       $t_b \leftarrow \tau$ 
11       $e\_idx \leftarrow c\_idx$ 
12    $n\_idx \leftarrow \text{FindNCol}(t_b, e\_idx)$ 
13   if  $n\_idx = -1$  then  $col_{new} \leftarrow \text{BCol}(t_b[e\_idx], t_0[0])$ 
14   else  $col_{new} \leftarrow \text{NCol}(t_b, e\_idx, n\_idx, t_0[0])$ 
15    $t \leftarrow \text{AddCol}(t, col_{new})$ 
16 Return  $t$ 

```

Figure 7. Operationalizing Ontological Concepts to Corresponding Data.

to save and load a model, making the model reusable and supporting later refinement of the model.

## Evaluation

Given the motivation underlying this work, two separate evaluations were conducted. The first was conducted to show that the generalization of a model built from a data set constructed by our system is improved over one based purely on data. Results from this evaluation provide support for the claim that the framework can help a nonexpert build an accurate model. The second evaluation was conducted to assess quality of input attributes selected through use of the ontology's background knowledge. Results from these evaluations show that the framework can help a nonexpert user avoid problems associated with data dredging.

We conducted these evaluations by building models of four different problem domains: Nuclear power, gross domestic product (GDP), poverty, and homelessness.<sup>3</sup> The model is applied to countries for the first three domains; U.S. states for the last one. For each problem domain, we formed two data sets

in order to build two different models for comparison. The first data set is constructed using our implementation as described in the preceding section. We denote this data set by  $t_{\text{ont}}$ . The second baseline data set, denoted  $t_{\text{base}}$ , was constructed by processing a URL of a Wikipedia category page containing links to many articles about the model's domain.<sup>4</sup> We visit every article in the category that has the term *List of* (countries/U.S. states) appearing in its URL and has not yet been visited when constructing  $t_{\text{ont}}$ . Algorithm 2 is executed on these articles to create a temporary data set denoted by  $t_{\text{temp}}$ . The data set  $t_{\text{base}}$  is then constructed by concatenating all columns from  $t_{\text{ont}}$  and  $t_{\text{temp}}$ .

Before using  $t_{\text{ont}}$  and  $t_{\text{base}}$  for learning, the issue of missing data in these data sets had to be addressed. For each problem domain, we examined  $t_{\text{base}}$  to remove any columns (except columns from  $t_{\text{ont}}$ ) and rows where 70 percent or more of their data are missing. We also removed the same set of rows from  $t_{\text{ont}}$ . Then, we manually filled in the remaining missing data for each column in both data sets by using an average value of data in that column. Finally, we invoked a learning method from the Scikit-learn

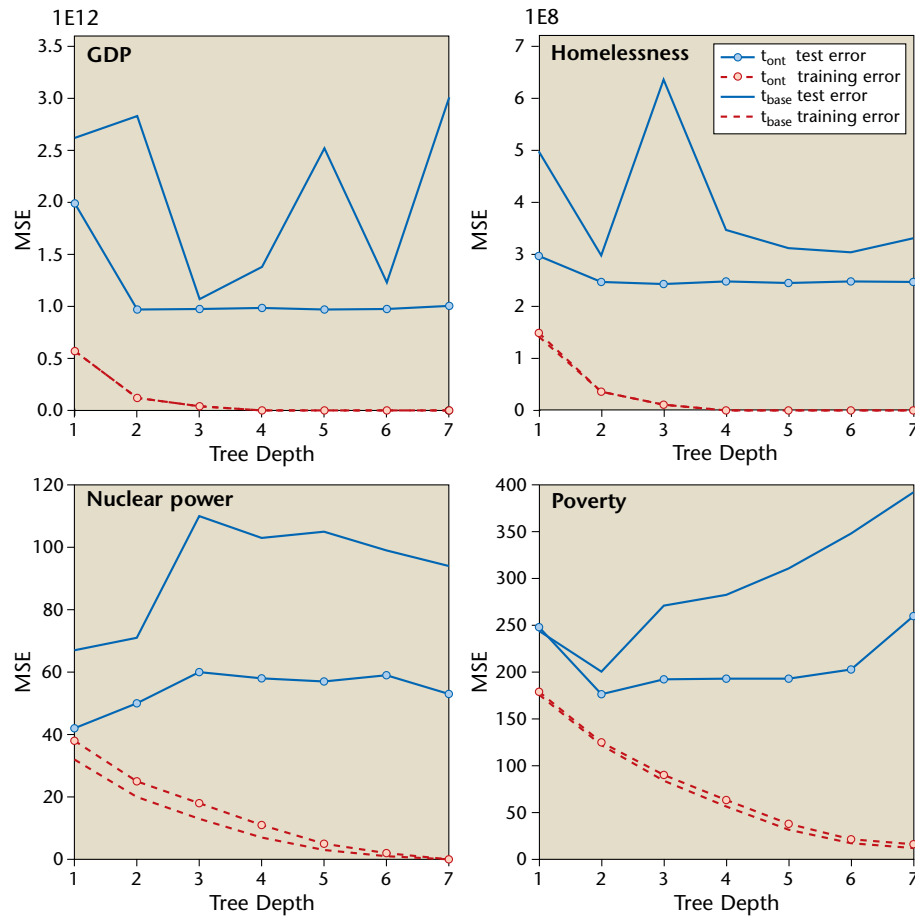


Figure 8. Comparing Generalization of Models.

For each problem domain, decision trees with different depths were built using data from  $t_{ont}$  and  $t_{base}$ . Average MSE values from 10-fold cross validation when testing these trees on training and test sets are shown for each depth. The models built from  $t_{base}$  show the occurrence of overfitting, while generalization of learned models from  $t_{ont}$  is improved.

library to build models from  $t_{ont}$  and  $t_{base}$  for each problem domain. The mean square error (MSE) is used to assess the quality of these models.

### Improving Generalization of a Model

Using knowledge in the ontology to select input attributes for learning could help improve generalization beyond given examples if the ranking incorporates (either implicitly or explicitly) causal and/or independence assumptions. To test this claim, the data sets  $t_{ont}$  and  $t_{base}$  for each problem domain were divided into training (80 percent of examples) and test sets. Training sets from  $t_{ont}$  and  $t_{base}$  contained the same set of examples (that is, both test sets also contain identical instances in the rows, but  $t_{ont}$  has

strictly fewer columns). Two decision trees for regression were learned from these training sets and then each tree was tested with the corresponding training and test sets to calculate MSE values. Finally, 10-fold cross validation was used to find the average MSE. We performed this evaluation repeatedly while increasing the depth of both trees, so as to examine overfitting and related phenomena.

From the results summarized in figure 8, we observe that in all problem domains trees learned using  $t_{base}$  are prone to overfitting (that is, when the complexity of the tree increases, prediction error of the tree on the training set decreases rapidly, but increases on the test set), whereas the learned trees from  $t_{ont}$  produce lower prediction errors on test sets

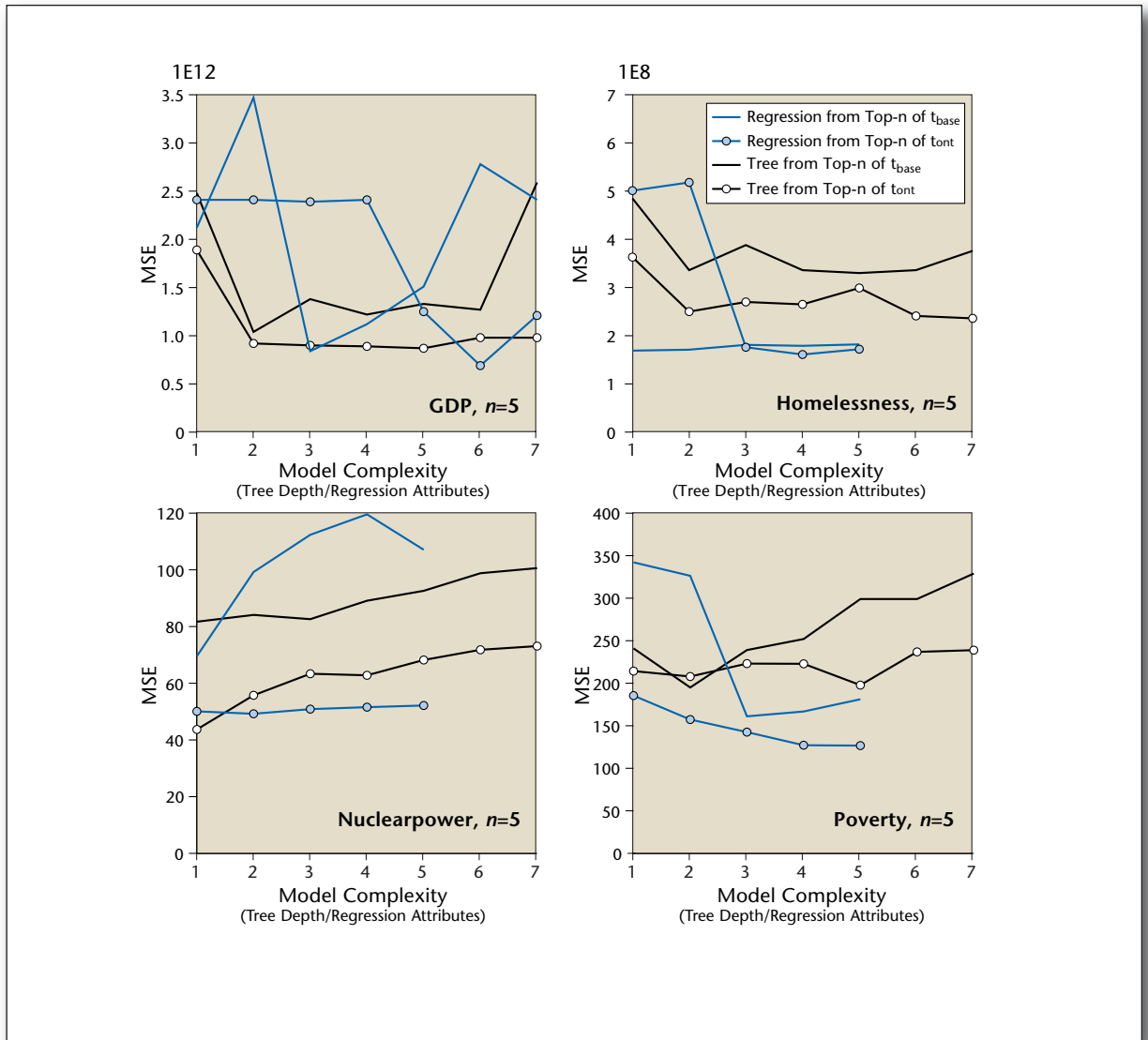


Figure 9. Assessing the Quality of the Input Attributes Selected by Using the Ontology's Background Knowledge.

This figure compares test set prediction errors of models learned from top- $n$  attributes from  $t_{ont}$  and  $t_{base}$  on the four problem domains. The models were constructed by using decision tree and linear regression methods. The results from both methods show that selecting attributes for learning by using prior domain knowledge helps improve generalization of the learned model.

when examining comparable depth. The results show that the framework helps improve generalization of a learned model so it predicts more accurately on unseen examples.

### Quality of the Top Ranked Attributes

We conducted a further experiment to show that the attributes automatically selected using the ontology's background knowledge are superior to attributes selected by correlations in a data set. For each problem domain, we constructed two new data sets to carry out this experiment. The first data set copies the first- $n$  columns from  $t_{ont}$ , where  $n = 5$  for all domains

except GDP where  $n = 8$  is used. Since each column in  $t_{ont}$  is ordered based on ranking of its corresponding concepts, this data set captures the top- $n$  concepts from our algorithm. The second data set is constructed by using a univariate feature-selection technique that selects columns based on correlations in the data set to select the top- $n$  columns from  $t_{base}$ . Then, we performed the same evaluations as before to assess performance of trees built from these new data sets. The results in figure 9 show that trees learned by using top- $n$  columns from  $t_{ont}$  produced prediction errors on test sets lower than trees learned by using top- $n$  columns from  $t_{base}$  in all problem

domains. These results suggest that selecting input attributes by using prior knowledge helps improve generalization of the learned model.

We also examined linear regression models to show that our results are independent of the learning method used to build a model. One expects that the generalization of a model built from  $t_{\text{ont}}$  should still be better than using  $t_{\text{base}}$  even if the learning method is changed. For each problem domain, we constructed two data sets. The first contained the first- $n$  columns of  $t_{\text{ont}}$ . The second data set is constructed by building a linear model from all attributes in  $t_{\text{base}}$ , ranking attributes based on absolute value of their coefficient (from high to low), and then selecting the top- $n$  attributes. Linear models with different complexities were built by limiting the number of attributes used. We started with all attributes in the data set and iteratively removed the lowest ranked attribute. The results in figure 9, especially in the nuclear power and poverty domains, support the same conclusion as the decision tree results.

We note that GDP and homelessness domains are challenging domains, and the learned models all have high prediction errors ( $\text{MSE} \times 10^{12}$  and  $10^8$ , respectively). These errors indicate that current attributes fail to capture the complexity of these problem domains. Improving our framework to enhance the quality of these models is part of our future work.

## Related Work

Several learning algorithms, such as knowledge-based artificial neural network (Shavlik and Towell 1989) and Bayesian belief networks (Pearl 1988; Russell and Norvig 2010), employ background knowledge to form an initial model and then use data to validate that model. Even though these algorithms have been demonstrated to outperform purely inductive learning (Mitchell 1997), their main limitation is that they can accommodate only a specific knowledge representation and learning method. In this work we present a framework that makes use of existing knowledge bases and data sources to build models of different problem domains. Also, this framework is designed to be independent of the learning method itself.

Semantic web technology provides data models for publishing background knowledge in a structured format so that a machine can automatically interpret and make use of the knowledge. Searching for elements that are relevant to a given query from structured knowledge is one of the main topics in the field of information retrieval (Franz et al. 2009; Cheng et al. 2008; Blanco, Mika, and Vigna 2011). Most of these works, however, require some preprocessing effort. None of them are specifically concerned with finding relevant ontological concepts to select attributes for learning.

Google Fusion Tables (Sarma et al. 2012) and Wikidata (Bhagavatula, Noraset, and Downey 2013) include an operation termed “Relevant Join” which uses data published in tabular format and aims to find suitable columns from different tables for joining to a given table. Our work can be viewed as a system that automatically performs a Relevant Join to construct a data set for learning. The main difference in our approach is that relevance of a column is justified by using prior domain knowledge rather than context in a table. Our system, moreover, need not be limited to data appearing in tabular formats. Data in another format, such as list, text, or query results, can also be used in the framework.

## Conclusion and Future Work

This article describes a learning framework that semi-automatically constructs a model using relevant ontological concepts and data attributes corresponding to those concepts. The attributes used in learning are selected by exploiting high-level knowledge separate from correlations within the data itself. As a consequence, the learned model is expected to generalize better than standard feature-selection approaches. We implemented this framework with DBpedia and Wikipedia and then used the implementation to build four models from four different problem domains. Prediction errors on unseen examples from these models are shown to support our claim. Moreover, the implementation helped build the models with very little human involvement.

What we present in this work is an attempt to address the changing needs of science: making it easier to produce models opens up vistas for inexperienced users, and helping automate the process of making sense of — and providing new interpretations for — existing data is one way to tame the deluge of data.

We believe that we are just starting to find uses of knowledge in an ontology to improve model generalization. Even though the results from the evaluations point out the possibility of using ontologies as background knowledge to help an inexperienced user build an accurate model from data, some important questions remain. Most obvious is that no formal explanation of how and why the system works was provided. Addressing these aspects is ongoing.

## Notes

1. Although we have emphasized its use by Nelson and Sprecher, the approach represents a standard approach in several sciences.
2. [en.wikipedia.org/wiki/Coal](http://en.wikipedia.org/wiki/Coal).
3. These four were chosen because the authors encountered articles talking about the desire to build models of these problem domains.



Visit AAI on  
LinkedIn™

AAI is on LinkedIn! If you are a current member of AAI, you can join us! We welcome your feedback at [info16@aaai.org](mailto:info16@aaai.org).

4. For example, [en.wikipedia.org/wiki/Category:Lists\\_of\\_countries\\_for\\_countries](http://en.wikipedia.org/wiki/Category:Lists_of_countries_for_countries).

## References

- Bhagavatula, C. S.; Noraset, T.; and Downey, D. 2013. Methods for Exploring and Mining Tables on Wikipedia. In *Proceedings of the ACM SIGKDD Workshop on Interactive Data Exploration and Analytics*, 18–26. New York: Association for Computing Machinery. [dx.doi.org/10.1145/2501511.2501516](https://doi.org/10.1145/2501511.2501516)
- Bizer, C.; Lehmann, J.; Kobilarov, G.; Auer, S.; Becker, C.; Cyganiak, R.; and Hellmann, S. 2009. DBpedia — A Crystallization Point for the Web of Data. *Journal of Web Semantics: Science, Services and Agents on the World Wide Web* 7(3): 154–165. [dx.doi.org/10.1016/j.websem.2009.07.002](https://doi.org/10.1016/j.websem.2009.07.002)
- Blanco, R.; Mika, P.; and Vigna, S. 2011. Effective and Efficient Entity Search in RDF Data. In *The Semantic Web: Semantics and Big Data: Proceedings of the 10th International Semantic Web Conference*, Lecture Notes in Computer Science Volume 7882, 83–97. Berlin: Springer. [dx.doi.org/10.1007/978-3-642-25073-6\\_6](https://doi.org/10.1007/978-3-642-25073-6_6)
- Cheng, G.; Ge, W.; Wu, H.; and Qu, Y. 2008. Searching Semantic Web Objects Based on Class Hierarchies. In *Proceedings of Linked Data on the Web Workshop*, Ceur Workshop Proceedings volume 369. Aachen, Germany: RWTH Aachen University.
- Etzioni, O.; Banko, M.; Soderland, S.; and Weld, D. 2008. Open Information Extraction from the Web. *Communications of the ACM* 51(12): 68–74. [dx.doi.org/10.1145/1409360.1409378](https://doi.org/10.1145/1409360.1409378)
- Franz, T.; Schultz, A.; Sizov, S.; and Staab, S. 2009. TripleRank: Ranking Semantic Web Data by Tensor Decomposition.

In *The Semantic Web: Proceedings of the 8th International Semantic Web Conference*, Lecture Notes in Computer Science Volume 5823, 213–228. Berlin: Springer.

Freitas, A.; Oliveira, J. G.; O’Riain, S.; Curry, E.; and Silva, J. C. P. D. 2011. Querying Linked Data Using Semantic Relatedness: A Vocabulary Independent Approach. In *Natural Language Processing and Information Systems*, Lecture Notes in Computer Science Volume 6716, 40–51. Berlin: Springer. [dx.doi.org/10.1007/978-3-642-22327-3\\_5](https://doi.org/10.1007/978-3-642-22327-3_5)

Kleinberg, J. M. 1999. Authoritative Sources in a Hyperlinked Environment. *Journal of the ACM* 46(5): 604–632. [dx.doi.org/10.1145/324133.324140](https://doi.org/10.1145/324133.324140)

Mitchell, T. M. 1997. *Machine Learning*. New York: McGraw-Hill.

Nelson, P., and Sprecher, C. M. 2008. What Determines the Extent of National Reliance on Civil Nuclear Power? In *Proceedings of the 49th Annual Meeting of the Institute of Nuclear Materials Management*, 72–79. Deerfield: Institute of Nuclear Materials Management.

Pearl, J. 1988. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann Series in Representation and Reasoning. San Francisco: Morgan Kaufmann.

Pirro, G. 2012. REWOrd: Semantic Relatedness in the Web of Data. In *Proceedings of the 26th AAAI Conference on Artificial Intelligence*, 129–135. Palo Alto, CA: AAAI Press.

Russell, S., and Norvig, P. 2010. *Artificial Intelligence: A Modern Approach*, Third Edition. Engelwood Cliffs, NJ: Prentice Hall.

Sarma, A. D.; Fang, L.; Gupta, N.; Halevy, A.; Lee, H.; Wu, F.; Xin, R.; and Yu, C. 2012. Finding Related Tables. In *Proceedings of the International Conference on Management of Data*, 817–828. New York: Association for Computing Machinery.

Shavlik, J., and Towell, G. 1989. An Approach to Combining Explanation-Based and Neural Learning Algorithms. *Connection Science* 1(3): 233–255. [dx.doi.org/10.1080/09540098908915640](https://doi.org/10.1080/09540098908915640)

Smith, G., and Shah, E. 2002. Data Dredging, Bias, or Confounding: They Can All Get You into the BMJ and the Friday Papers. *BMJ* 325(7378): 1437–1438. [dx.doi.org/10.1136/bmj.325.7378.1437](https://doi.org/10.1136/bmj.325.7378.1437)

Tukey, J. 1977. *Exploratory Data Analysis*. Boston: Addison-Wesley.

**Sasin Janpuangtong** is a Ph.D. candidate in the Department of Computer Science and Engineering at Texas A&M University in College Station, Texas. His research focus is on developing a system that leverages both background knowledge from existing ontologies and data from online data sources to help nonexpert users easily and accurately build models from data. He is also interested in semantics for big data and data analytics for smart city applications.

**Dylan A. Shell** is an associate professor in the Department of Computer Science and Engineering at Texas A&M University in College Station, Texas. He and his students work on problems in distributed AI and robotics; in addition to the material in this article, his interests include algorithms for multirobot task allocation, approaches for modeling and reasoning about groups of autonomous agents, and forms of self-organization.