

COMBINE TAG AND VALUE SIMILARITY FOR DATA EXTRACTION AND ALIGNMENT

D.Phani Sri Lakshmi, D.N.S.B.Kavitha

Department of Computer Science Engineering
Bhimavaram, India

Abstract— Based on a user's query Web databases create query result pages. For many applications, such as data integration, which need to cooperate with multiple web databases there is a need to automatically extract the data from these query result pages. So we present a data extraction and alignment method called CTVS which combines both tag and value similarity. The data values from the same attribute are put into the similar column in which CTVS automatically extract data from query result pages by first identifying and segmenting the query result records (QRRs) in query result pages and then align the segmented QRRs into a table. Specially, we advise new techniques to switch the case when the QRRs are not secure, which may be due to the presence of main information, such as a commentary, proposal or advert, and for handling any nested structure that may exist in the QRRs. By CTVS, we create novel record alignment algorithms that align the attributes in a record, in pair wise first and then holistically. Experimental results show that CTVS achieves high precision and outperforms alive state-of-the-art data extraction methods.

Index Terms—Automatic wrapper generation, data extraction, data record alignment, information integration, nested structure processing etc (*key words*)

I. INTRODUCTION.

ONLINE database, called web database, include pages in the deep web are dynamically generate in response to a user query submitted through the query interface of a web database, which Compared with webpages in the surface web, that can be accessed by a unique URL. Upon receiving a user's query, a web database returns the related data, either structured or semi structured, encoded in HTML pages. Many web applications need the data from multiple web databases, such as meta querying, data integration and comparison shopping. For these applications to further use the data embedded in HTML pages, automatic data extraction is necessary. Only when the data are extracted and organized in a structured way such as tables, they can be compared and aggregated. Hence, accurate data extraction is very important for these applications to perform correctly.

This paper focuses on the problem of automatically extract data records that are encoded in the query result pages

generate by web databases. The goal of web database data extraction is to eliminate any unrelated information from the query result page, extract the query result records (referred to as QRRs in this paper) from the page, and align the extracted QRRs into a table such that the data value belong to the same characteristic are placed into the same table column.

The following two-step method, called combine Tag and Value Similarity (CTVS), to extract the QRRs from a query result page p.

- Record extraction: Identifies the QRRs in p and involve two sub steps: data region identification and the actual segmentation.
- Record alignment: Aligns the data values of the QRRs in p into a table so that values for the same attribute are aligned into the same table column.

CTVS precisely extracts and aligns the QRRs in query result pages if there are at least two records in the page. Compared with existing data extraction methods, CTVS improve data pulling out precision in three ways.

- New methods to imagine that the QRRs are existing contiguously in only one data region in a page. However, it is not best guess be true for many web databases where maintain in order separate the QRRs. We scan 10 websites to which the QRRs in the query result pages are non-contiguous, but indicate that non-contiguous data regions are rather familiar. Furthermore, 10 out 4 of websites have non-contiguous QRRS with the same parent in the page HTML tag tree. This problem is concentrate on two methods according to the outline of the QRRs and the major in series in the result HTML tag trees (i.e., DOM trees). a. An off-the-peg data region recognition method is likely to recognize the non-contiguous QRRs that have the similar parents according to their tag similarity. b. A combine method is proposed to join different data regions that contain the QRRs (with or without the same parent) into a single.

- In this we proposed to align the data values in the recognized QRRs, first pairwise then holistically. Together tag structure similarity and data value comparison are used in the pairwise alignment procedure i.e, first to join tag structure and data value similarity to achieve the alignment.

- A new nested-structure dealing out algorithm is future to handle any nested structure in the QRRs after the

holistic alignment i.e, CTVS uses both tag and data value similarity in sequence to get better nested structure processing exactness.

II. QRR EXTRACTIONS

Fig.1 shows the framework for QRR extraction. In this query result page, the Tag Tree structure module first construct a tag tree for the page rooted in the <HTML>tag. Each node correspond to a tag in the HTML page as well as its children are tags covered inside it. In every inner node n of the tag tree has a tag string tsn, which include the tags of n and all tags of n's descendants, and a tag path tpn, which include the tags from the root to n.

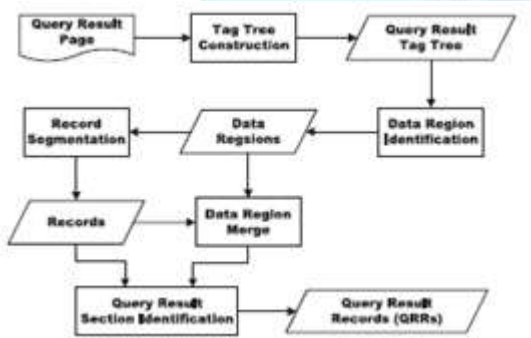


Fig.1. QRR extraction framework.

A. Data Region Identification

The Data Region Identification modules identify all likely data region, which usually contain dynamically generates data, top down starting from the root node. We first suppose that some child subtrees of the same parent node form alike data records, which assemble a data region. Thus, we propose a new process applied to more web databases. The data region identification algorithm is apply to a node n and recursively to its children $n_i, i = 1 \dots m$ as follows:

- Compute the similarity computation method in the data region identification algorithm is recursively apply to the children of n_i only if it does not have any alike siblings. Several data regions may be recognized in this.
- Section the data region into data records using the record segmentation algorithm.
- Suppose that the tag tree has n internal nodes and a node has a maximum of m children and a maximum tag string length of l. The time complexity of the data region identification algorithm is $O(nm2l^2)$.

B. Record Segmentation

The Record Segmentation module then segments the identified data regions into data records according to the tag

patterns in data regions. If only one tandem replicate is create in a data region, we assume that each frequent instance in the tandem replicate corresponds to a record. If many tandem replicate are found in a data region, we need to select one to indicate the record. The two heuristics are used for the tandem replicate collection

- If there is supplementary information, which corresponds to nodes between record instances, within a data region, the tandem replicate that stops at the supplementary information is the correct tandem replicate since information usually is not inserted into the middle of a record.
- The optical gap between two records in a data region is usually larger than any optical gap within a record. Hence, the tandem replicate that satisfies this restriction is selected.
- If the earlier two heuristics cannot be used, we select the tandem replicate that start the data region.

C. Data Region Merge

Given the segmented data records, the Data Region Merge module merges the data regions contain similar records. If any two data regions, we treat them as similar if the segmented records they contain are similar. The similarity between any two records from two data regions is measured by the resemblance of their tag strings. The resemblance between two data regions is calculated as the average record similarity. Two data regions can be merged into a merged data region if the records in the two data regions have an average similarity greater or equal to 0.6. Given, that two data regions d_1 and d_2 with n_1 and n_2 records and most record tag string length of l_1 and l_2 , Respectively, the time complexity of the data region merges algorithm is $O(n_1n_2l_1l_2)$.

D. Query Result Section Identification

The Query Result Section Identification module select one of the merged data regions as the one that contains the QRRs, there may still be multiple data regions in a query result page. Three heuristics are used to recognize this data region, called the query result section.

- The query result section generally occupies a large space in the query result page. For each data region d,
- an area weight, which is calculated as d's area divided by the largest area of all identified data regions, is assigned for d.
- The query result section is generally located at the center of the query result page. For each data region d, a center distance is considered among its center and the center of the page, and a center distance weight, which is calculated as the smallest center distance among all identified regions divided by d's center distance, is assigned for d. If a merged data region d contains multiple regions d_1, \dots, d_n to be found in dissimilar parts of the page, then first find the region d_i that has the largest space in the query result page in the middle of d_1, \dots, d_n and

assume that the center distance weight of d_i is the center distance weight of d .

- d) Each QRR generally contains more raw data strings than the raw data strings in other sections. For each data region d , a value weight, which is calculate as the average number of raw data strings in the records of d divided by the largest average number of data values in all recognized regions, is assign for d . All the above three weights are summed and the data

III. QRR ALIGNMENT

QRR alignment is performing by novel three-step data alignment methods that combine tag and value similarity.

A. Pairwise QRR Alignment

The pairwise QRR alignment aligns the data values in a pair of QRRs to present the proof for how the data values must be aligned among all QRRs. A pairwise alignment of r_1 and r_2 is composed of a set of data value alignments, each of which assumes that the corresponding data values from r_1 and r_2 belong to the same attribute.

Every QRR includes two kinds of information: the text string for the i th value and the tag path for the i th value. Throughout the pairwise alignment, we involve that the data value alignments must suit the following three constraints:

- Same record path constraint: The record path of a data value f comprises the tag from the root of the record to the node that contains f in the tag tree of the query result page. Each pair of corresponding values have the similar tag path. Hence, if f_{1i} has an altered tag path with f_{2j} , then s_{ij} is assigned a small negative value to prevent the pair of values from being aligned.
- Unique constraint: Each data value can be aligned for the most part one data value from the other QRR.
- No cross alignment constraint: If f_{1i} is matched to f_{2j} , then there should be no data value alignment between f_{1k} and f_{2l} such that $k < i$ and $l > j$ or $k > i$ and $l < j$.

Based on these constraints, a dynamic programming algorithm aligns the two records. The similarity is 0 if one of the QRRs is empty or else, if f_{1i} and f_{2j} have the same tag path, then just one of the following three data value alignments is possible.

- The first $(i - 1)$ values of r_1 can be aligned with the first $(j - 1)$ values of r_2 plus the data value alignment between f_{1i} and f_{2j} , which has the summing similarity score $L_{(i-1)(j-1)} + s_{ij}$.
- f_{1i} can be ignored and the first $(i - 1)$ values of r_1 can be aligned with the first j values of r_2 , which has the summing similarity score $L_{(i-1)j}$.
- f_{2j} can be ignored and the first i values of r_1 can be aligned with the first $j - 1$ values of r_2 , which has the summing similarity score L_{ij-1} .

region that has the biggest summed weight is selected as the query result section. Records in this data region are supposed to be QRRs.

A restriction of this approach is that if a query result page has more than one data region that contains query result records and the records in the different data regions, then we will choose only one of the data regions and discard the others.

The alignment with the largest summing similarity score among these three alternatives is chosen. That is

$$L_{ij} = \max(L_{(i-1)(j-1)} + s_{ij}, L_{(i-1)j}, L_{ij-1}) \quad (1)$$

With this dynamic programming method, the time complexity of the pairwise alignment algorithm is $O(l_1 l_2)$ where l_1 and l_2 are the number of data values in the two QRRs. Hence, given a data region with n records, the time complexity of the pairwise alignment algorithm is $O(n^2 l_2)$ in which l is the largest number of data values in a record.

i. Data Value Similarity Calculation

Given two data values f_1 and f_2 from dissimilar QRRs, we need their similarity, s_{12} , to be a real value in $[0, 1]$. The data value similarity is calculated according to the data type tree shown in Fig. 2.

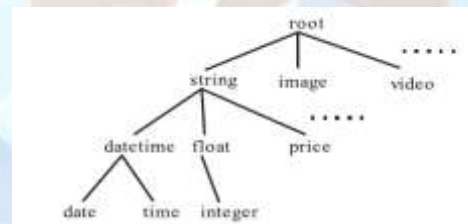


Fig. 2. Data type tree

Each child node is a subset of its parent node. For example, the “string” type includes several children data types, which are frequent on the web such as “datetime,” “float,” and “price.” The highest depth of the data type tree is 4. We will transfer to a nonstring data type as a specific data type.

Given two data values f_1 and f_2 , we first moderator data types and then fit them as extremely as possible into the nodes n_1 and n_2 of the data type tree. For example, given a string “567,” we will put it in node “integer.” The similarity s_{12} between two data values f_1 and f_2 with data type nodes n_1 and n_2 is defined as nonstring data type as a specific data type.

$$s_{12} = \begin{cases} 0:5 & n_1 = p(n_2) \ \& \ n_1 \neq \text{String} \ \text{OR} \\ 1 & n_2 = p(n_1) \ \& \ n_2 \neq \text{String} \\ \text{cosine similarity} & n_1 = n_2 \neq \text{String} \\ 0 & \text{otherwise,} \end{cases}$$

where $p(n_i)$ refers to the parent node of n_i in the data type tree. The similarity between data values f_1 and f_2 is

- . 0.5, if they belong to different specific data types that have a common parent.
- . 1, if they belong to the same specific data type.
- . string cosine similarity of f_1 and f_2 , if both f_1 and f_2 belong to the string data type.
- . 0 otherwise, which occurs when one of f_1 and f_2 belongs to the string data type and the other one belongs to a specific data type, or f_1 and f_2 belong to different specific data types without any direct parent.

As Table 1 shows, data values with the same data type usually have larger similarity.

TABLE 1
Example of Data Value Similarity

f_1 / n_1	f_2 / n_2	s_{12}
23/int	876/int	1
Jul 1, 2007/date	2005.12.3, 1:30/datetime	0.5
\$32.5/price	\$18/price	1
Harry Potter and the Goblet of Fire/string	Harry Potter and the Prisoner of Azkaban/string	0.714
Harry Potter and the Goblet of Fire/string	876/int	0

B. Holistic Alignment

Holistic alignment align the data values in all the QRRs and this step of holistic alignment performs the alignment worldwide with all QRRs to create a table in which all data values of the same attribute are aligned in the same table column. Thus, holistic alignment problem is equal to that of finding connected components in an undirected graph. Each connected component of the graph represents a table column inside which the linked data values from different records are aligned vertically. We need to consider two application constraints that are specific to our holistic alignment problem.

- a) Vertices from the same record are not allowed to be included in the same connected component as they are considered to come from two different attributes of the record. If two vertices from the same record breach this constraint, a path must exist between the two, which we call a breach path.
- b) Connected components are not allowed to intersect each other. If C_1 and C_2 are two connected components, then vertices in C_1 should be either all on the left side of C_2 or all on the right side of C_2 , and vice versa.

So, we design a 3-steps for the holistic alignment problem. First, we traverse the graph once by a depth-first search to discover the preliminary connected components in the graph (the Traverse and Visit functions). Throughout the traversal, a color array is use to indicate whether every vertex has been

visited or not (WHITE for unvisited, GRAY for under processing, and BLACK for processing finished). In the Visit function, when a new vertex is encounter, we add it into the current connected component. Second, At the same time we also mark those components containing breach paths. If a connected component is start containing breach paths, the BreakBreachPath function is called to break it by remove the edges with the smallest sum of pairwise similarity scores. Given two recognized nodes v_i and v_j from the same record, the problem of breaking a breach path is accurately the max-flow/min-cut problem in which v_i and v_j are the source and sink nodes, respectively. Third, we traverse the components contain breach paths to eliminate some edges so as to break the breach paths (i.e., enforce the first constraint). Finally, we use a divide-and-conquer method to recognize and split up the intersecting components to inflict the second constraint.

C. Nested Structure Processing

This identifies the nested structures that exist in the QRRs. If a QRR contains a nested structure such that an attribute has multiple values, then various of the values might not be aligned to any other values. Hence nested structure giving out identifies the data values of a QRR that are generated by nested structures. Relying only on HTML tags to recognize nested structures, as is done by approximately all existing methods, may incorrectly recognize a plain structure as a nested one. To overcome this problem, CTVS uses both the HTML tags and the data values to identify the nested structures.

Given an aligned table, a nested column comprises at least two ordered sets signifying the data values that aregenerated by repetitive parts in the template. A nested column set C is comprised of a set of nested columns.

Given columns cp in a holistic alignment and a similarity threshold S_{nest} as input, the procedure nested decides, using the similarities of the data values in cp , whether cp contains a repetitive tag pattern that is formed by a nested structure. We assume that two columns are generated by the same attribute if there is a large data value similarity between these two columns. Given a column c_1 , which contains m data values, we define the intracolumn similarity sim_{intra} to be the average data value similarity within each column in c_1 .

$$sim_{intra} = \frac{\sum_{j=1}^{m-1} \sum_{i=j+1}^m s_{ij}}{m(m-1)} \quad (2)$$

In (2), s_{ij} is the data value similarity among the i th and j th data values of c_1 . For cp , its intracolumn similarity is the average of the intracolumn similarity of all columns in cp .

For two columns c_1 and c_2 , which have m and n data values, respectively, the intercolumn similarity sim_{inter} is defined to be the average data value similarity of every pair of data values in c_1 and c_2 .

$$\text{siminter} = \frac{\sum_{j=1}^n \sum_{i=1}^m s_{ij}}{mn} \quad (3)$$

In (3), s_{ij} is the data value similarity among the i th data values of c_1 and j th data values of c_2 using the data values similarity calculation described in Section i.

After siminter and simintra are calculated for identified columns cp , if $\text{siminter} = \text{simintra} > S_{\text{nest}}$, where S_{nest} is a threshold that is set to 0.5, cp is assumed to be a nested column set, which means that the data values in it are generated from a nested structure.

Given data columns cp and the nested column set C as input, the method `add_nested_column` adds the nested columns cp to C . Then c_i in C is replaced with cp . Otherwise, cp is simply added as a new element into C .

Given n records with a maximum of m data values and a maximum tag string length of l , the time complexity of the nested structure processing algorithm is $O(nl^2m^2)$. For each record, at most $O(l^2)$ time is needed to identify repetitive tag patterns; if a pattern is found, at most $O(m^2)$ time is required to calculate the intra-/intercolumn similarity. compare the nested structure processing methods in DeLa [29] and NET [20], the nested structure processing technique in CTVS has the following advantages.

- a) CTVS processes the nested structures after the data records are aligned rather than before as is the case in DeLa and NET. The nested structure before the records are aligned makes them weak to optional attributes that makes the tag structure irregular. This difficulty does not arise in CTVS.
- b) In CTVS the data value similarity information efficiently prevents a flat structure from being identified as a nested structure. It shares similar tag structures, a flat structure by several columns having the same tag structure, may be wrongly identified as a nested structure one can have serious consequences in DeLa and NET. DeLa condenses all the values into one parent value and then aligns them to other records, which makes the alignment much more complex. If NET wrongly identifies a simple structure as a nested structure, it will create a new row in the table for each data value of the simple structure.

IV. EXPERIMENTAL RESULTS

The act of the data extraction methods is compare in three different ways. General data set evaluation present the act on the first three data sets, which display a variety of properties and have been used in earlier work by others. The other two evaluations focus on exact properties of the query result pages. Noncontiguous QRR evaluation compares the act for query result pages in which the QRRs are contiguous and non-contiguous is shown in table 2. Nested-structure evaluation

compare the performance for query result pages with and without a nested structure is shown in table 3.

TABLE 2
Data Extraction Performance by AUXI Data Set

	Contiguous Pages			Non-contiguous Pages		
# QRRs	510			543		
Method	CTVS	ViNTs	DeLa	CTVS	ViNTs	DeLa
# extracted QRRs	506	502	503	530	432	437
# correct extracted QRRs	499	486	479	510	418	415
record-level precision	98.6%	96.8%	95.2%	96.2%	96.8%	95.4%
record-level recall	97.8%	95.2%	93.9%	93.9%	77%	76.5%
page-level precision	92.5%	90%	85%	85%	35%	37.5%

TABLE 3
Data Extraction Achieved by NESTED Data Set

	Non-nested Pages			Nested Pages		
# QRRs	420			449		
Method	CTVS	ViNTs	DeLa	CTVS	ViNTs	DeLa
# extracted QRRs	421	420	414	446	452	445
# correct extracted QRRs	411	400	398	428	388	417
record-level precision	97.6%	95.2%	96.1%	95.9%	86.4%	93.7%
record-level recall	97.9%	95.2%	94.7%	95.3%	85.8%	92.9%
page-level precision	95%	90%	90%	90%	75%	85%

V. RELATED WORK

In wrapper induction, extraction systems are copied based on inductive learning. This not scalable to a large number of web databases. Hence, the wrapper induction approach involves two additional difficult problems: monitoring changes in a page's format and maintaining a wrapper when a page's format changes. To conquer the problems of wrapper induction, some unsupervised learning methods, such as, IEPAD, ExAlg, DeLa, have been planned to repeatedly extract the data from the query result pages. To conquer these shortcomings, methods such as ViPER and ViNTs make use of extra information in the query result pages.

All the works make use of only the information in the query result pages to execute the data extraction. There are works that make use of extra information, specifically ontologies, to aid in the data extraction. While these approaches can overcome some of the limitations of CTVS (e.g., that a query result page should contain at least two

QRRs) and can get high accuracy, they need the availability of extra resources to construct an ontology as well as the additional step of actually constructing the ontology.

VI. CONCLUSIONS AND FUTURE DIRECTIONS

We presented a new data extraction method, CTVS, to repeatedly extract QRRs from a query result page. CTVS employs two steps i.e. The first step identify and segment the QRRs. This improve on alive technique by allowing the QRRs in a data region to be non-contiguous. The second step aligns the data values between the QRRs. A novel alignment method is proposed in which the alignment is perform in three successive steps: pairwise alignment, holistic alignment, and nested structure processing.

Although CTVS has been shown to be an correct data extraction method, it still suffer from some restrictions. First, it require at least two QRRs in the query result page. Second, any optional attribute that appear as the start node

in a data region will be treat as auxiliary information. Third, CTVS mostly depends on tag structures to find out data values. Therefore, CTVS does not hold the case where multiple data values from more than one attribute are clustered inside one leaf node of the tag tree, as well as the case where one data value of a single element spans multiple leaf nodes.

REFERENCES

- [1] B. Liu and Y. Zhai, "NET - A System for Extracting Web Data from Flat and Nested Data Records," Proc. Sixth Int'l Conf. Web Information Systems Eng., pp. 487-495, 2005.
- [2] J. Wang and F.H. Lochovsky, "Data Extraction and Label Assignment for Web Databases," Proc. 12th World Wide Web Conf., pp. 187-196, 2003.
- [3] K. Simon and G. Lausen, "ViPER: Augmenting Automatic Information Extraction with Visual Perceptions," Proc. 14th ACM Int'l Conf. Information and Knowledge Management, pp. 381-388, 2005.
- [4] Y. Zhai and B. Liu, "Structured Data Extraction from the Web Based on Partial Tree Alignment," IEEE Trans. Knowledge and Data Eng., vol. 18, no. 12, pp. 1614-1628, Dec. 2006.
- [5] H. Zhao, W. Meng, Z. Wu, V. Raghavan, and C. Yu, "Fully Automatic Wrapper Generation for Search Engines," Proc. 14th World Wide Web Conf., pp. 66-75, 2005.
- [6] M.K. Bergman, "The Deep Web: Surfacing Hidden Value," White Paper, BrightPlanet Corporation, <http://www.brightplanet.com/resources/details/deepweb.html>, 2001.
- [7] K.C.-C. Chang, B. He, C. Li, M. Patel, and Z. Zhang, "Structured Databases on the Web: Observations and

Implications," SIGMOD Record, vol. 33, no. 3, pp. 61-70, 2004.

[8] C.H. Chang and S.C. Lui, "IEPAD: Information Extraction Based on Pattern Discovery," Proc. 10th World Wide Web Conf., pp. 681- 688, 2001.

[9] A. Arasu and H. Garcia-Molina, "Extracting Structured Data from Web Pages," Proc. ACM SIGMOD Int'l Conf. Management of Data, pp. 337-348, 2003.

[10] H. Snoussi, L. Magnin, and J.-Y. Nie, "Heterogeneous Web Data Extraction Using Ontologies," Proc. Fifth Int'l Conf. Agent-Oriented Information Systems, pp. 99-110, 2001.

[11] W. Su, J. Wang, and F.H. Lochovsky, "ODE: Ontology-Assisted Data Extraction," ACM Trans. Database Systems, vol. 34, no. 2, article 12, p. 35, 2009.

[12] D.W. Embley, D.M. Campbell, Y.S. Jiang, S.W. Liddle, D.W. Lonsdale, Y.-K. Ng, and R.D. Smith, "Conceptual-Model-Based Data Extraction from Multiple-Record Web Pages," Data and Knowledge Eng., vol. 31, no. 3, pp. 227-251, 1999.